

# Probabilistic Graphical Models

10-708

## Towards Complex Graphical Models and Approximate Inference

Eric Xing

Lecture 16, Nov 7, 2005

Reading: MJ-Chap. 21

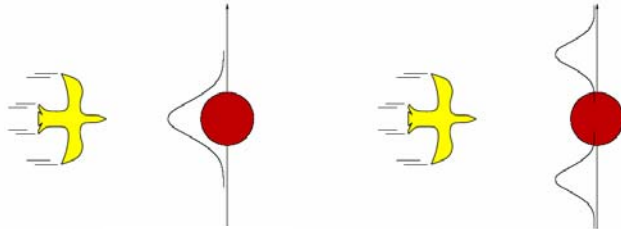


## The need for multimodal belief states in dynamic models

- An LDS defines only unimodal belief states

$$\hat{\mathbf{x}}_{t+1|t+1} = \hat{\mathbf{x}}_{t+1|t} + K_{t+1}(y_{t+1} - C\hat{\mathbf{x}}_{t+1|t})$$

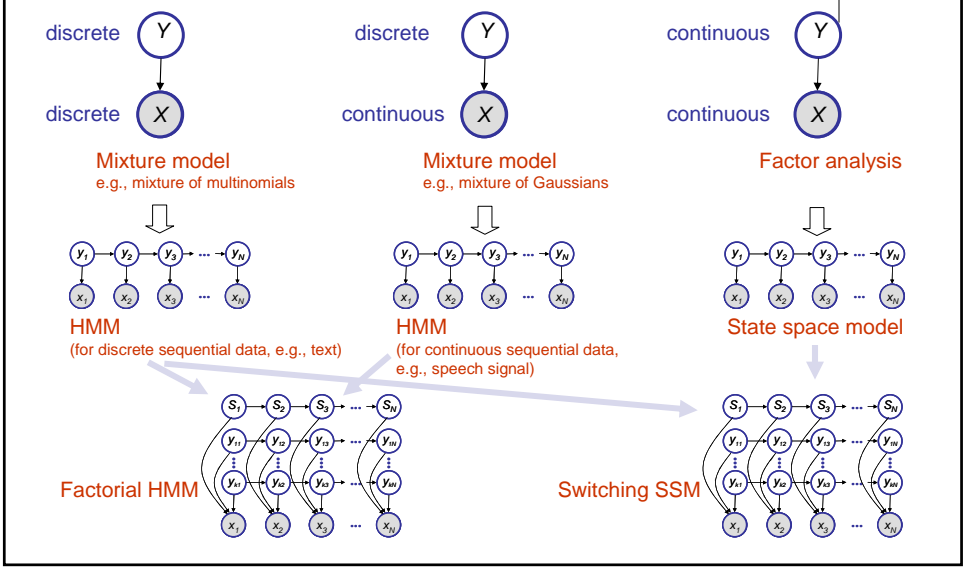
$$P_{t+1|t+1} = P_{t+1|t} - K C P_{t+1|t}$$



- (a) A Kalman filter will predict the location of the bird using a single Gaussian centered on the obstacle.
- (b) A more realistic model allows for the bird's evasive action, predicting that it will fly to one side or the other.



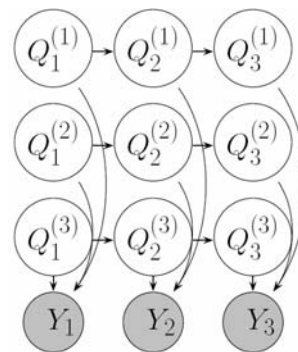
# A road map to more complex dynamic models



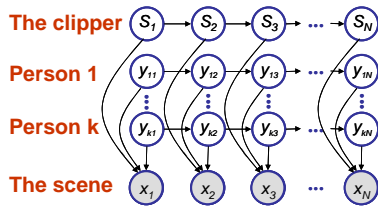
# Factorial HMM



- The belief state at each time is  $X_t = \{Q_t^{(1)}, \dots, Q_t^{(k)}\}$  and in the most general case has a state space  $O(d^k)$  for  $k$   $d$ -nary chains
- The common observed child  $Y_t$  couples all the parents (explaining away).
- But the parameterization cost for fHMM is  $O(kd^2)$  for  $k$  chain-specific transition models  $p(Q_t^{(i)} | Q_{t-1}^{(i)})$  rather than  $O(d^k)$  for  $p(X_t | X_{t-1})$



