

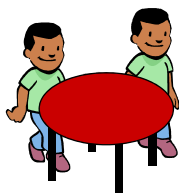
Machine Learning

10-701/15-781, Fall 2011

Infinite Mixture Models

Eric Xing

Lecture 10, October 12, 2011



© Eric Xing @ CMU, 2006-2011

1

Clustering



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

© Eric Xing @ CMU, 2006-2011

2

Image Segmentation

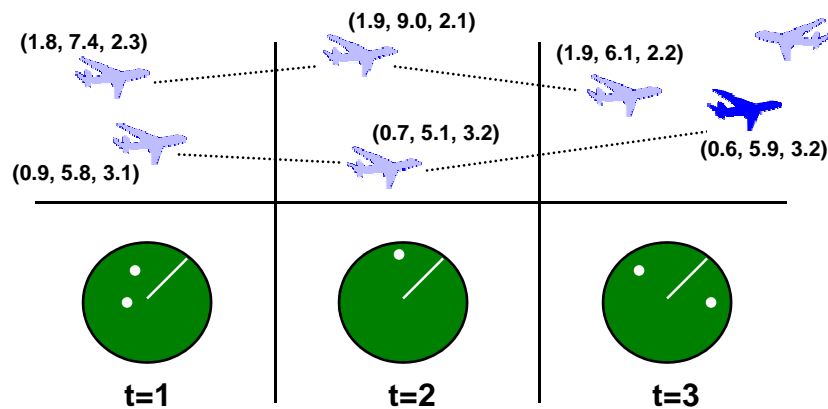


- How to segment images?
 - Manual segmentation (very expensive)
 - Algorithm segmentation
 - K-means
 - Statistical mixture models
 - Spectral clustering
- Problems with most existing algorithms
 - Ignore the spatial information
 - Perform the segmentation one image at a time
 - Need to specify the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

© Eric Xing @ CMU, 2006-2011

3

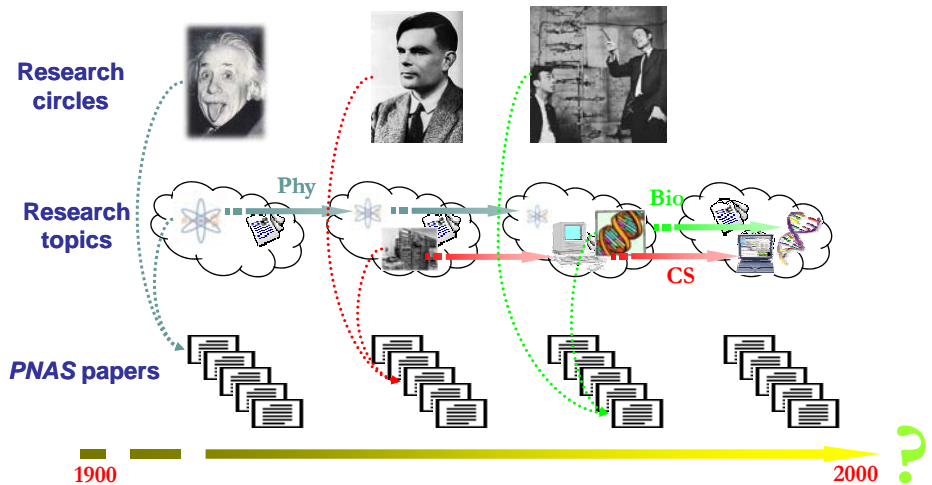
Object Recognition and Tracking



© Eric Xing @ CMU, 2006-2011

4

The Evolution of Science



© Eric Xing @ CMU, 2006-2011

5

Clustering as (nonparametric) Bayesian Inference

- Cluster = Mixture
- Each data point is generated from one Mixture

$$z_i \sim P(z|\pi)$$

$$x_i | z_i \sim P(x|\theta_{z_i})$$



How Many clusters

- Cross Validation
 - Data hungry!
- Information theoretic
 - AIC, MDL, etc.
- Non-Parametric
 - Manage model uncertainty
 - Integrate over different clustering configurations
 - Base for many interesting extensions

© Eric Xing @ CMU, 2006-2011

6

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

Clustering

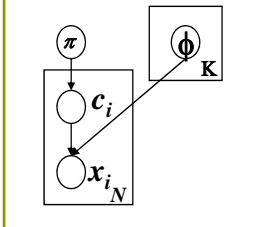


- How to label them ?
- How many clusters ???

Mixture model: The Generative Process

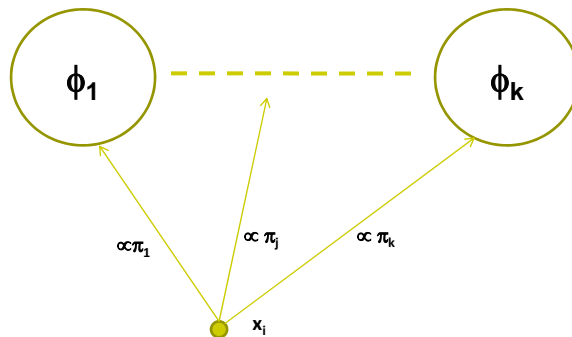


Mixture of G



Generative Process

- For data point x_i
 - Sample $c_i \sim \text{Multi}(\pi)$
 - Sample $x_i \sim f(\phi_{c_i})$



Number of clusters DOES NOT grow with the data

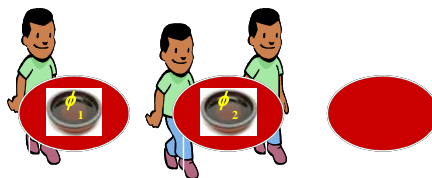
© Eric Xing @ CMU, 2006-2011

9

The Chinese Restaurant Process



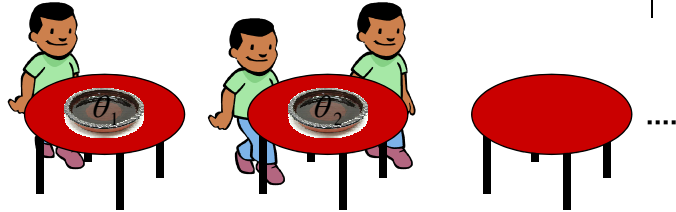
- Customers correspond to data points
- Tables correspond to clusters/mixture components
- Dishes correspond to parameter of the mixtures



© Eric Xing @ CMU, 2006-2011

10

The Chinese Restaurant Process



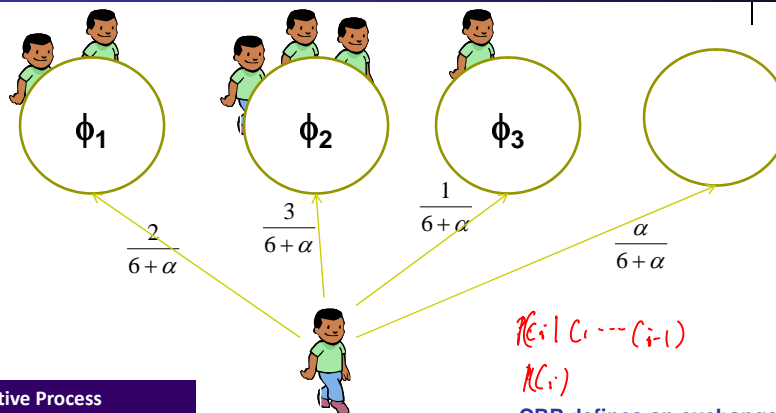
$$P(c_i = k | \mathbf{c}_{-i}) = \begin{array}{ccc} \frac{1}{1+\alpha} & \frac{0}{1+\alpha} & \frac{0}{1+\alpha} \\ \frac{1}{2+\alpha} & \frac{1}{2+\alpha} & \frac{\alpha}{2+\alpha} \\ \frac{1}{3+\alpha} & \frac{2}{3+\alpha} & \frac{\alpha}{3+\alpha} \\ \frac{m_1}{i+\alpha-1} & \frac{m_2}{i+\alpha-1} & \dots \frac{\alpha}{i+\alpha-1} \end{array}$$

