









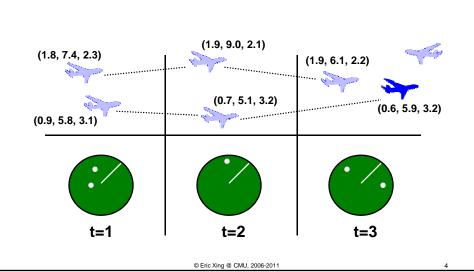
- How to segment images?
  - Manual segmentation (very expensive)
  - Algorithm segmentation
    - K-means
    - Statistical mixture models
    - Spectral clustering
- Problems with most existing algorithms
  - Ignore the spatial information
  - Perform the segmentation one image at a time
  - Need to specify the number of segments a priori

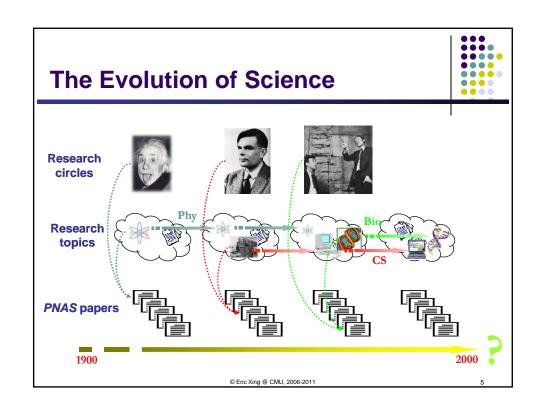
© Eric Xing @ CMU, 2006-2011

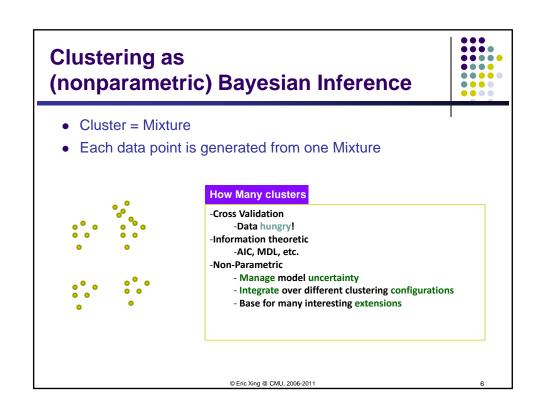
3

## **Object Recognition and Tracking**









### **Outline**



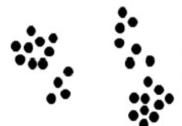
- Dirichlet Process: From finite mixture to Infinite mixture model
  - The Chinese Restaurant Process
  - Example: heliotype inference
- Intro to Markov Chain Monte Carlo
  - Gibbs sampling
- Dynamic Dirichlet Process
  - The recurrent Chinese Restaurant Process
  - Example: Application: evolutionary clustering of documents

© Eric Xing @ CMU, 2006-2011

7

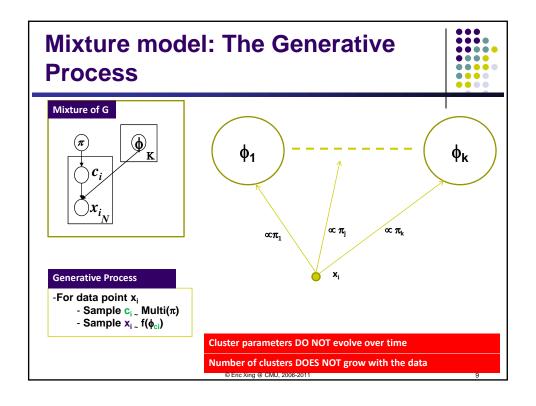
# **Clustering**





- How to label them?
- How many clusters ???