## Special case: switching HMM

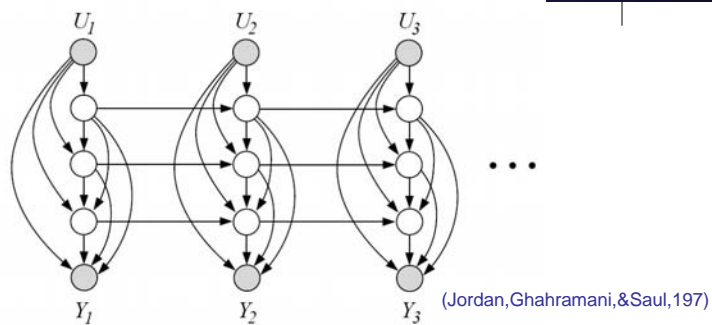


Multi-View Face Tracking with Factorial and Switching HMM

- Different chains have different state space and different semantics
- The exact calculation is intractable and we must use approximate inference methods

Peng Wang, Qiang Ji  
 Department of Electrical, Computer and System Engineering  
 Rensselaer Polytechnic Institute  
 Troy, NY 12180

## Hidden Markov decision trees



(Jordan, Ghahramani, & Saul, 1997)

- A combination of decision trees with factorial HMMs
- This gives a "command structure" to the factorial representation
- Appropriate for multi-resolution time series
- Again, the exact calculation is intractable and we must use approximate inference methods

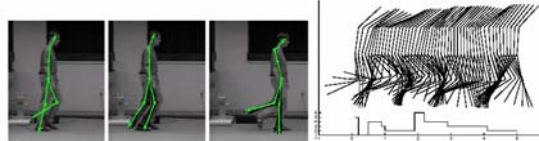
# Switching LDS



- Possible world:
  - multiple motion state:



- Task:
  - Trajectory prediction



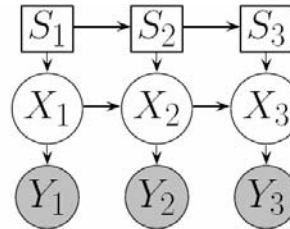
- Model:
  - Combination of HMM and LDS

$$p(X_t = x_t | X_{t-1} = x_{t-1}, S_t = i) = \mathcal{N}(x_t; A_i x_{t-1}, Q_i)$$

$$p(Y_t = y_t | X_t = x_t) = \mathcal{N}(y_t; C x_t, R)$$

$$p(S_t = j | S_{t-1} = i) = M(i, j)$$

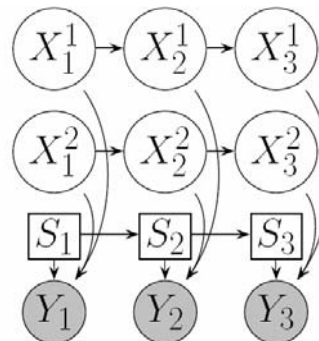
- Belief state has  $O(k^t)$  Gaussian modes:



# Data association (correspondence problem)



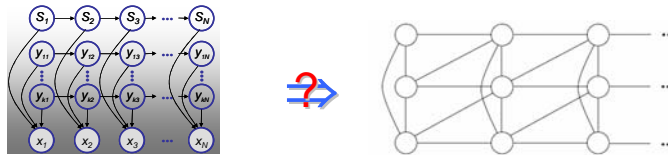
- Optimal belief state has  $O(k^t)$  modes.
- Common to use nearest neighbor approximation.
- For each time slice, can enforce that at most one source causes each observation
- Correspondence problem also arises in shape matching and stereo vision.



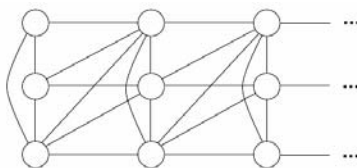
# Triangulating fHMM



- Is the following triangulation correct?