The rich gets richer effect

© Eric Xing @ CMU, 2006-2011

11

The Chinese Restaurant Process



Generative Process

- For data point x_i
 - Choose table $j \propto N_j$ and Sample $x_i \sim f(\phi_j)$
 - Choose a new table $K+1 \propto \alpha$
 - Sample $\phi_{K+1} \sim G_0$ and Sample $x_i \sim f(\phi_{K+1})$

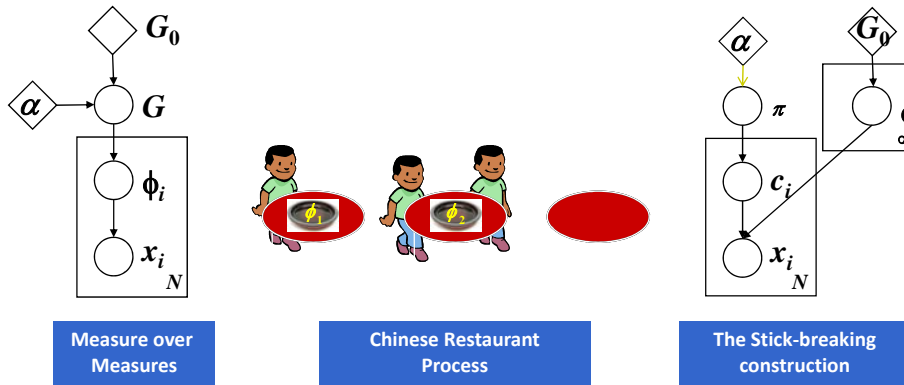
CRP defines an exchangeable distribution on partitions over an (infinite) sequence of samples, such a distribution is formally known as the Dirichlet Process (DP)

© Eric Xing @ CMU, 2006-2011

12

More formally ...

- This is a Dirichlet Process Mixture Model
- Three equivalent constructions:



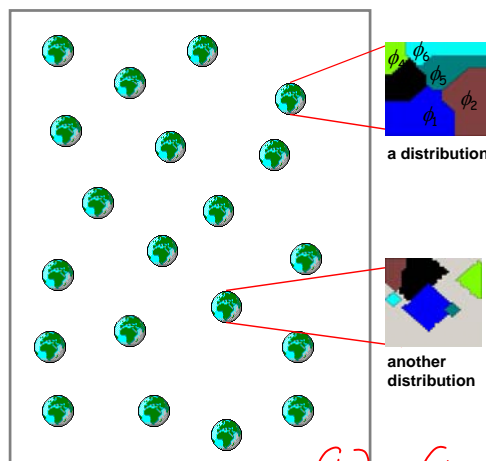
© Eric Xing @ CMU, 2006-2011

13

Dirichlet Process

$$D_{\alpha} \equiv P(\theta_1, \theta_2, \theta_3, \dots | \alpha)$$

$$1 > \theta_i > 0 \quad \sum \theta_i = 1$$



- A CDF, G , on possible worlds of random partitions follows a Dirichlet Process if for any measurable finite partition $(\phi_1, \phi_2, \dots, \phi_m)$:

$$(G(\phi_1), G(\phi_2), \dots, G(\phi_m)) \sim \text{Dirichlet}(\alpha G_0(\phi_1), \dots, \alpha G_0(\phi_m))$$

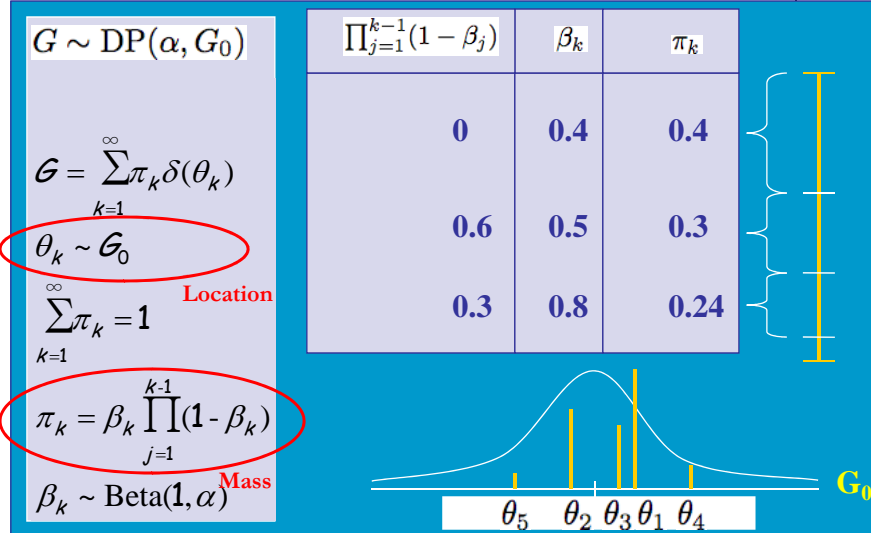
where G_0 is the base measure and α is the scale parameter

Mixture: Thus a Dirichlet Process G defines a distribution of distribution

© Eric Xing @ CMU, 2006-2011

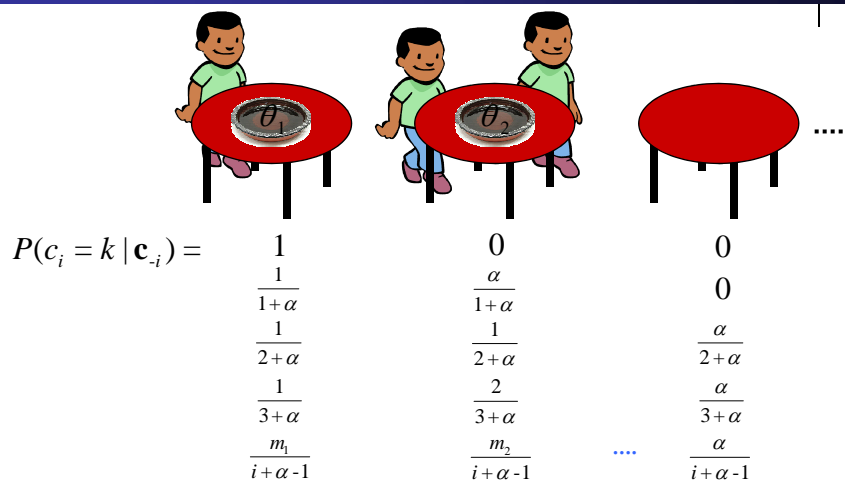
14

Stick-breaking Process



15

Chinese Restaurant Process



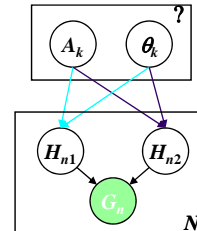
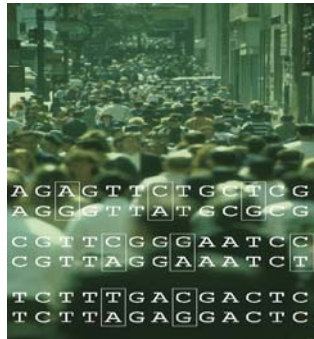
CRP defines an exchangeable distribution on partitions over an (infinite) sequence of samples, such a distribution is formally known as the Dirichlet Process (DP)