© Eric Xing @ CMU, 2006-2011



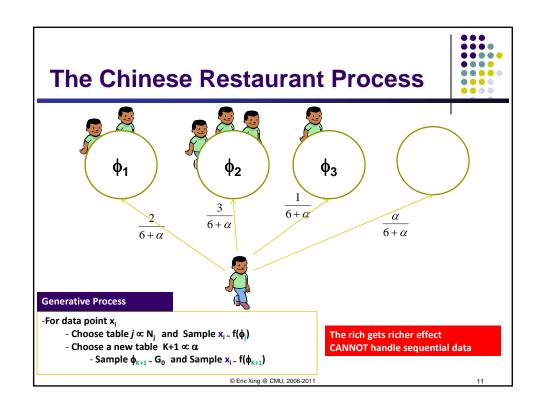
#### **The Chinese Restaurant Process**

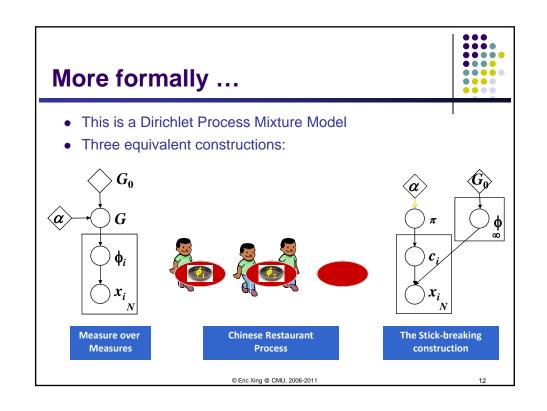


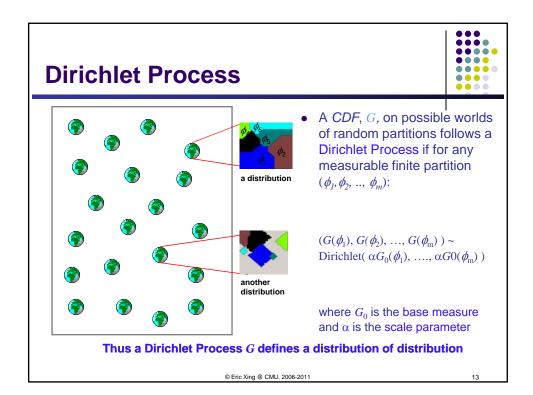
- Customers correspond to data points
- Tables correspond to clusters/mixture components
- Dishes correspond to parameter of the mixtures

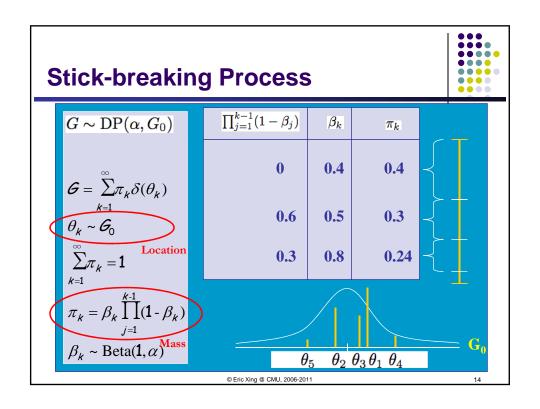


© Eric Xing @ CMU, 2006-2011



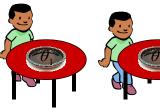












$$P(c_{i} = k \mid \mathbf{c}_{-i}) = 1$$

$$\frac{1}{1+\alpha}$$

$$\frac{1}{2+\alpha}$$

$$\frac{1}{3+\alpha}$$

$$m$$

$$\begin{array}{ccc}
0 & 0 \\
\frac{\alpha}{1+\alpha} & 0 \\
\frac{1}{2+\alpha} & \frac{\alpha}{2+\alpha} \\
\frac{2}{3+\alpha} & \frac{\alpha}{3+\alpha} \\
\frac{m_2}{i+\alpha-1} & \cdots & \frac{\alpha}{i+\alpha-1}
\end{array}$$

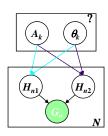
CRP defines an exchangeable distribution on partitions over an (infinite) sequence of samples, such a distribution is formally known as the Dirichlet Process (DR)

© Eric Xing @ CMU, 2006-2011

## **Case study: ancestral Inference**

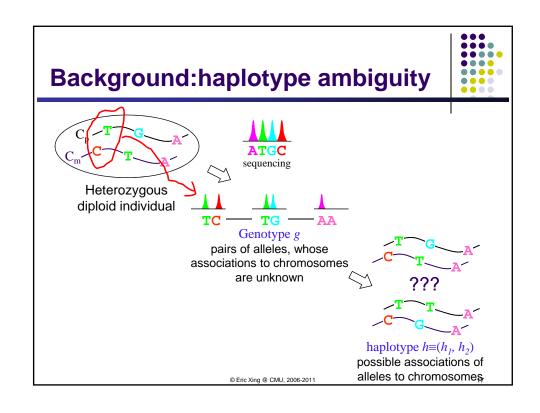


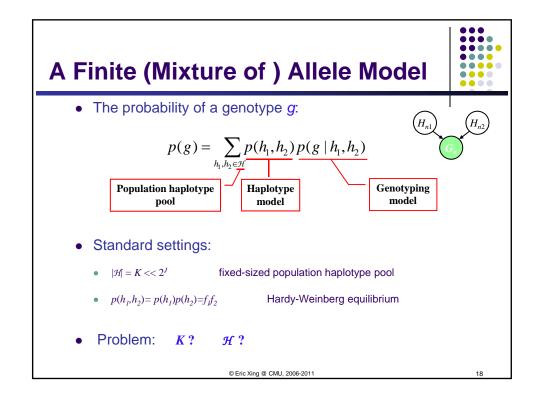




Essentially a clustering problem, but ...

- Better recovery of the ancestors leads to better haplotyping results (because of more accurate grouping of common haplotypes)
- True haplotypes are obtainable with high cost, but they can validate model more subjectively (as opposed to examining saliency of clustering)
- Many other biological/scientific utilities
   © Eric Xing @ CMU, 2006-2011

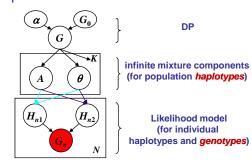




### **Example: DP-haplotyper** [Xing et al, 2004]



Clustering human populations



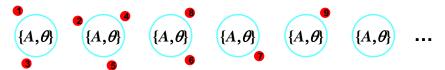
- Inference: Markov Chain Monte Carlo (MCMC)
  - Gibbs sampling
  - Metropolis Hasting