- Here is a triangulation

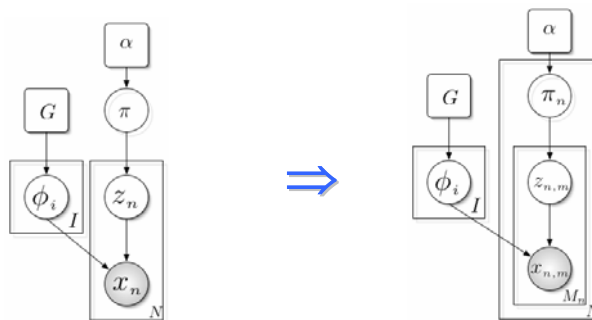


- We have created cliques of size  $k+1$ , and there are  $O(kT)$  of them. The junction tree algorithm is not efficient for factorial HMMs.

# Mixed Membership Model ( $M^3$ )



- Mixture versus admixture



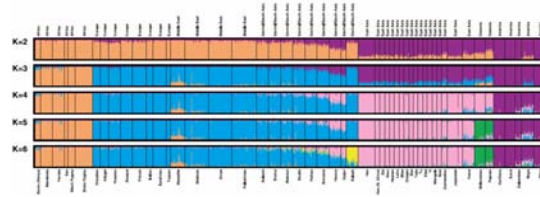
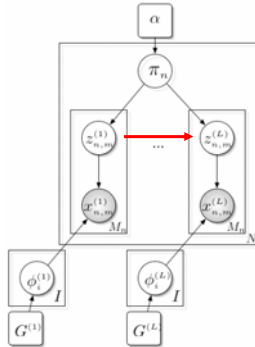
A Bayesian mixture model

A Bayesian admixture model:  
Mixed membership model

# Population admixture: $M^3$ in genetics



- The genetic materials of each modern individual are inherited from multiple ancestral populations, each DNA locus may have a different generic origin ...



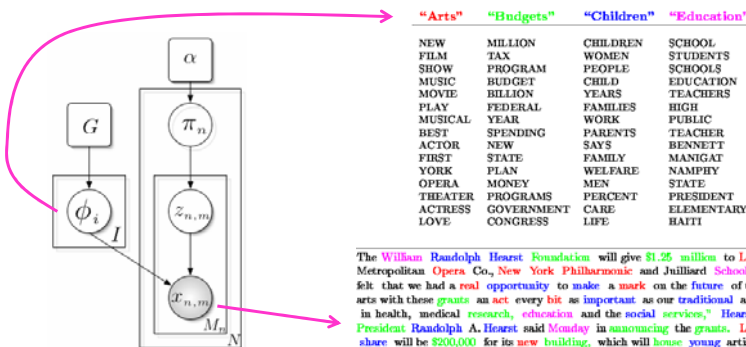
**Genetic Structure of Human Populations**  
 Noah A. Rosenberg,<sup>1\*</sup> Jonathan K. Pritchard,<sup>2</sup> James L. Weber,<sup>3</sup>  
 Howard M. Cann,<sup>4</sup> Kenneth K. Kidd,<sup>5</sup> Lev A. Zhivotovskiy,<sup>6</sup>  
 Marcus W. Feldman<sup>7</sup>  
 SCIENCE VOL 298 20 DECEMBER 2002

- Ancestral labels may have (e.g., Markovian) dependencies

# Latent Dirichlet Allocation: $M^3$ in text mining



- A document is a bag of words each generated from a randomly selected topic



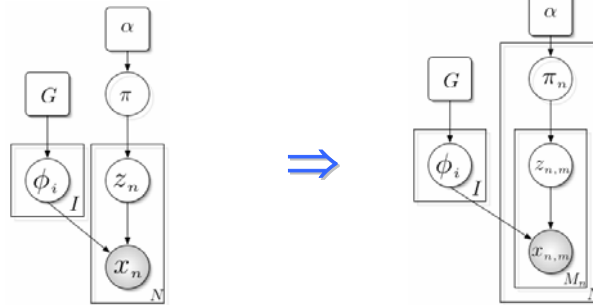
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Raudolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services." Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