© Eric Xing @ CMU, 2006-2011

16

Case study: ancestral Inference

Xing et al 2004



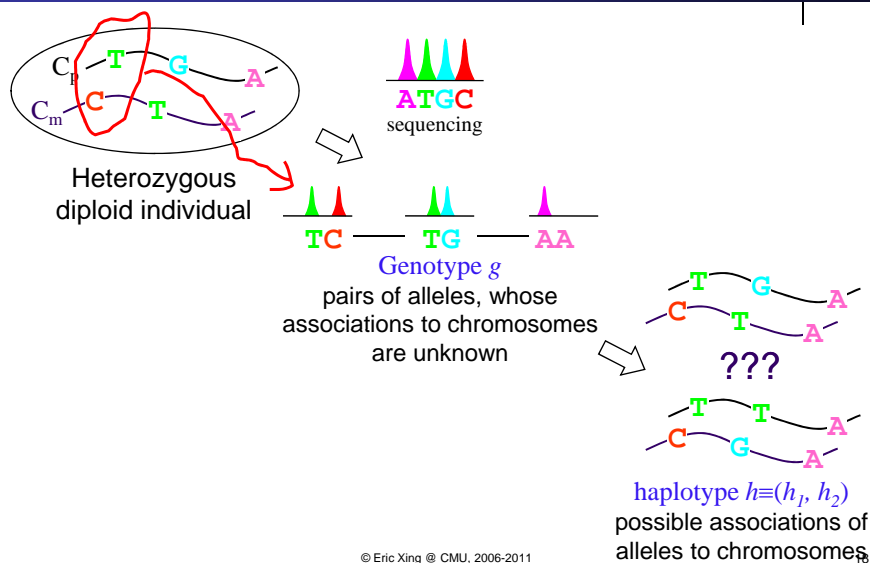
Essentially a clustering problem, but ...

- 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

17

Background:haplotype ambiguity



© Eric Xing @ CMU, 2006-2011

A Finite (Mixture of) Allele Model

- The probability of a genotype g :

$$p(g) = \sum_{h_1, h_2 \in \mathcal{H}} p(h_1, h_2) p(g | h_1, h_2)$$

- Standard settings:

- $|\mathcal{H}| = K \ll 2^J$ fixed-sized population haplotype pool
- $p(h_1, h_2) = p(h_1)p(h_2) = f_1 f_2$ Hardy-Weinberg equilibrium

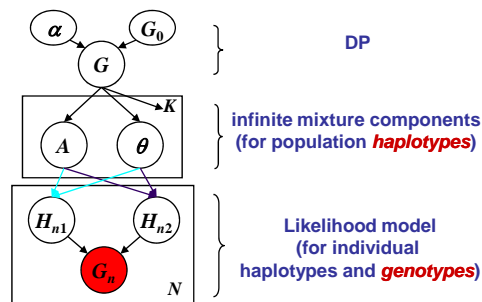
- Problem: $K?$ $\mathcal{H}?$

© Eric Xing @ CMU, 2006-2011

19

Example: DP-haplotype [Xing et al, 2004]

- Clustering human populations



- Inference: Markov Chain Monte Carlo (MCMC)

- Gibbs sampling
- Metropolis Hasting

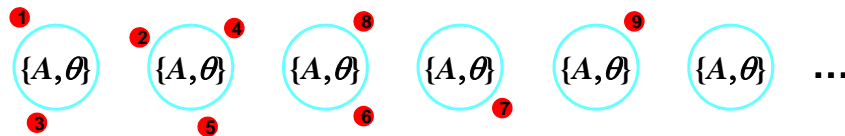
© Eric Xing @ CMU, 2006-2011

20

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

21

MCMC for Haplotype Inference



- Gibbs sampling for exploring the posterior distribution under the proposed model
 - Integrate out the parameters such as θ or λ , and sample c_{i_e}, a_k and h_{i_e}

$$p(c_{i_e} = k | \mathbf{c}_{[-i_e]}, \mathbf{h}, \mathbf{a}) \propto \underbrace{p(c_{i_e} = k | \mathbf{c}_{[-i_e]})}_{\text{Prior}} \times \underbrace{p(h_{i_e} | a_k, \mathbf{h}_{[-i_e]}, \mathbf{c})}_{\text{Likelihood}}$$

⋮

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

22

MCMC for Haplotype Inference



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

$$p(c_{ie}^{(j)} = k | \mathbf{c}^{[-j, ie]}, \mathbf{h}, \mathbf{a})$$

$$\propto p(c_{ie}^{(j)} = k | \mathbf{c}^{[-j, ie]}, \mathbf{m}, \mathbf{n}) p(h_{ie}^{(j)} | a_k, \mathbf{c}, \mathbf{h}^{[-j, ie]})$$

$$\propto (m_{jk}^{[-j, ie]} + \tau \beta_k) p(h_{ie}^{(j)} | a_k, \mathbf{l}_k^{[-j, ie]}), \text{ for } k = 1, \dots, K + 1$$
 2. Sample a_k from

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

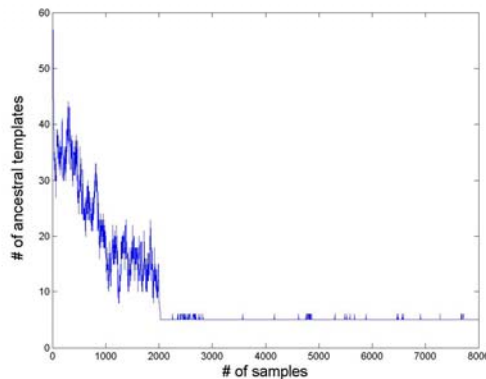
$$= \frac{\Gamma(\alpha_h + l_{k,t}) \Gamma(\beta_h + l'_{k,t})}{\Gamma(\alpha_h + \beta_h + m_k) (|B| - 1)^{l_{k,t}}} R(\alpha_h, \beta_h)$$
 3. Sample $h_{ie}^{(j)}$ from

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

© Eric Xing @ CMU, 2006-2011

23

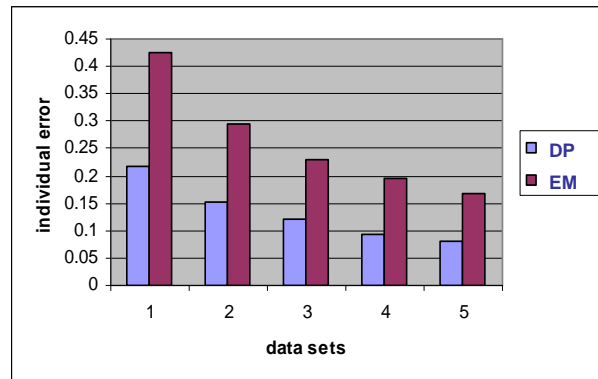
Convergence of Ancestral Inference



© Eric Xing @ CMU, 2006-2011

24

DP vs. Finite Mixture via EM



© Eric Xing @ CMU, 2006-2011

25

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

$$P(x_i | \theta)$$

$$\approx N(\cdot)$$

$$P(h | \theta, a)$$

$$P(x | \phi_i)$$

© Eric Xing @ CMU, 2006-2011

26

Monte Carlo methods

$$F = \int f(x) p(x) dx$$



- 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

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

$$x \sim p(x)$$

- **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?

