# Advanced Machine Learning
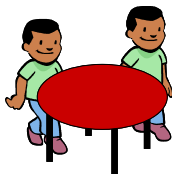
## Nonparametric Bayesian Models

### --Learning/Reasoning in Open Possible Worlds

**Eric Xing**

**Lecture 17, August 14, 2009**
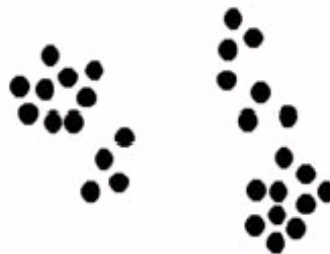
**Reading:**

---

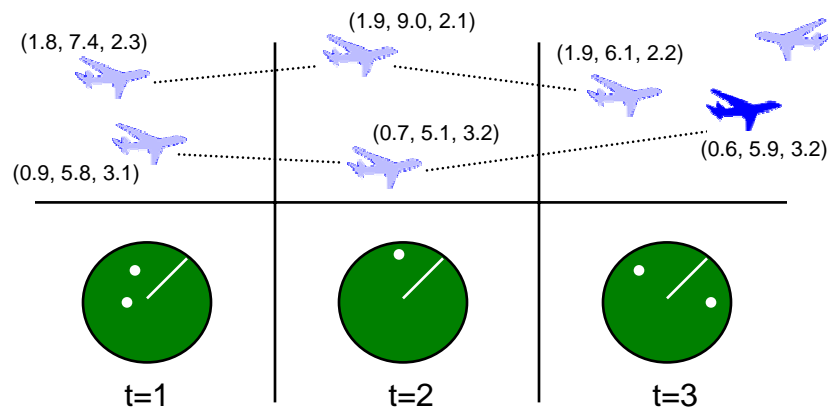# Clustering

# 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*

# Object Recognition and Tracking



(1.9, 9.0, 2.1)

(1.8, 7.4, 2.3)

(1.9, 6.1, 2.2)

(0.7, 5.1, 3.2)

(0.6, 5.9, 3.2)

(0.9, 5.8, 3.1)

t=1          t=2          t=3

# Modeling The Mind …

**Latent brain processes:**



- View picture
- Read sentence
- Decide whether consistent

**fMRI scan:**



$\Sigma$

t=1          · · ·          t=T

---

# The Evolution of Science



**Research circles**

**Research topics**

Phy

Bio

CS

*PNAS* papers

1900          2000

# A Classical Approach

- Clustering as Mixture Modeling



- Then "model selection"

# Partially Observed, Open and Evolving Possible Worlds

- Unbounded # of objects/trajectories
- Changing attributes
- Birth/death, merge/split
- Relational ambiguity
- The parametric paradigm:

**Event model**
$$p\left(\left\{\phi_k^0\right\}\right) \text{ or } p\left(\left\{\phi_k^{1:T}\right\}\right)$$

**motion model**
$$p\left(\left\{\phi_k^{t+1}\right\} \mid \left\{\phi_k^t\right\}\right)$$

*Entity space*

*observation space*

**Sensor model**
$$p\left(\mathbf{x} \mid \left\{\phi_k\right\}\right)$$

- Finite
- Structurally unambiguous

**How to open it up?**

4

# Model Selection vs. Posterior Inference

- Model selection
  - "intelligent" guess: ???
  - cross validation: data-hungry ☹
  - information theoretic:
    - AIC
    - TIC
    - MDL :    $\left.\begin{array}{c}\end{array}\right\}$ $\arg\min KL\left(f(\cdot)\mid g(\cdot\mid\hat{\theta}_{ML},K)\right)$
    
    Parsimony, Ockam's Razor
  - Bayes factor:    need to compute data likelihood

- Posterior inference:
  we want to handle uncertainty of model complexity explicitly
  
  $$p(M\mid D)\propto p(D\mid M)p(M)$$
  
  $$M\equiv\{\theta,K\}$$
  
  - we favor a distribution that does not constrain $M$ in a "closed" space!

---

# Two "Recent" Developments

- First order probabilistic languages (FOPLs)
  - Examples: PRM, BLOG …
  - Lift graphical models to "open" world (#rv, relation, index, lifespan …)
  - Focus on complete, consistent, and operating rules to instantiate possible worlds, and formal language of expressing such rules
  - Operational way of defining distributions over possible worlds, via sampling methods

- Bayesian Nonparametrics
  - Examples: Dirichlet processes, stick-breaking processes …
  - From finite, to infinite mixture, to more complex constructions (hierarchies, spatial/temporal sequences, …)
  - Focus on the laws and behaviors of both the generative formalisms and resulting distributions
  - Often offer explicit expression of distributions, and expose the structure of the distributions --- motivate various approximate schemes

# Clustering



- How to label them ?