# Inference in Mixed Membership Models



- Mixture versus admixture



$$p(D) = \sum_{\{z_{n,m}\}} \int \dots \int \left( \prod_n \left( \prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | \mathcal{G}) d\pi_1 \dots d\pi_N d\phi$$

- Inference is very hard in  $M^3$ , all hidden variables are coupled and not factorizable!

$$p(\pi_n | D) \sim \sum_{\{z_{n,m}\}} \int \left( \prod_n \left( \prod_m p(x_{n,m} | \phi_{z_n}) p(z_{n,m} | \pi_n) \right) p(\pi_n | \alpha) \right) p(\phi | \mathcal{G}) d\pi_{-n} d\phi$$

# Approaches to inference

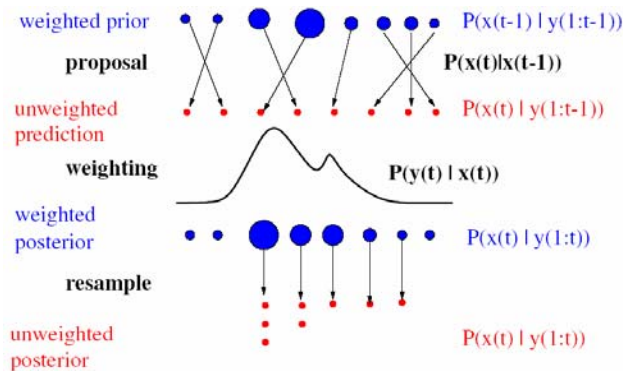


- Exact inference algorithms
  - The elimination algorithm
  - The junction tree algorithms
- Approximate inference techniques
  - Monte Carlo algorithms:
    - Stochastic simulation / sampling methods
    - Markov chain Monte Carlo methods
  - Variational algorithms:
    - Belief propagation
    - Assumed density filtering
    - Variational inference

## Example: Particle filtering (sequential Monte Carlo)



- Represent belief state as weighted set of samples (non-parametric).
- Can handle nonlinear transition/emission and multi-modality.
- Easy to implement.
- Only works well in small dimensions.



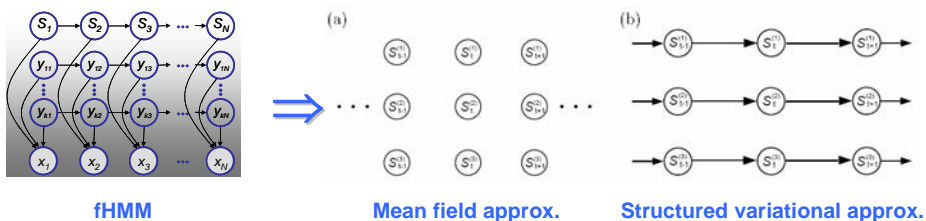
## Example: Structured Variational approximation



- Finds an optimal  $q^*(\cdot)$  in a **tractable family** to approximate the original joint  $p(\cdot)$

$$q^*(\cdot) \in \arg \min_{q \in \mathcal{F}} F(q \| p)$$

- There can be many different choices of  $\mathcal{F}$  and  $F(\cdot)$ .



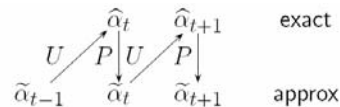


## Example: Assumed density filtering (ADF)



- ADF forces the **belief state** to live in some restricted family  $\mathcal{F}$ , e.g., product of histograms, Gaussian.
- Given a prior  $\tilde{\alpha}_{t-1} \in \mathcal{F}$ , do one step of exact Bayesian updating to get  $\hat{\alpha}_t \notin \mathcal{F}$ . Then do a projection step to find the closest approximation in the family:

$$\tilde{\alpha}_t \in \arg \min_{q \in \mathcal{F}} \text{KL}(\hat{\alpha}_t \parallel q)$$



- The Boyen-Koller (BK) algorithm is ADF applied to a DBN
  - e.g., let  $\mathcal{F}$  be a product of (singleton) marginals:
- This is also a variational method, and the updating step can still be intractable