© Eric Xing @ CMU, 2006-2011

## The DP Mixture of Ancestral **Haplotypes**



- The customers around a table in CRP form a cluster
  - associate a mixture component (i.e., a population haplotype) with a table
  - sample  $\{a, \theta\}$  at each table from a base measure  $G_0$  to obtain the population haplotype and nucleotide substitution frequency for that component



With  $p(h/\{A, \theta\})$  and  $p(g/h_1, h_2)$ , the CRP yields a posterior distribution on the number of population haplotypes (and on the haplotype configurations and the nucleotide substitution frequencies)

© Eric Xing @ CMU, 2006-2011

## **MCMC** for Haplotype Inference



- · Gibbs sampling for exploring the posterior distribution under the proposed model
  - Integrate out the parameters such as  $\theta$  or  $\lambda$ , and sample  $c_{i_e}$ ,  $a_k$ and  $h_{i_a}$

$$p(c_{i_e} = k \mid \mathbf{c}_{[-i_e]}, \mathbf{h}, \mathbf{a}) \propto p(c_{i_e} = k \mid \mathbf{c}_{[-i_e]}) \ p(h_{i_e} \mid a_{k,} \mathbf{h}_{[-i_e]}, \mathbf{c})$$
Posterior
Prior x Likelihood

CRP

Gibbs sampling algorithm: draw samples of each random variable to be sampled given values of all the remaining variables

© Eric Xing @ CMU, 2006-2011

### **MCMC** for Haplotype Inference



Sample  $c_{ie}^{(j)}$ , from

$$\begin{split} & p(c_{i_e}^{(j)} = k | \mathbf{c}^{[-j,i_e]}, \mathbf{h}, \mathbf{a}) \\ & \propto p(c_{i_e}^{(j)} = k | \mathbf{c}^{[-j,i_e]}, \mathbf{m}, \mathbf{n}) p(h_{i_e}^{(j)} | a_k, \mathbf{c}, \mathbf{h}^{[-j,i_e]}) \\ & \propto (m_{jk}^{[-j,i_e]} + \tau \beta_k) p(h_{i_e}^{(j)} | a_k, \mathbf{l}_k^{[-j,i_e]}), \text{ for } k = 1, ..., K + 1 \end{split}$$

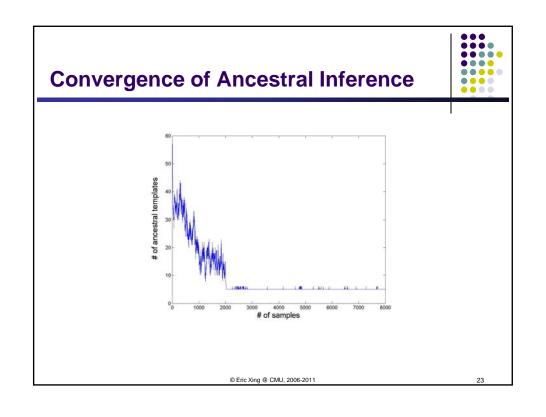
Sample  $a_k$  from

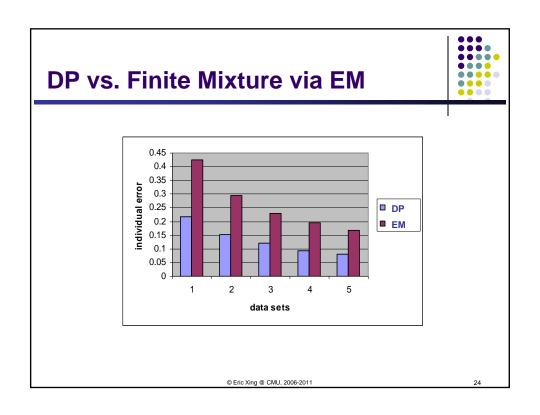
$$p(a_{k,t}|\mathbf{c}, \mathbf{h}) \propto \prod_{j,i_e|c_{i_e,t}^{(j)} = k} p(h_{i_e,t}^{(j)}|a_{k,t}, l_{k,t}^{(j)})$$

$$= \frac{\Gamma(\alpha_h + l_{k,t})\Gamma(\beta_h + l_{k,t}')}{\Gamma(\alpha_h + \beta_h + m_k)(|B| - 1)^{l_{k,t}'}} R(\alpha_h, \beta_h)$$

- $\text{Sample } h_{ie}^{(j)} \text{ from } \qquad p(h_{i_e,t}^{(j)}|\mathbf{h}_{[-i_e,t]}^{(j)},\mathbf{c},\mathbf{a},\mathbf{g})$
- For DP scale parameter  $\alpha$ : a vague inverse Gamma prior

© Eric Xing @ CMU, 2006-2011





### **Outline**



- Dirichlet Process: From finite mixture to Infinite mixture model
  - The Chinese Restaurant Process
  - Example: heliotype inference
- Intro to Markov Chain Monte Carlo
  - Gibbs sampling
- Dynamic Dirichlet Process
  - The recurrent Chinese Restaurant Process
  - Example: Application: evolutionary clustering of documents

© Eric Xing @ CMU, 2006-2011

25

### **Monte Carlo methods**



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

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

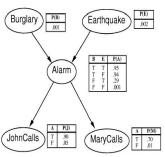
- Asymptotically exact and easy to apply to arbitrary models
- Challenges:
  - how to draw samples from a given dist. (not all distributions can be trivially sampled)?
  - how to make better use of the samples (not all sample are useful, or eqally useful, see an example later)?
  - how to know we've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

© Eric Xing @ CMU, 2006-2011

## **Example: naive sampling**



 Sampling: Construct samples according to probabilities given in a BN.