- How many clusters ???

# Random Partition of Probability Space



$\{\phi_6, \pi_6\}$

$\{\phi_4, \pi_4\}$

$\{\phi_5, \pi_5\}$

(event $p_{event}$)

centroid := $\phi$

$\{\phi_2, \pi_2\}$

Image ele. $= (x, \theta)$

. . .

# Stick-breaking Process

$$G \sim \mathrm{DP}(\alpha, G_0)$$

$$G = \sum_{k=1}^{\infty} \pi_k \delta(\theta_k)$$

$$\theta_k \sim G_0$$ **Location**

$$\sum_{k=1}^{\infty} \pi_k = 1$$

$$\pi_k = \beta_k \prod_{j=1}^{k-1}(1-\beta_k)$$

$$\beta_k \sim \mathrm{Beta}(1,\alpha)$$ **Mass**

| $\prod_{j=1}^{k-1}(1-\beta_j)$ | $\beta_k$ | $\pi_k$ |
|---|---|---|
| 0 | 0.4 | 0.4 |
| 0.6 | 0.5 | 0.3 |
| 0.3 | 0.8 | 0.24 |

$\theta_5 \quad \theta_2 \; \theta_3 \theta_1 \; \theta_4$ $\qquad G_0$

---

# Chinese Restaurant Process

$\theta_1$     $\theta_2$    .....

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

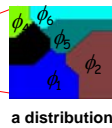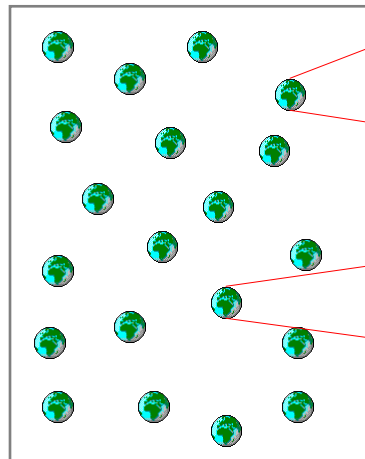| | | |
|---|---|---|
| $1$ | $0$ | $0$ |
| $\dfrac{1}{1+\alpha}$ | $\dfrac{\alpha}{1+\alpha}$ | $0$ |
| $\dfrac{1}{2+\alpha}$ | $\dfrac{1}{2+\alpha}$ | $\dfrac{\alpha}{2+\alpha}$ |
| $\dfrac{1}{3+\alpha}$ | $\dfrac{2}{3+\alpha}$ | $\dfrac{\alpha}{3+\alpha}$ |
| $\dfrac{m_1}{i+\alpha-1}$ | $\dfrac{m_2}{i+\alpha-1}$ .... | $\dfrac{\alpha}{i+\alpha-1}$ |

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

# Dirichlet Process



**a distribution**

**another distribution**

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

$(G(\phi_1), G(\phi_2), …, G(\phi_m)) \sim$ Dirichlet$(\alpha G_0(\phi_1), …., \alpha G0(\phi_m))$

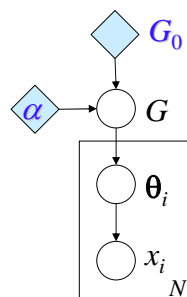where $G_0$ is the base measure and $\alpha$ is the scale parameter

Thus a Dirichlet Process $G$ defines a distribution of distribution

---

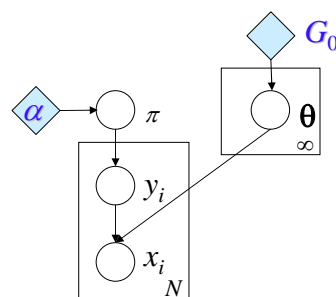# Graphical Model Representations of DP



**The CRP construction**

**The Stick-breaking construction**

# 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)

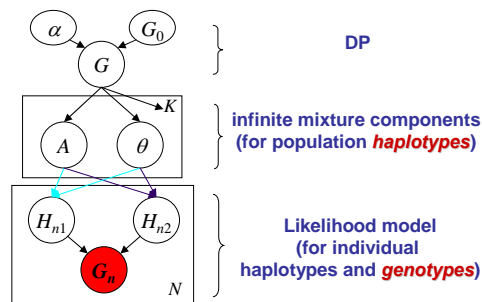- Many other biological/scientific utilities

---

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

- Clustering human populations



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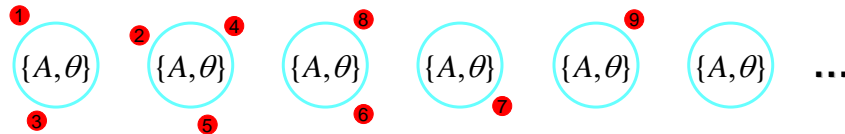  - Gibbs sampling
  - Metropolis Hasting

# 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



$\{A,\theta\}$  $\{A,\theta\}$  $\{A,\theta\}$  $\{A,\theta\}$  $\{A,\theta\}$  $\{A,\theta\}$  ...