## Monte Carlo methods



- Draw random samples from the desired distribution
- Yield a stochastic representation of a complex distribution
  - marginals and other expectations can be approximated using **sample-based averages**

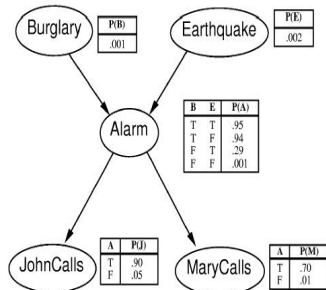
$$E[f(\mathbf{x})] = \frac{1}{N} \sum_{t=1}^N f(\mathbf{x}^{(t)})$$

- **Asymptotically** exact and easy to apply to arbitrary models
- Challenges:
  - how to draw samples from a given dist. (not all distributions can be trivially sampled)?
  - how to make better use of the samples (not all sample are useful, or equally useful, see an example later)?
  - how to know we've sampled enough?

## Example: naive sampling



- Construct samples according to probabilities given in a BN.



E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E1	B0	A1	M1	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0

**Alarm example:** (Choose the right sampling sequence)

1) Sampling:  $P(B) = \langle 0.001, 0.999 \rangle$  suppose it is false, B0. Same for E0.  $P(A|B0, E0) = \langle 0.001, 0.999 \rangle$  suppose it is false...

2) Frequency counting: In the samples right,  $P(J|A0) = P(J, A0) / P(A0) = \langle 1/9, 8/9 \rangle$ .

## Example: naive sampling



- Construct samples according to probabilities given in a BN.

**Alarm example:** (Choose the right sampling sequence)

3) what if we want to compute  $P(J|A1)$ ?  
we have only one sample ...  
 $P(J|A1) = P(J, A1) / P(A1) = \langle 0, 1 \rangle$ .

4) what if we want to compute  $P(J|B1)$ ?  
**No such sample available!**  
 $P(J|A1) = P(J, B1) / P(B1)$  can not be defined.

For a model with hundreds or more variables, rare events will be very hard to garner enough samples even after a long time or sampling ...

E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E1	B0	A1	M1	J1
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0
E0	B0	A0	M0	J0

## Monte Carlo methods (cond.)



- Direct Sampling
  - We have seen it.
  - Very difficult to populate a high-dimensional state space
- Rejection Sampling
  - Create samples like direct sampling, only count samples which is consistent with given evidences.
- Likelihood weighting, ...
  - Sample variables and calculate evidence weight. Only create the samples which support the evidences.
- Markov chain Monte Carlo (MCMC)
  - Metropolis-Hasting
  - Gibbs

## Rejection sampling



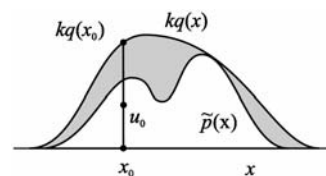
- Suppose we wish to sample from dist.  $\Pi(X)=\Pi'(X)/Z$ .
  - $\Pi(X)$  is difficult to sample, but  $\Pi'(X)$  is easy to evaluate
  - Sample from a simpler dist  $Q(X)$
  - Rejection sampling

$$x^* \sim Q(X), \quad \text{accept } x^* \text{ w.p. } \Pi'(x^*)/kQ(x^*)$$

- Correctness:

$$\begin{aligned}
 p(x) &= \frac{[\Pi'(x)/kQ(x)]Q(x)}{\int [\Pi'(x)/kQ(x)]Q(x)dx} \\
 &= \frac{\Pi'(x)}{\int \Pi'(x)dx} = \Pi(x)
 \end{aligned}$$

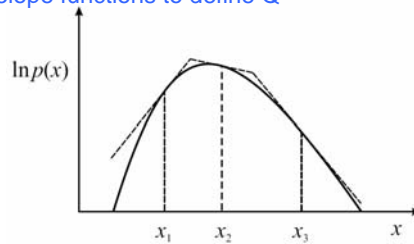
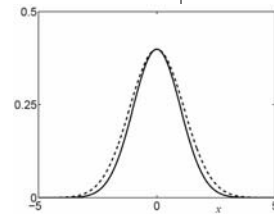
- Pitfall ...