Alarm example: (Choose the right sampling sequence)

1) Sampling:P(B)=<0.001, 0.999> suppose it is false, B0. Same for E0. P(A|B0, E0)=<0.001, 0.999> suppose it is false...

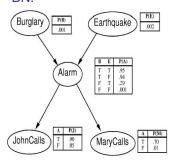
© Eric Xing @ CMU, 2006-2011

27

### **Example: naive sampling**



 Sampling: Construct samples according to probabilities given in a BN.



Alarm example: (Choose the right sampling sequence)

- 1) Sampling:P(B)=<0.001, 0.999> suppose it is false, B0. Same for E0. P(A|B0, E0)=<0.001, 0.999> suppose it is false...
- 2) Frequency counting: In the samples right,

P(J|A0)=P(J,A0)/P(A0)=<1/9, 8/9>.

	E0	B0	A0	MO	J0	
	E0	В0	A0	M0	J0	
	E0	В0	A0	MO	J1	
	E0	В0	A0	M0	J0	
	E0	B0	A0	M0	J0	
M)	E0	В0	A0	MO	J0	
ing	E1	В0	A1	M1	J1	
e it is false,	E0	В0	A0	M0	J0	
99>	E0	В0	A0	M0	J0	
ght,	E0	В0	A0	M0	J0	
© Eric Xing @ CMU, 2006-2011 28						

## **Example: naive sampling**



 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)=<0, 1>.

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 evough samples even after a long time or sampling ...

E0	B0	A0	MO	J0
E0	В0	A0	MO	J0
E0	В0	A0	MO	J1
E0	B0	A0	MO	J0
E0	B0	A0	MO	J0
E0	B0	A0	MO	J0
E1	B0	A1	M1	J1
E0	В0	A0	M0	J0
E0	B0	A0	MO	J0
E0	В0	A0	MO	J0

© Eric Xing @ CMU, 2006-2011

29

### **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.
- ....
- Markov chain Monte Carlo (MCMC)

© Eric Xing @ CMU, 2006-2011

#### **Markov chain Monte Carlo**



- Samples are obtained from a Markov chain (of sequentially evolving distributions) whose stationary distribution is the desired p(x)
- Construct a Markov chain whose stationary distribution is the target density = 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.

© Eric Xing @ CMU, 2006-2011

31

### **Gibbs sampling**



- Gibbs sampling is an MCMC algorithm that is especially appropriate for inference in graphical models.
- The procedue
  - we have variable set  $X=\{x_1, x_2, x_3,..., x_N\}$  for a GM
  - at each step one of the variables  $X_i$  is selected (at random or according to some fixed sequences), denote the remaining variables as  $X_i$ , and its current value as  $X_i^{(f-1)}$ 
    - Using the "alarm network" as an example, say at time t we choose  $\mathcal{X}_{E}$  and we denote the current value assignments of the remaining variables,  $\mathcal{X}_{\mathcal{E}}$ , obtained from previous samples, as  $\mathcal{X}_{\mathcal{E}}^{(t-1)} = \left\{ x_{\mathcal{B}}^{(t-1)}, x_{\mathcal{A}}^{(t-1)}, x_{\mathcal{A}}^{(t-1)}, x_{\mathcal{A}}^{(t-1)} \right\}$
  - the conditional distribution  $p(X_i | \mathbf{x}_i^{(t-1)})$  is computed
  - a value  $x_i^{(t)}$  is sampled from this distribution
  - the sample  $\mathbf{x}_{i}^{(f)}$  replaces the previous sampled value of  $X_{i}$  in  $X_{i}$ .
    - i.e.,  $\mathbf{X}^{(t)} = \mathbf{X}_{-F}^{(t-1)} \cup \mathbf{X}_{F}^{(t)}$

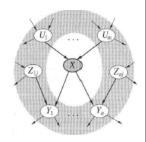
© Eric Xing @ CMU, 2006-2011

## Why Gibbs Sampling is Simple



- Markov-Blanket
  - A variable is independent from others, given its parents, children and children's parents. d-separation.

 $\Rightarrow p(X_i \mid X_j) = p(X_i \mid MB(X_j))$ 



- Gibbs sampling
  - Create a random sample.
     Every step, choose one
     variable and sample it by
     P(X|MB(X)) based on previous sample.