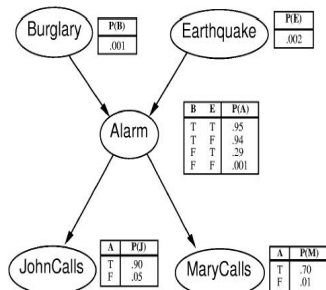
© Eric Xing @ CMU, 2006-2011

27

Example: naive sampling



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



Alarm example: (Choose the right sampling sequence)

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

© Eric Xing @ CMU, 2006-2011

28



- | | | | | |
|----|----|----|----|----|
| E0 | B0 | A0 | M0 | J0 |
| E0 | B0 | A0 | M0 | J0 |
| E0 | B0 | A0 | M0 | J0 |
| 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 |

29



- | | | | | |
|----|----|----|----|----|
| 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 |
| 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.
-
- Markov chain Monte Carlo (MCMC)

© Eric Xing @ CMU, 2006-2011

31

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

32

Gibbs sampling

- Gibbs sampling is an MCMC algorithm that is especially appropriate for inference in graphical models.
- The procedure
 - we have variable set $\mathbf{X} = \{X_1, X_2, X_3, \dots, X_N\}$ for a GM $p(\mathbf{X})$
 - at each step one of the variables X_i is selected (at random or according to some fixed sequences), denote the remaining variables as \mathbf{X}_{-i} , and its current value as $x_{-i}^{(t-1)}$
 - Using the "alarm network" as an example, say at time t we choose X_E and we denote the current value assignments of the remaining variables, \mathbf{X}_{-E} , $p(\vec{X})$ obtained from previous samples, as $\mathbf{x}_{-E}^{(t-1)} = \{x_B^{(t-1)}, x_A^{(t-1)}, x_J^{(t-1)}, x_M^{(t-1)}\}$
 - 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 $\mathbf{x}^{(t)}$ replaces the previous sampled value of X_i in $\mathbf{x}^{(t-1)}$
 - i.e., $\mathbf{x}^{(t)} = \mathbf{x}_{-E}^{(t-1)} \cup x_E^{(t)}$

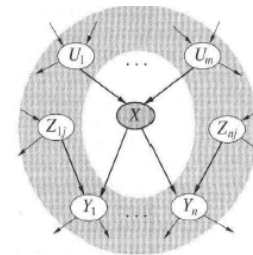
© Eric Xing @ CMU, 2006-2011

33

Why Gibbs Sampling is Simple

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

$$\Rightarrow p(X_i | \mathbf{X}_{-i}) = p(X_i | \text{MB}(X_i))$$



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

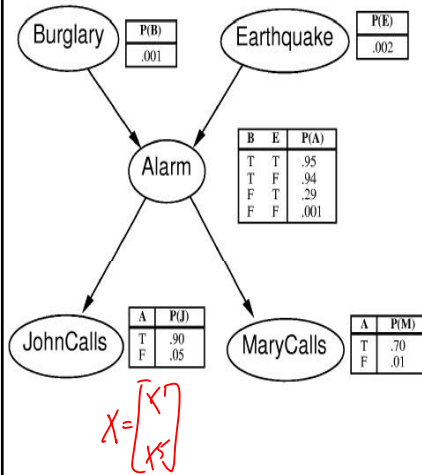
$$\text{MB}(A) = \{B, E, J, M\}$$

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

© Eric Xing @ CMU, 2006-2011

34

Example: alarm network again

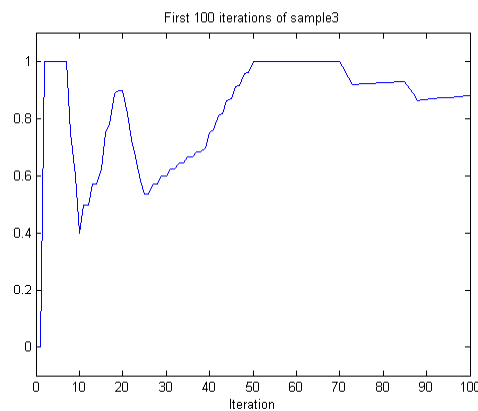


- To calculate $P(J|B1, M1)$
- Choose $(B1, E0, A1, M1, J1)$ as a start
- **Evidences** are $B1, M1$, **variables** are A, E, J .
- Choose next variable as A
- Sample A by $P(A|MB(A)) = P(A|B1, E0, M1, J1)$ suppose to be false.
- $(B1, E0, A0, M1, J1)$
- Choose next random variable as E , sample $E \sim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

© Eric Xing @ CMU, 2006-2011

35

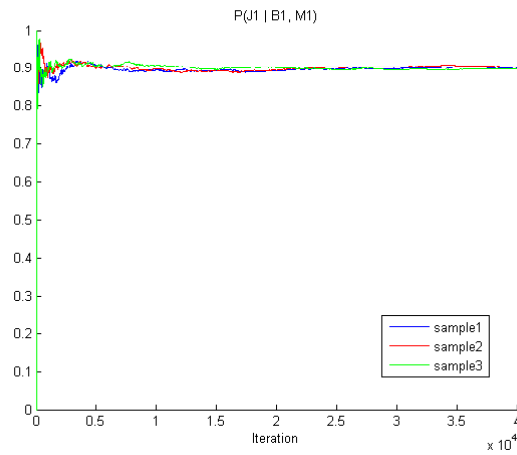
Example



© Eric Xing @ CMU, 2006-2011

36

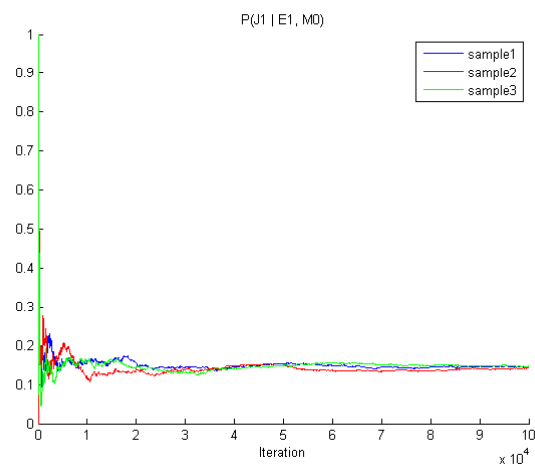
Example:



© Eric Xing @ CMU, 2006-2011

37

Example



© Eric Xing @ CMU, 2006-2011

38

Example

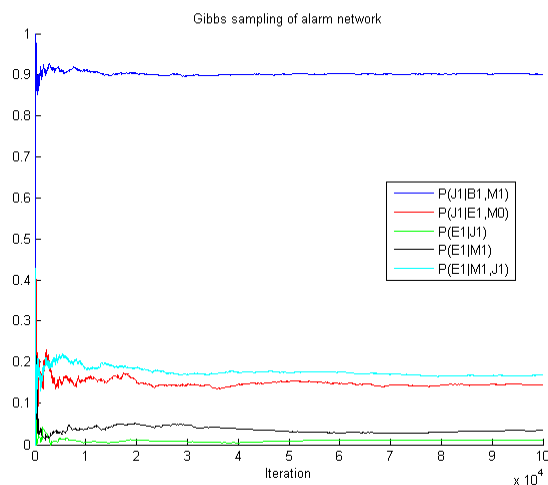
$$P(J1 \mid B1, M1) = 0.90$$

$$P(J1 \mid E1, M0) = 0.14$$

$$P(E1 \mid J1) = 0.01$$

$$P(E1 \mid M1) = 0.04$$

$$P(E1 \mid M1, J1) = 0.17$$



© Eric Xing @ CMU, 2006-2011

39

Outline

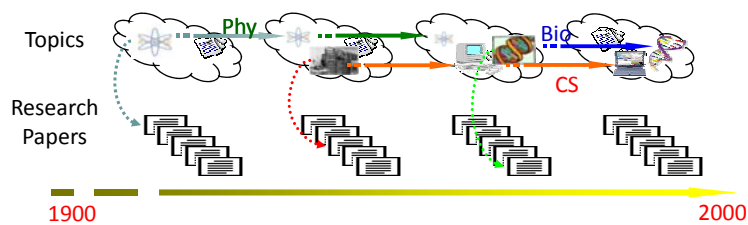
- Dirichlet Process: From finite mixture to Infinite mixture model
 - The Chinese Restaurant Process
 - Example: heliotype inference
- Intro to Markov Chain Monte Carlo
 - Gibbs sampling
- Dynamic Dirichlet Process
 - The recurrent Chinese Restaurant Process
 - Example: Application: evolutionary clustering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