## Rejection sampling



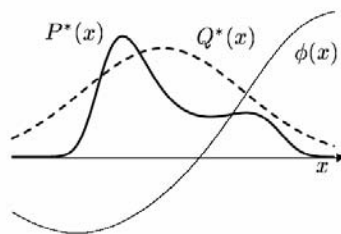
- Pitfall:
  - Using  $Q = \mathcal{N}(\mu, \sigma_q)$  to sample  $P = \mathcal{N}(\mu, \sigma_p)$
  - If  $\sigma_q$  exceeds  $\sigma_p$  by 1%, and  $\text{dimensional} = 1000$ ,
  - The optimal acceptance rate  $k = (\sigma_q/\sigma_p)^d \approx 1/20,000$
  - Big waste of samples!
- Adaptive rejection sampling
  - Using envelope functions to define  $Q$



## Unnormalized importance sampling



- Suppose sampling from  $P(\cdot)$  is hard.
- Suppose we can sample from a "simpler" proposal distribution  $Q(\cdot)$  instead.
- If  $Q$  dominates  $P$  (i.e.,  $Q(x) > 0$  whenever  $P(x) > 0$ ), we can sample from  $Q$  and reweight:



$$\begin{aligned}
 \langle f(X) \rangle &= \int f(x) P(x) dx \\
 &= \int f(x) \frac{P(x)}{Q(x)} Q(x) dx \\
 &\approx \frac{1}{M} \sum_m f(x^m) \frac{P(x^m)}{Q(x^m)} \quad \text{where } x^m \sim Q(X) \\
 &= \frac{1}{M} \sum_m f(x^m) w^m
 \end{aligned}$$



## Normalized importance sampling

- Suppose we can only evaluate  $P'(x) = \alpha P(x)$  (e.g. for an MRF).
- We can get around the nasty normalization constant  $\alpha$  as follows:

- Let  $r(X) = \frac{P'(x)}{Q(x)} \Rightarrow \langle r(X) \rangle_Q = \int \frac{P'(x)}{Q(x)} Q(x) dx = \int P'(x) dx = \alpha$

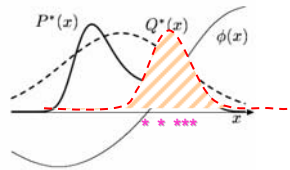
- Now

$$\begin{aligned} \langle f(X) \rangle_P &= \int f(x) P(x) dx = \frac{1}{\alpha} \int f(x) \frac{P'(x)}{Q(x)} Q(x) dx \\ &= \frac{\int f(x) r(x) Q(x) dx}{\int r(x) Q(x) dx} \\ &\approx \frac{\sum_m f(x^m) r^m}{\sum_m r^m} \quad \text{where } x^m \sim Q(X) \\ &= \sum_m f(x^m) w^m \quad \text{where } w^m = \frac{r^m}{\sum_m r^m} \end{aligned}$$



## Weighted resampling

- Problem of importance sampling: depends on how well  $Q$  matches  $P$ 
  - If  $P(x)f(x)$  is strongly varying and has a significant proportion of its mass concentrated in a small region,  $r_m$  will be dominated by a few samples



- Note that if the high-prob mass region of  $Q$  falls into the low-prob mass region of  $P$ , the variance of  $r^m = P(x^m)/Q(x^m)$  can be small even if the samples come from low-prob region of  $P$  and potentially erroneous .

- Solution

- Use heavy tail  $Q$ .