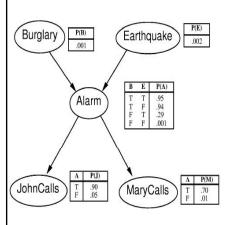
 $MB(A)=\{B, E, J, M\}$  $MB(E)=\{A, B\}$ 

© Eric Xing @ CMU, 2006-2011

33

### **Example: alarm network again**

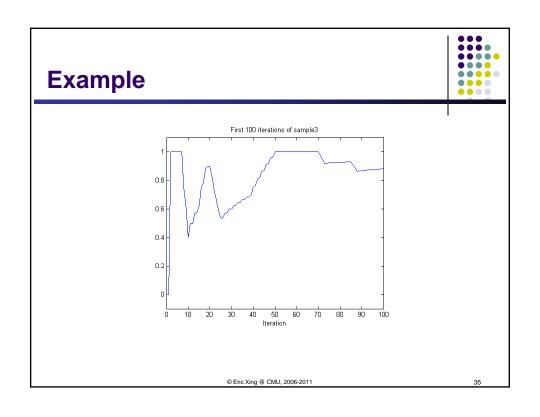


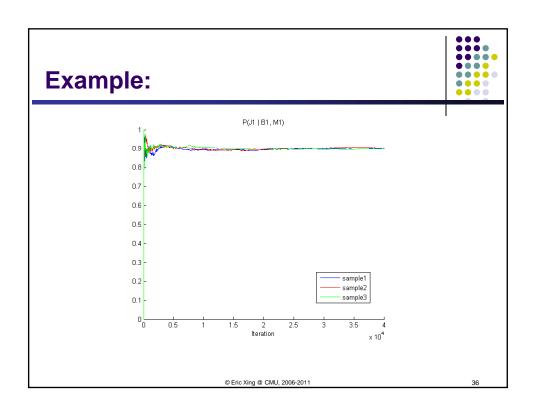


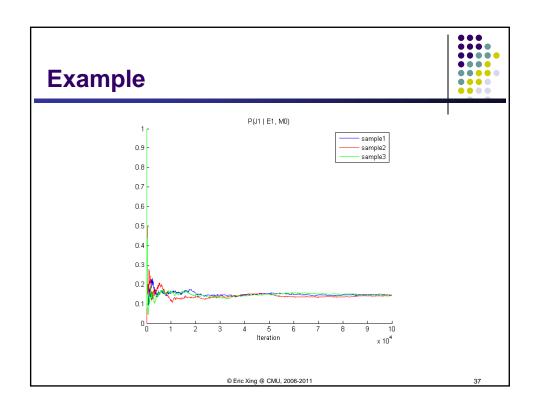
- To calculate P(J|B1,M1)
- Choose (B1,E0,A1,M1,J1) as a start
- Evidences are B1, M1, variables are A, E, J.
- Choose next variable as A
- Sample A by P(A|MB(A))=P(A|B1, E0, M1, J1) suppose to be false.
- (B1, E0, A0, M1, J1)
- Choose next random variable as E, sample E~P(E|B1,A0)

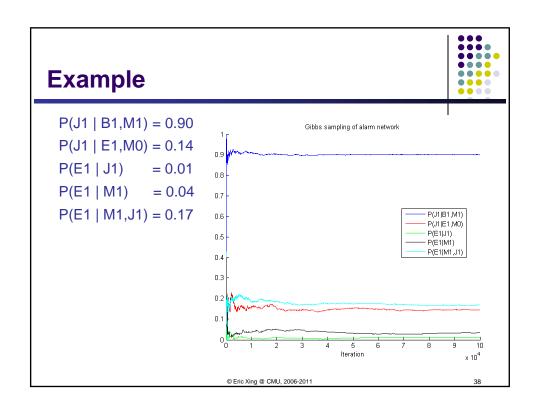
• ..

© Eric Xing @ CMU, 2006-2011









### **Outline**



- Dirichlet Process: From finite mixture to Infinite mixture model
  - The Chinese Restaurant Process
  - Example: heliotype inference
- Intro to Markov Chain Monte Carlo
  - Gibbs sampling
- Dynamic Dirichlet Process
  - The recurrent Chinese Restaurant Process
  - Example: Application: evolutionary clustering of documents

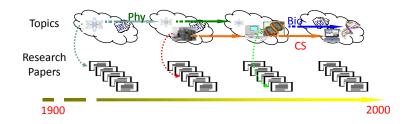
© Eric Xing @ CMU, 2006-2011

39

### **Evolutionary Clustering**



- Adapts the number of mixture components over time
  - Mixture components can die out
  - New mixture components are born at any time
  - Retained mixture components parameters evolve according to a Markovian dynamics



© Eric Xing @ CMU, 2006-2011

## **Temporal DPM**

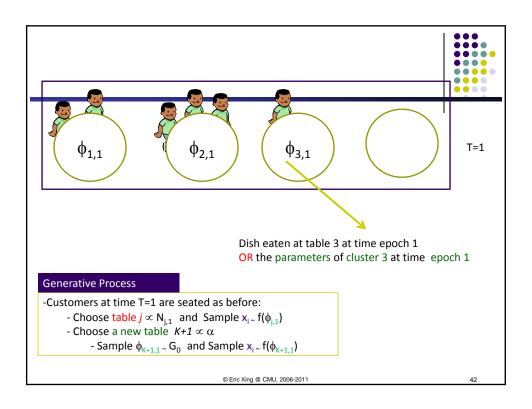
[Ahmed and Xing 2008]