© Eric Xing @ CMU, 2006-2011

40

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

41

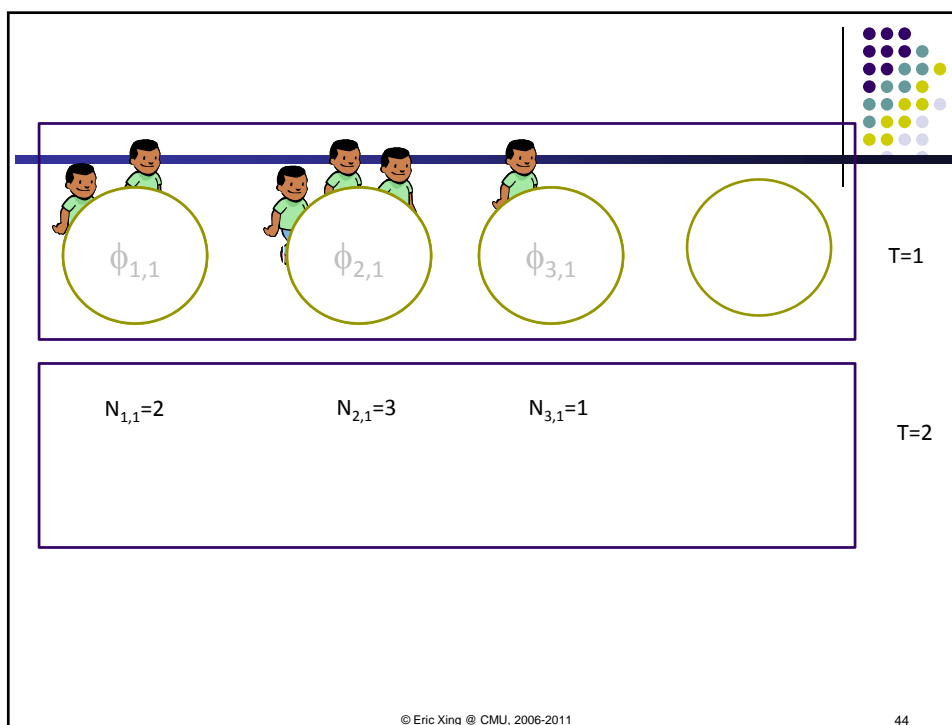
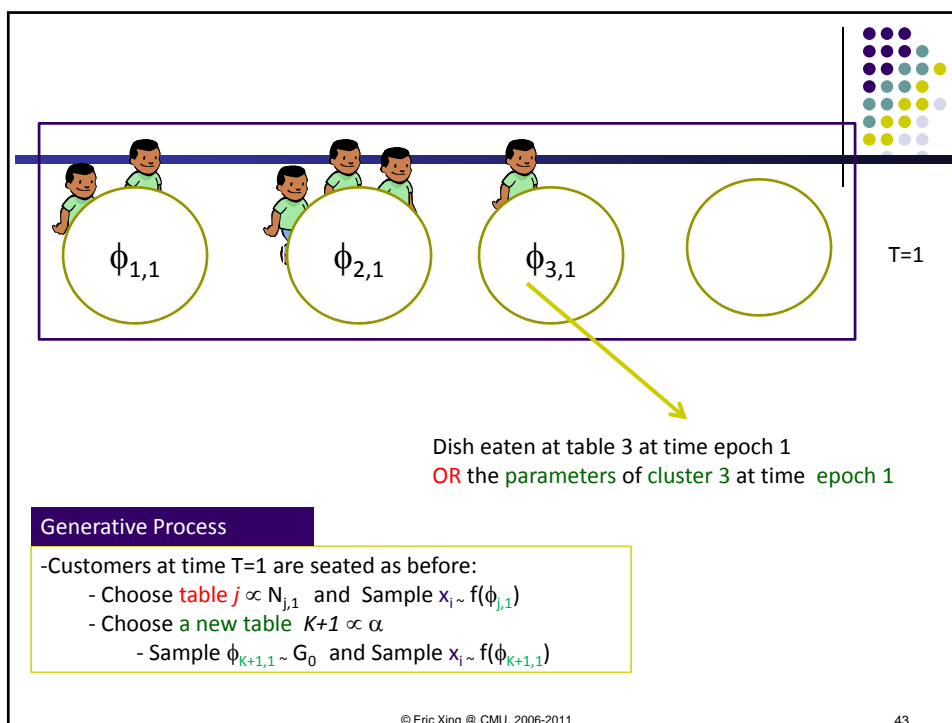
Temporal DPM

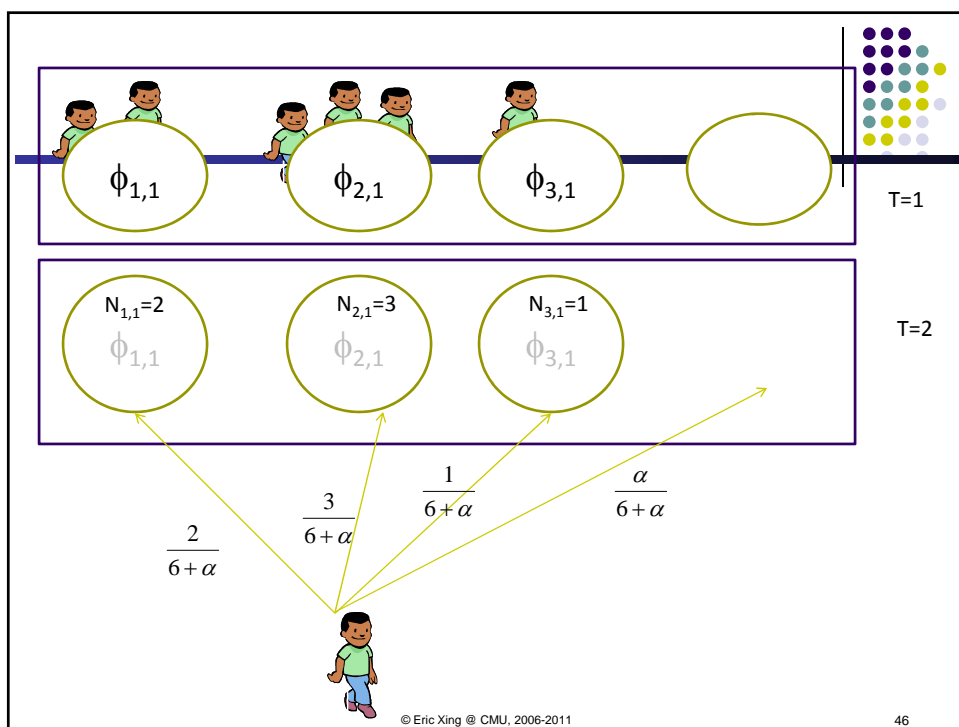
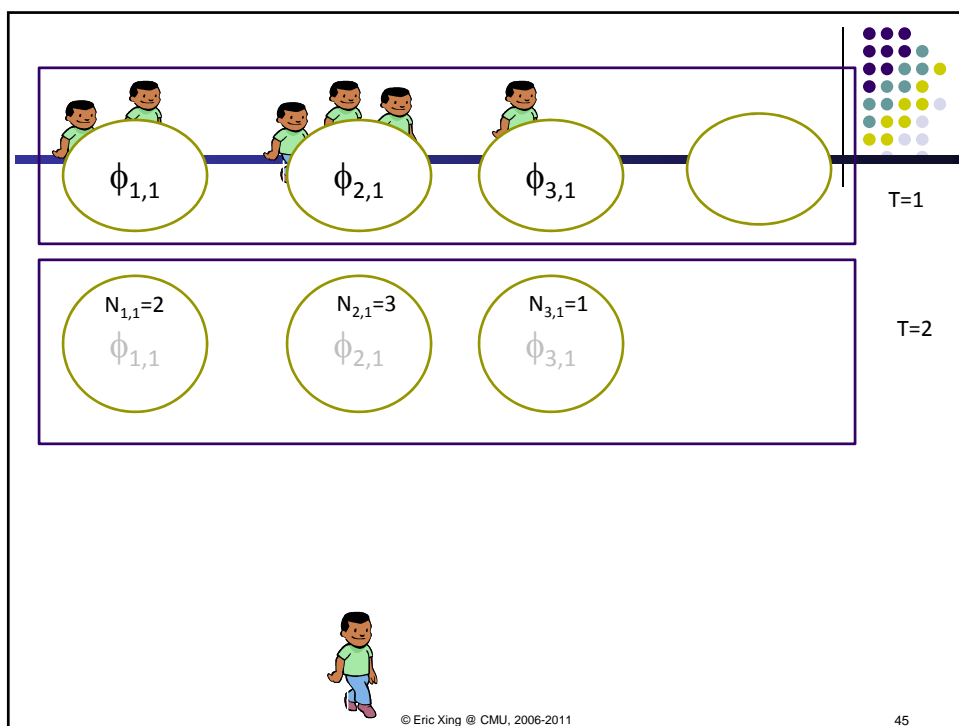
[Ahmed and Xing 2008]