- 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)

---

# Inheritance and Observation Models

- Single-locus mutation model

$$A_{C_{i_e}} \rightarrow H_{i_e}$$

$$P_H(h_t \mid a_t, \theta) = \begin{cases} \theta & \text{for } h_t = a_t \\ \dfrac{1-\theta}{|B|-1} & \text{for } h_t \neq a_t \end{cases}$$
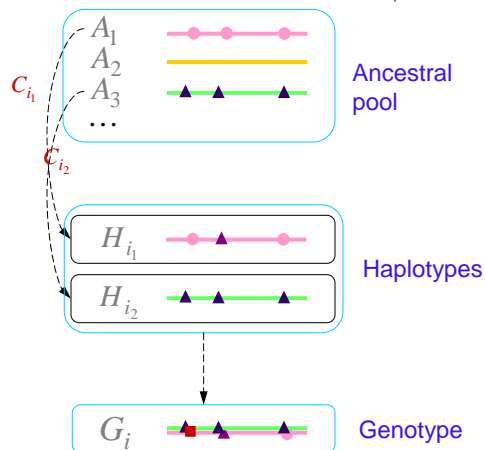
$$\rightarrow h_t = a_t \;\; with \; prob . \, \theta$$

- Noisy observation model

$$H_{i_1}, H_{i_2} \rightarrow G_i$$

$P_G(g \mid h_1, h_2):$
$$g_t = h_{1,t} \oplus h_{2,t} \;\; with \; prob. \; \lambda$$

$A_1$
$A_2$
$A_3$    Ancestral pool
...
$C_{i_1}$
$C_{i_2}$

$H_{i_1}$    Haplotypes
$H_{i_2}$

$G_i$    Genotype

10

# 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 $h_{i_e}$

$$p(c_{i_e} = k \mid \mathbf{c}_{[-i_e]}, \mathbf{h}, \mathbf{a}) \propto p(c_{i_e} = k \mid \mathbf{c}_{[-i_e]})\, p(h_{i_e} \mid a_k, \mathbf{h}_{[-i_e]}, \mathbf{c})$$

  Posterior           Prior     x     Likelihood

  CRP

  :

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

---

# MCMC for Haplotype Inference

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

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

2. Sample $a_k$ from

$$p(a_{k,t} | \mathbf{c}, \mathbf{h}) \propto \prod_{j, i_e | c_{i_e, t}^{(j)} = k} p(h_{i_e, t}^{(j)} | a_{k,t}, l_{k,t}^{(j)})$$
$$= \frac{\Gamma(\alpha_h + l_{k,t})\Gamma(\beta_h + l_{k,t}')}{\Gamma(\alpha_h + \beta_h + m_k)(|B| - 1)^{l_{k,t}'}} R(\alpha_h, \beta_h)$$

3. Sample $h_{ie}^{(j)}$ from

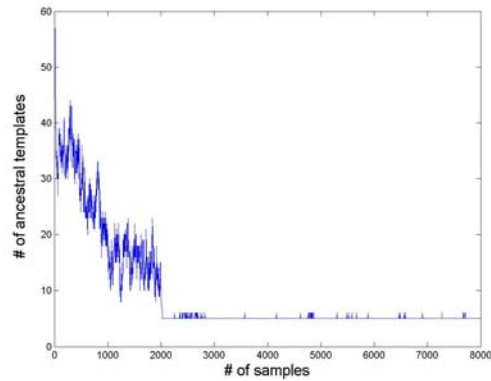$$p(h_{i_e, t}^{(j)} | \mathbf{h}_{[-i_e, t]}^{(j)}, \mathbf{c}, \mathbf{a}, \mathbf{g})$$

- For DP scale parameter $\alpha$: a vague inverse Gamma prior

# Convergence of Ancestral Inference

# DP vs. Finite Mixture via EM

# Variational Inference [Blei & Jordan 2005, Kurihara et al 2007]

- Gibbs sampling solution is not efficient enough to scale up to the large scale problems.
- Truncated stick-breaking approximation can be formulated in the space of explicit, non-exchangeable cluster labels.
- Variational inference can now be applied to such a finite-dimensional distribution