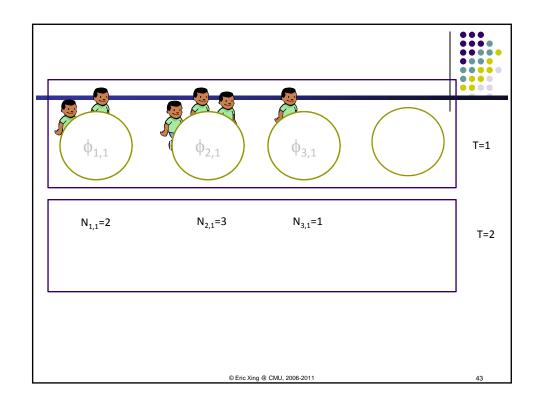


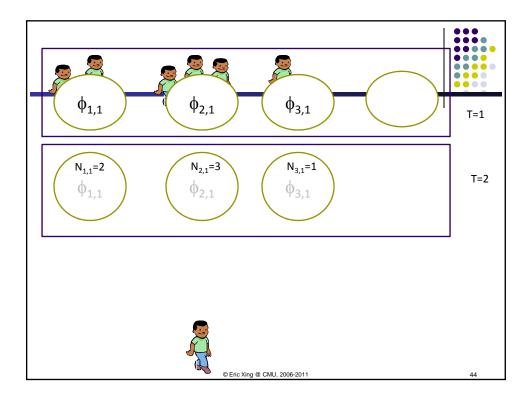
#### • The Recurrent Chinese Restaurant Process

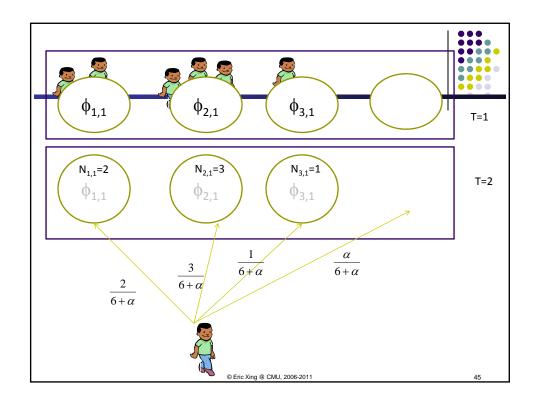
- The restaurant operates in epochs
- The restaurant is closed at the end of each epoch
- The state of the restaurant at time epoch t depends on that at time epoch t-1
  - Can be extended to higher-order dependencies.