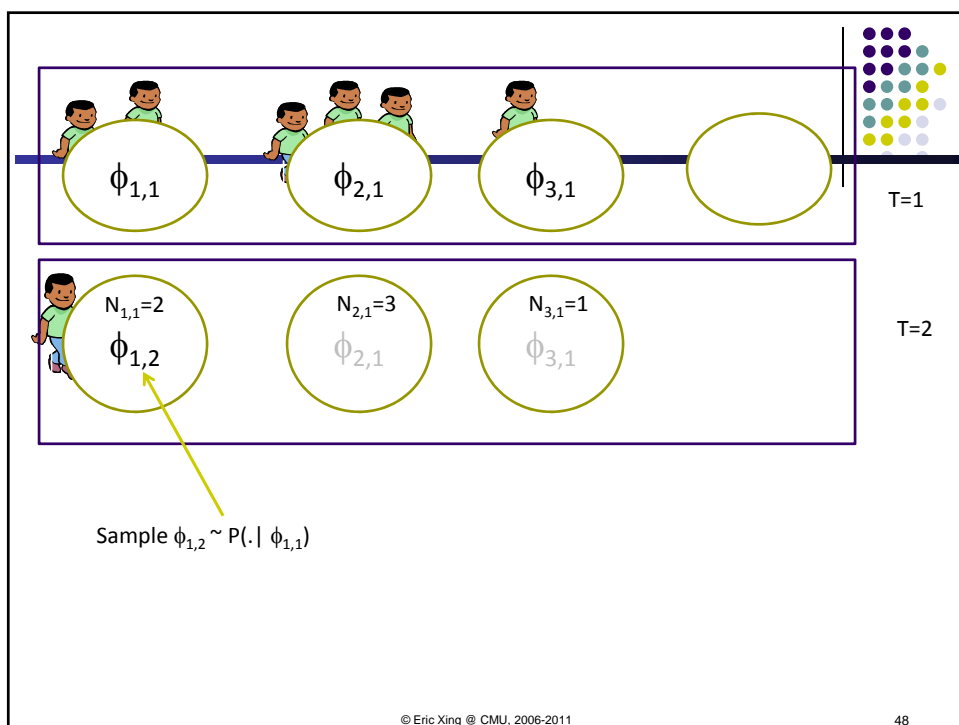
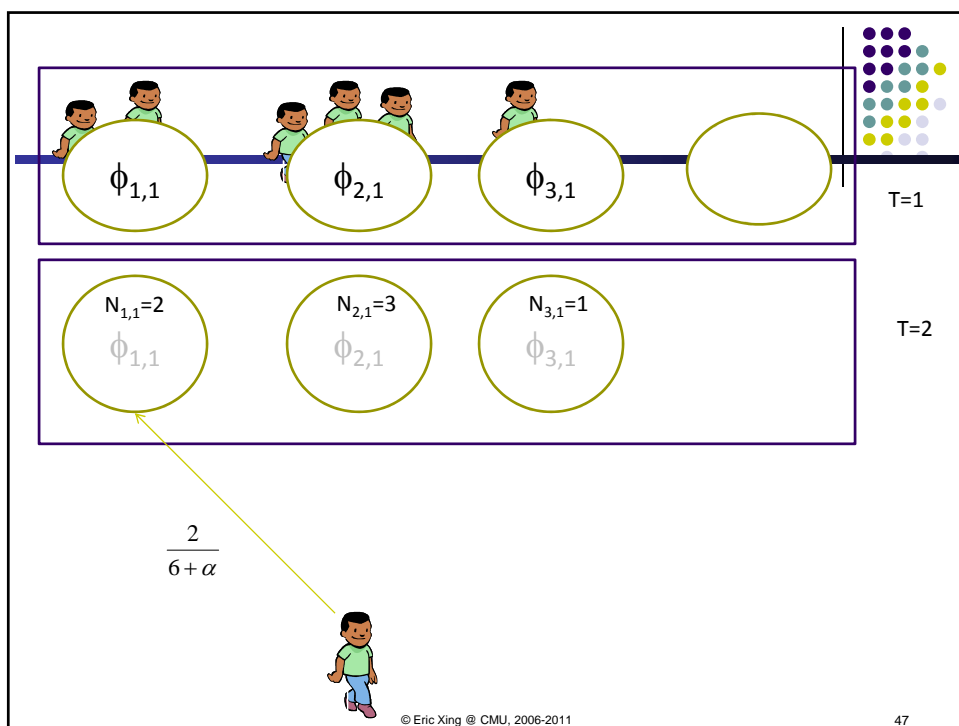
- The Recurrent Chinese Restaurant Process
 - The restaurant operates in epochs
 - The restaurant is closed at the end of each epoch
 - The state of the restaurant at time epoch t depends on that at time epoch $t-1$
 - Can be extended to higher-order dependencies.