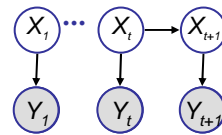
- Weighted resampling

$$w^m = \frac{P(x^m)/Q(x^m)}{\sum_l P(x^l)/Q(x^l)} = \frac{r^m}{\sum_m r^m}$$

# Weighted resampling



- Sampling importance resampling (SIR):
  1. Draw  $N$  samples from  $Q$ :  $X_1 \dots X_N$
  2. Constructing weights:  $w_1 \dots w_N$ ,  $w^m = \frac{p(x^m)/Q(x^m)}{\sum_j p(x^j)/Q(x^j)} = \frac{r^m}{\sum_m r^m}$
  3. Sub-sample  $x$  from  $\{X_1 \dots X_N\}$  w.p.  $(w_1 \dots w_N)$
- Particular Filtering



- A special weighted resampler
- Yield samples from posterior  $p(X_t | Y_{1:t})$

# Sketch of Particle Filters



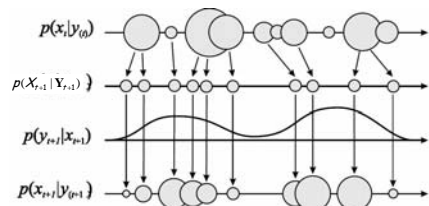
- The starting point
 
$$p(X_t | Y_{1:t}) = p(X_t | Y_t, Y_{1:t-1}) = \frac{p(X_t | Y_{1:t-1})p(Y_t | X_t)}{\int p(X_t | Y_{1:t-1})p(Y_t | X_t)dX_t}$$
  - Thus  $p(X_t | Y_{1:t})$  is represented by
 
$$\left\{ X_t^m \sim p(X_t | Y_{1:t-1}), w_t^m = \frac{p(Y_t | X_t^m)}{\sum_{m=1}^M p(Y_t | X_t^m)} \right\}$$

- A sequential weighted resampler

- Time update
 
$$p(X_{t+1} | Y_{1:t}) = \int p(X_{t+1} | X_t)p(X_t | Y_{1:t})dX_t$$

$$= \sum_m w_t^m p(X_{t+1} | X_t) \text{ (sample from a mixture model)}$$
- Measurement update
 
$$p(X_{t+1} | Y_{1:t+1}) = \frac{p(X_{t+1} | Y_{1:t})p(Y_{t+1} | X_{t+1})}{\int p(X_{t+1} | Y_{1:t})p(Y_{t+1} | X_{t+1})dX_{t+1}}$$

$$\Rightarrow \left\{ X_{t+1}^m \sim p(X_{t+1} | Y_{1:t}), w_{t+1}^m = \frac{p(Y_{t+1} | X_{t+1}^m)}{\sum_{m=1}^M p(Y_{t+1} | X_{t+1}^m)} \right\} \text{ (reweight)}$$





## Rao-Blackwellised sampling

- Sampling in high dimensional spaces causes high variance in the estimate.
- RB idea: sample some variables  $X_p$ , and conditional on that, compute expected value of rest  $X_d$  analytically:

$$\begin{aligned} E_{p(x|e)}[f(X)] &= \int p(x_p, x_d | e) f(x_p, x_d) dx_p dx_d \\ &= \int p(x_p | e) \left( \int p(x_d | x_p, e) f(x_p, x_d) dx_d \right) dx_p \\ &= \int p(x_p | e) E_{p(x_d|x_p,e)}[f(x_p, X_d)] dx_p \\ &= \frac{1}{M} \sum_m E_{p(x_d|x_p^m,e)}[f(x_p^m, X_d)] \quad x_p^m \sim p(x_p | e) \end{aligned}$$

- This has lower variance, because of the identity:

$$\text{var}[\tau(X_p, X_d)] = \text{var}[E[\tau(X_p, X_d) | X_p]] + E[\text{var}[\tau(X_p, X_d) | X_p]]$$

- Hence  $\text{var}[E[\tau(X_p, X_d) | X_p]] \leq \text{var}[\tau(X_p, X_d)]$ , so  $\tau(X_p, X_d) = E[f(X_p, X_d) | X_p]$  is a lower variance estimator.



## Markov chain Monte Carlo (MCMC)

- Importance sampling does not scale well to high dimensions.
- Rao-Blackwellisation not always possible.
- MCMC is an alternative.
- Construct a Markov chain whose stationary distribution is the target density  $= \mathcal{P}(X|e)$ .
- Run for  $T$  samples (burn-in time) until the chain converges/mixes/reaches stationary distribution.
- Then collect  $M$  (correlated) samples  $x_m$ .
- Key issues:
  - Designing proposals so that the chain mixes rapidly.
  - Diagnosing convergence.