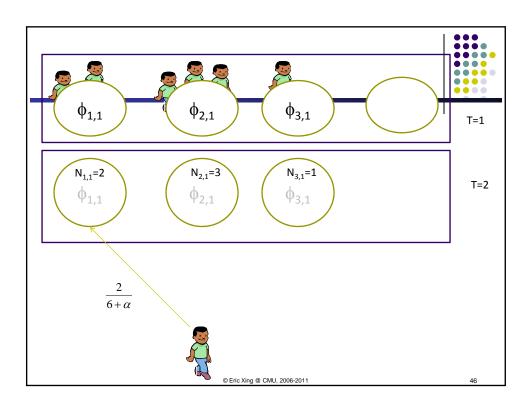
© Eric Xing @ CMU, 2006-2011

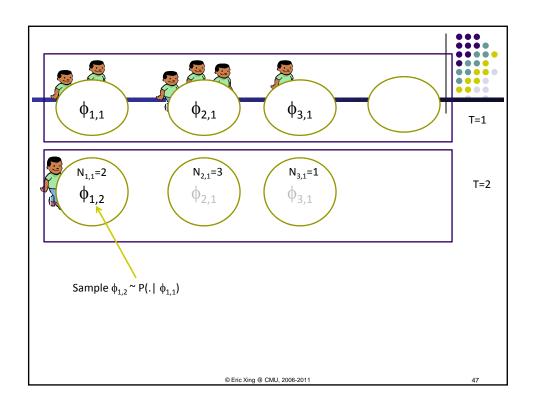


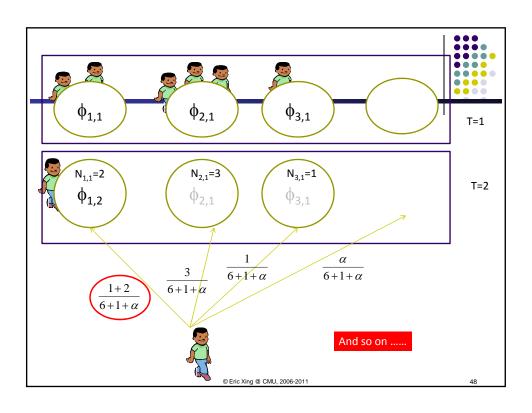


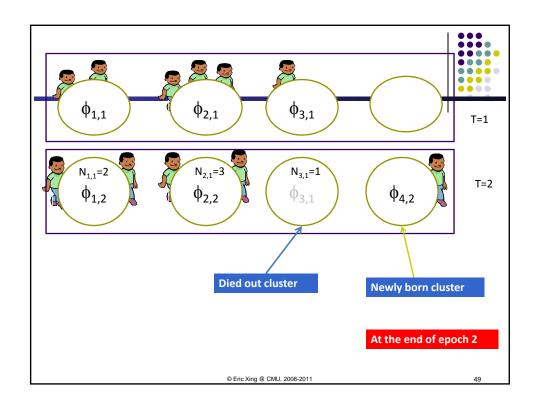


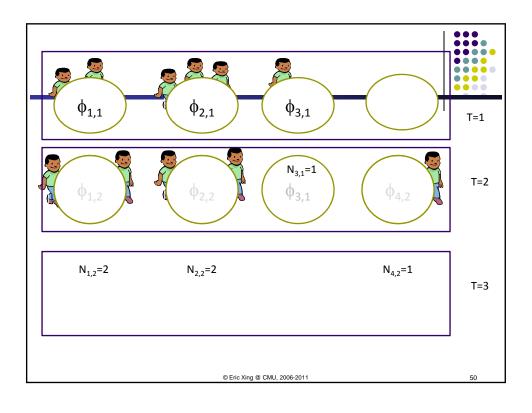


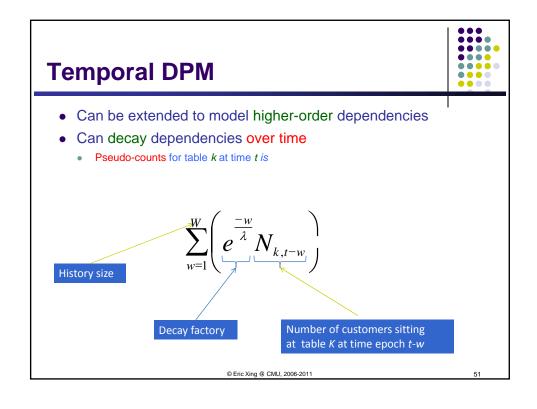


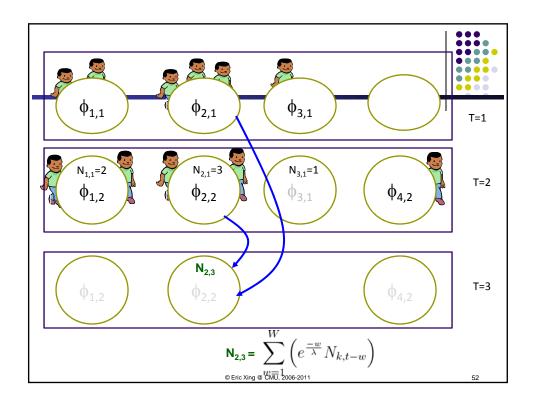


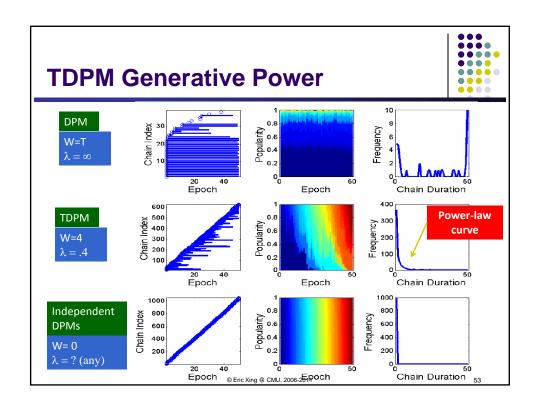










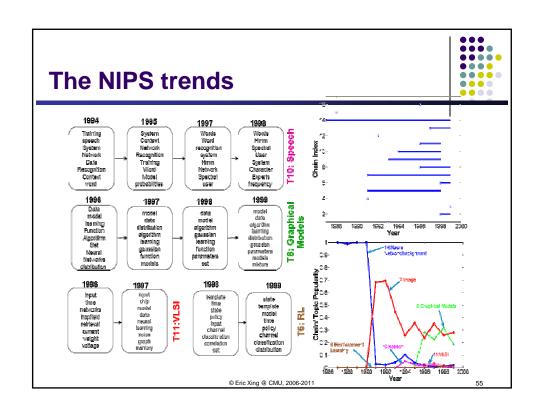


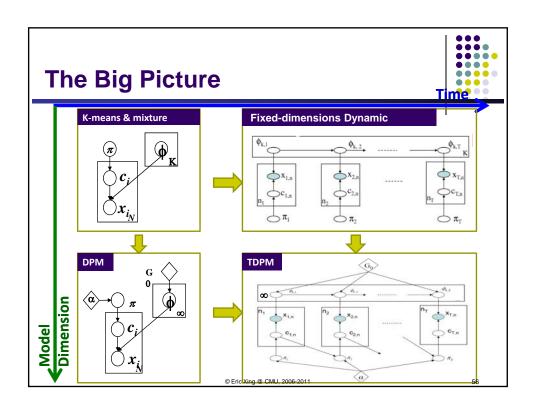
#### **Results: NIPS 12**



- Building a simple dynamic topic model
- Chain dynamics is as before
- Emission model for document  $x_{k,t}$  is:
  - Project  $\phi_k$ , over the simplex
  - Sample  $x_{k,l}|c_{t,i} \sim \text{Multinomial}(.|\text{Logisitic}(\phi_{k,t}))$
- Unlike LDA here a document belongs to one topic
- Use this model to analyze NIPS12 corpus
  - Proceeding of NIPS conference 1987-1999