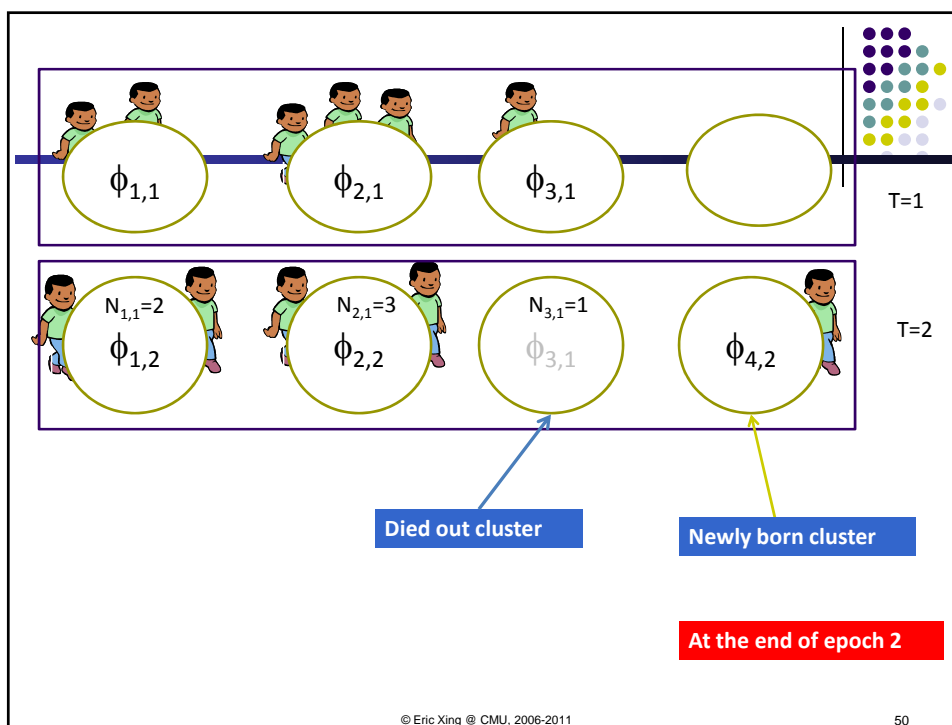
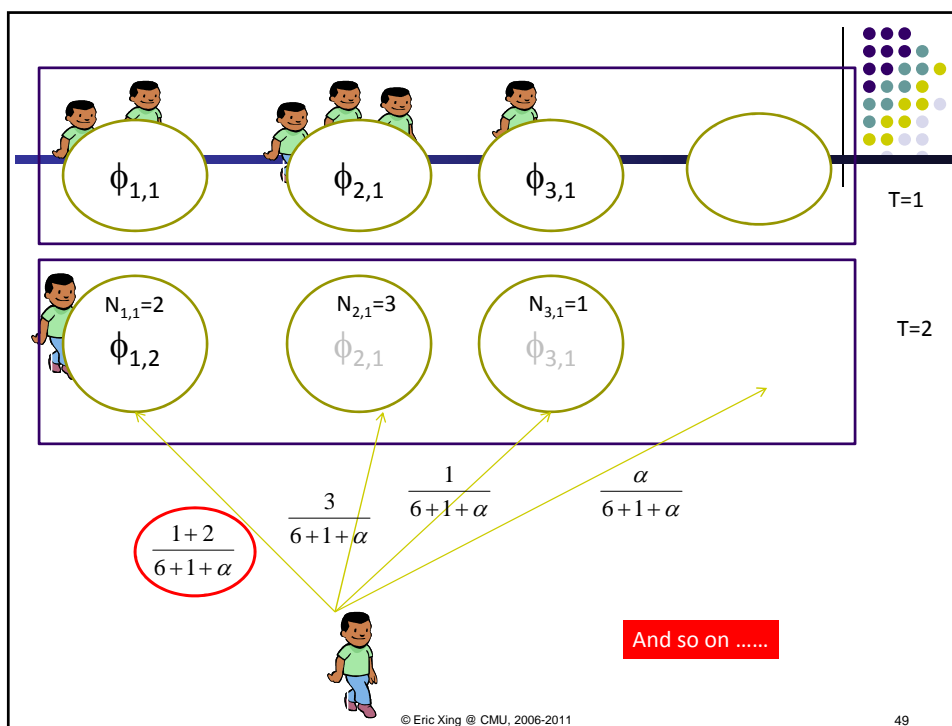
© Eric Xing @ CMU, 2006-2011

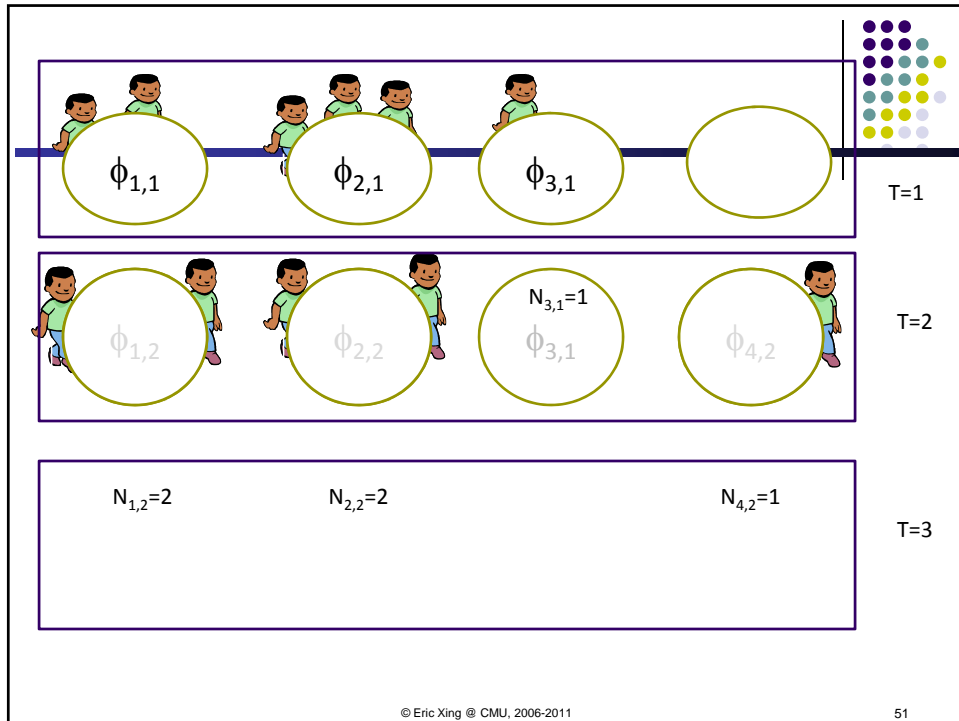
42











Temporal DPM

- Can be extended to model **higher-order** dependencies
- Can **decay** dependencies **over time**
 - Pseudo-counts for table k at time t is

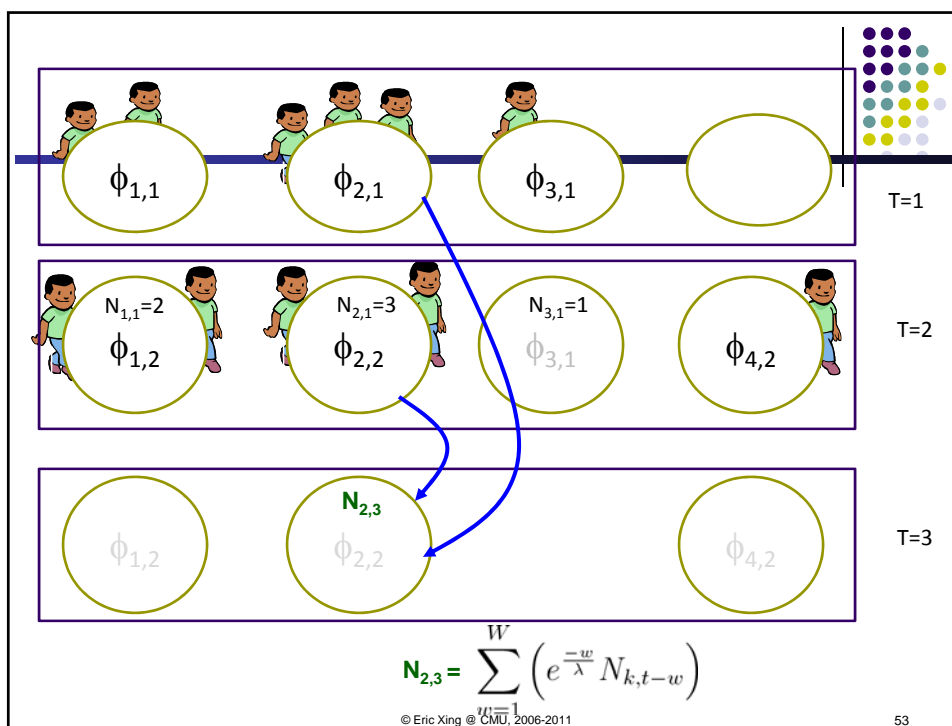
$$\sum_{w=1}^W \left(e^{\frac{-w}{\lambda}} N_{k,t-w} \right)$$

History size

Decay factory

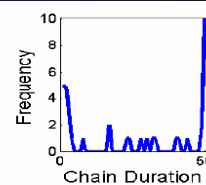
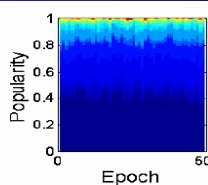
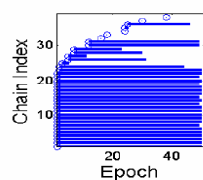
Number of customers sitting at table K at time epoch $t-w$

© Eric Xing @ CMU, 2006-2011 52

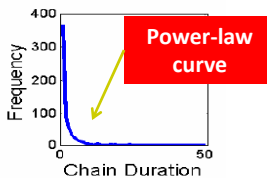
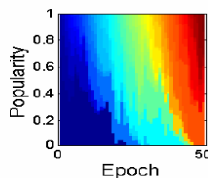
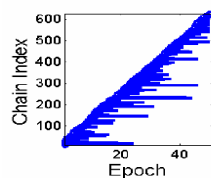


TDPM Generative Power

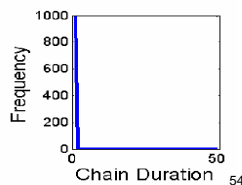
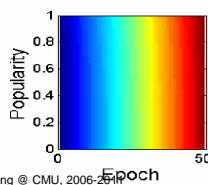
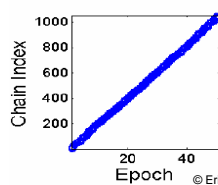
DPM
 $W=T$
 $\lambda = \infty$



TDPM
 $W=4$
 $\lambda = .4$



Independent DPMs
 $W=0$
 $\lambda = ?$ (any)



© Eric Xing @ CMU, 2006-2011

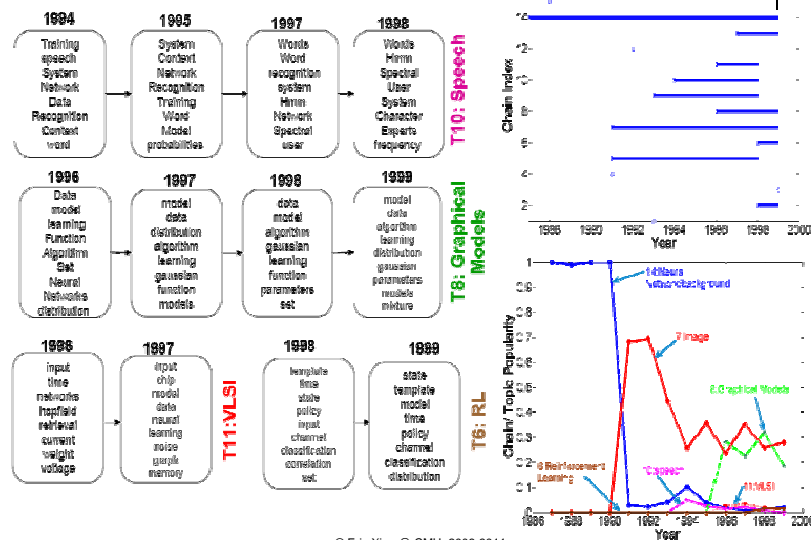
Results: NIPS 12

- Building a **simple** dynamic **topic** model
- Chain dynamics is as before
- Emission model for document $x_{k,t}$ is:
 - Project $\phi_{k,t}$ over the simplex
 - Sample $x_{k,t}|c_{k,t} \sim \text{Multinomial}(\cdot | \text{Logistic}(\phi_{k,t}))$
- 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

55

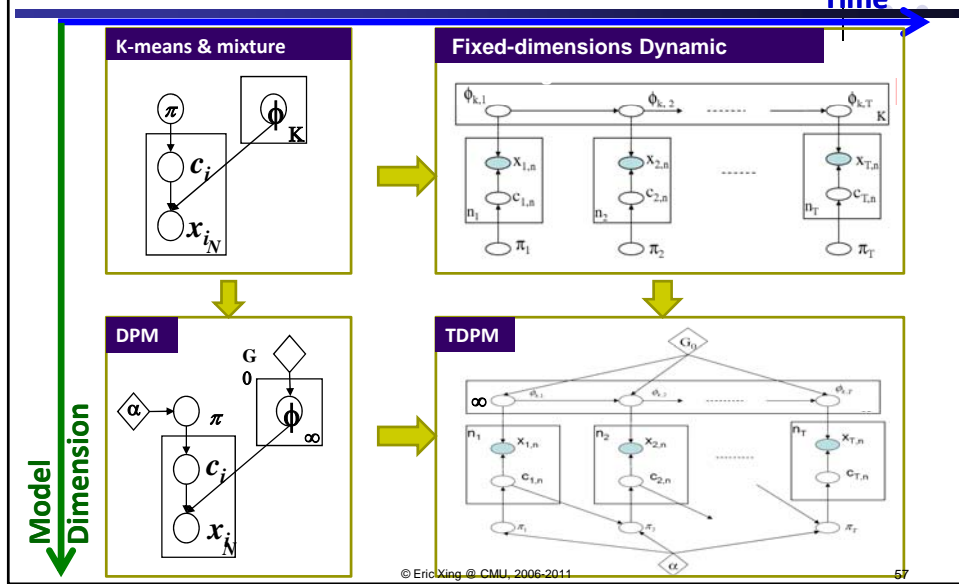
The NIPS trends



© Eric Xing @ CMU, 2006-2011

56

The Big Picture



Appendix:

Theory of MCMC (optional)

- **Definition:** Markov Chains

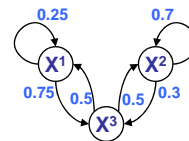
- Given an n-dimensional state space
- Random vector $\mathbf{X} = (X_1, \dots, X_n)$
- $\mathbf{x}^{(t)} = \mathbf{x}$ at time-step t
- $\mathbf{x}^{(t)}$ transitions to $\mathbf{x}^{(t+1)}$ with prob
 $P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t)} \rightarrow \mathbf{x}^{(t+1)})$

- **Homogenous:** chain determined by state $\mathbf{x}^{(0)}$, fixed *transition kernel* T (rows sum to 1)

- **Equilibrium:** $\pi(\mathbf{x})$ is a *stationary (equilibrium) distribution* if
 $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) T(\mathbf{x} \rightarrow \mathbf{x}')$.

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

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



© Eric Xing @ CMU, 2006-2011

59

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 \mathbf{x} to \mathbf{x}' in a finite number of steps.
- **Detailed balance:** $\text{prob}(\mathbf{x}^{(t)} \rightarrow \mathbf{x}^{(t+1)}) = \text{prob}(\mathbf{x}^{(t+1)} \rightarrow \mathbf{x}^{(t)})$

$$p(\mathbf{x}^{(t)}) T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = p(\mathbf{x}^{(t+1)}) T(\mathbf{x}^{(t)} | \mathbf{x}^{(t+1)})$$

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

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

- Detailed bal \rightarrow stationary dist exists

© Eric Xing @ CMU, 2006-2011

60

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 $T(x \rightarrow x') = Q(x'|x)A(x'|x)$
 - Detailed balance means $\pi(x)Q(x'|x)A(x'|x) = \pi(x')Q(x|x')A(x|x')$
 - 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

61

Gibbs sampling



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

$$Q((x_i, \mathbf{x}_{-i}) \rightarrow (x'_i, \mathbf{x}_{-i})) = p(x'_i | \mathbf{x}_{-i})$$

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

$$\begin{aligned} A((x_i, \mathbf{x}_{-i}) \rightarrow (x'_i, \mathbf{x}_{-i})) &= \min\left(1, \frac{p(x'_i | \mathbf{x}_{-i})Q((x'_i, \mathbf{x}_{-i}) \rightarrow (x_i, \mathbf{x}_{-i}))}{p(x_i | \mathbf{x}_{-i})Q((x_i, \mathbf{x}_{-i}) \rightarrow (x'_i, \mathbf{x}_{-i}))}\right) \\ &= \min\left(1, \frac{p(x'_i | \mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x_i | \mathbf{x}_{-i})}{p(x_i | \mathbf{x}_{-i})p(\mathbf{x}_{-i})p(x'_i | \mathbf{x}_{-i})}\right) = \min(1, 1) \end{aligned}$$

Thus $T((x_i, \mathbf{x}_{-i}) \rightarrow (x'_i, \mathbf{x}_{-i})) = p(x'_i | \mathbf{x}_{-i})$

- It is efficient since $p(x'_i | \mathbf{x}_{-i})$ only depends on the values in X_i 's Markov blanket

© Eric Xing @ CMU, 2006-2011

62