© Eric Xing @ CMU, 2006-2011





## **Appendix:**



© Eric Xing @ CMU, 2006-2011

57

### **Theory of MCMC (optional)**



- **Definition:** Markov Chains
  - Given an n-dimensional state space
  - Random vector  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
  - $\mathbf{x}^{(t)} = \mathbf{x}$  at time-step t
  - $\mathbf{x}^{(t)}$  transitions to  $\mathbf{x}^{(t+1)}$  with prob  $P(\mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = T(\mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t)} \Rightarrow \mathbf{x}^{(t+1)})$
- Homogenous: chain determined by state x<sup>(0)</sup>, fixed transition kernel T (rows sum to 1)
- Equilibrium:  $\pi(\mathbf{x})$  is a stationary (equilibrium) distribution if  $\pi(\mathbf{x}') = \Sigma_{\mathbf{x}} \pi(\mathbf{x}) \ \mathsf{T}(\mathbf{x} \rightarrow \mathbf{x}').$

i.e., is a left eigenvector of the transition matrix  $\pi^{T}T = \pi^{T}T$ .

$$(0.2 \quad 0.5 \quad 0.3) = (0.2 \quad 0.5 \quad 0.3) \begin{pmatrix} 0.25 & 0 & 0.75 \\ 0 & 0.7 & 0.3 \\ 0.5 & 0.5 & 0 \end{pmatrix}$$

0.25 0.7 (X<sup>1</sup>) (X<sup>2</sup>) 0.75 0.5 0.3

© Eric Xing @ CMU, 2006-2011

#### **Markov Chains**



- An MC is *irreducible* if transition graph connected
- An MC is aperiodic if it is not trapped in cycles
- An MC is ergodic (regular) if you can get from state x to x'
  in a finite number of steps.
- **Detailed balance**:  $prob(x^{(t)} \rightarrow x^{(i-1)}) = prob(x^{(t-1)} \rightarrow x^{(t)})$

$$p(\mathbf{x}^{(t)})T(\mathbf{x}^{(t-1)} \mid \mathbf{x}^{(t)}) = p(\mathbf{x}^{(t-1)})T(\mathbf{x}^{(t)} \mid \mathbf{x}^{(t-1)})$$

summing over  $\mathbf{x}^{(t-1)}$ 

$$p(\mathbf{x}^{(t)}) = \sum_{\mathbf{x}^{(t-1)}} p(\mathbf{x}^{(t-1)}) T(\mathbf{x}^{(t)} \mid \mathbf{x}^{(t-1)})$$

Detailed bal → stationary dist exists

© Eric Xing @ CMU, 2006-2011

59

### **MCMC Via Metropolis-Hastings**



- Treat the target distribution as stationary distribution
- Sample from an easier proposal distribution, followed by an acceptance test
- This induces a transition matrix that satisfies detailed balance
  - MH proposes moves according to Q(x \ x) and accepts samples with probability A(x \ x).
  - The induced transition matrix is

$$T(X \to X') = Q(X'|X)A(X'|X)$$

Detailed balance means

$$\pi(\boldsymbol{x})Q(\boldsymbol{x}'|\,\boldsymbol{x})A(\boldsymbol{x}'|\,\boldsymbol{x}) = \pi(\boldsymbol{x}')Q(\boldsymbol{x}\,|\,\boldsymbol{x}')A(\boldsymbol{x}\,|\,\boldsymbol{x}')$$

• Hence the acceptance ratio is

$$A(x'|x) = \min\left(1, \frac{\pi(x')Q(x|x')}{\pi(x)Q(x'|x)}\right)$$

© Eric Xing @ CMU, 2006-2011



- · Gibbs sampling is a special case of MH
- The transition matrix updates each node one at a time using the following proposal:

$$Q((\mathbf{X}_i, \mathbf{X}_{-i}) \to (\mathbf{X}_i', \mathbf{X}_{-i})) = p(\mathbf{X}_i' | \mathbf{X}_{-i})$$

- This is efficient since for two reasons
  - It leads to samples that is always accepted

$$\begin{split} A\Big((\boldsymbol{x}_{i}, \mathbf{x}_{-i}) \rightarrow (\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'})\Big) &= \min \left(1, \frac{p(\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'})Q\big((\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'}) \rightarrow (\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'})\big)}{p(\boldsymbol{x}_{i}, \mathbf{x}_{-i}^{'})Q\big((\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'}) \rightarrow (\boldsymbol{x}_{i}^{'}, \mathbf{x}_{-i}^{'})\big)}\right) \\ &= \min \left(1, \frac{p(\boldsymbol{x}_{i}^{'}|\mathbf{x}_{-i}^{'})p(\mathbf{x}_{-i}^{'})p(\boldsymbol{x}_{i}^{'}|\mathbf{x}_{-i}^{'})}{p(\boldsymbol{x}_{i}^{'}|\mathbf{x}_{-i}^{'})p(\mathbf{x}_{-i}^{'})p(\boldsymbol{x}_{-i}^{'}|\mathbf{x}_{-i}^{'})}\right) = \min(1,1) \end{split}$$

Thus

$$T((\mathbf{X}_i, \mathbf{X}_{-i}) \to (\mathbf{X}_i', \mathbf{X}_{-i})) = p(\mathbf{X}_i' | \mathbf{X}_{-i})$$

• It is efficient since  $p(\mathbf{x}_i^{\cdot} | \mathbf{x}_{-i})$  only depends on the values in  $\mathcal{X}_i^{\cdot}$ s Markov blanket

© Eric Xing @ CMU, 2006-2011