- Variational Inference:
  - For a complicated $P(X_1, X_2, \ldots X_n)$, approximate it with $Q(X)$:

$$Q(\mathbf{X}) = \prod_i Q(\mathbf{X}_{C_i})$$

$$\{Q^*(\mathbf{X}_{C_i})\} = \arg\min KL(Q(\mathbf{X})|P(\mathbf{X}))$$

---

# Approximations to DP

- Truncated stick-breaking representation

$$v_i \sim \mathcal{B}(v_i; 1, \alpha) \qquad i = 1, \ldots, T-1$$
$$v_T = 1$$
$$\pi_i = v_i \prod_{j<i}(1 - v_j) \qquad i = 1, \ldots, T$$
$$\pi_i = 0 \qquad i > T$$

- Finite symmetric Dirichlet approximation

$$\boldsymbol{\pi} \sim \mathcal{D}(\boldsymbol{\pi}; \tfrac{\alpha}{K}, \ldots, \tfrac{\alpha}{K})$$

- The joint distribution can be expressed as:

$$P(X, \mathbf{z}, \mathbf{v}, \boldsymbol{\eta}) = \left[\prod_{n=1}^{N} p(\mathbf{x}_n|\eta_{z_n})\, p(z_n|\boldsymbol{\pi}(\mathbf{v}))\right] \left[\prod_{i=1}^{T} p(\eta_i)\mathcal{B}(v_i; 1, \alpha)\right]$$

- The joint distribution can be expressed as:

$$P(X, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\eta}) = \left[\prod_{n=1}^{N} p(\mathbf{x}_n|\eta_{z_n})\, p(z_n|\boldsymbol{\pi})\right] \left[\prod_{i=1}^{K} p(\eta_i)\right] \mathcal{D}(\boldsymbol{\pi}; \tfrac{\alpha}{K}, \ldots, \tfrac{\alpha}{K})$$

13

# TDP vs. TSB



- TDP is size biased
- cluster labels is NOT interchangeable under TDP but is interchangeable under TSB

---

$$P(X, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\eta}) = \left[\prod_{n=1}^{N} p(\mathbf{x}_n|\eta_{z_n})\, p(z_n|\boldsymbol{\pi})\right]\left[\prod_{i=1}^{K} p(\eta_i)\right] \mathcal{D}(\boldsymbol{\pi}; \tfrac{\alpha}{K}, ..., \tfrac{\alpha}{K})$$

# Marginalization

- In variational Bayesian approximation, we assume a factorized form for the posterior distribution.
- However it is not a good assumption since changes in $\boldsymbol{\pi}$ will have a considerable impact on $\boldsymbol{z}$.

If we can integrate out $\boldsymbol{\pi}$, the joint distribution is given by

$$P(X, \mathbf{z}, \boldsymbol{\eta}) = \left[\prod_{n=1}^{N} p(\mathbf{x}_n|\eta_{z_n})\right] p(\mathbf{z}) \left[\prod_{i=1}^{\infty} p(\eta_i)\right]$$

For the TSB representation:

$$p_{\text{TSB}}(\mathbf{z}) = \prod_{i<T} \frac{\Gamma(1+N_i)\Gamma(\alpha+N_{>i})}{\Gamma(1+\alpha+N_{\geq i})} \ \alpha$$

For the FSD representation:

$$p_{\text{FSD}}(\mathbf{z}) = \frac{\Gamma(\alpha)\prod_{k=1}^{K}\Gamma(N_k+\frac{\alpha}{K})}{\Gamma(N+\alpha)\Gamma(\frac{\alpha}{K})^K}$$

14

# VB inference

- We can then apply the VB inference on the four approximations

$$\{Q^*(\mathbf{X}_{C_i})\} = \arg\min KL(Q(\mathbf{X})|P(\mathbf{X}))$$

The approximated posterior distribution for TSB and FSD are

$$Q_{\mathrm{TSB}}(\mathbf{z}, \boldsymbol{\eta}, \mathbf{v}) = \left[\prod_n^N q(z_n)\right]\left[\prod_{i=1}^T q(\eta_i)q(v_i)\right] \qquad Q_{\mathrm{FSD}}(\mathbf{z}, \boldsymbol{\eta}, \boldsymbol{\pi}) = \left[\prod_n^N q(z_n)\right]\left[\prod_{k=1}^K q(\eta_k)\right]q(\boldsymbol{\pi})$$

Depending on marginalization or not, **v** and $\pi$ may be integrated out.
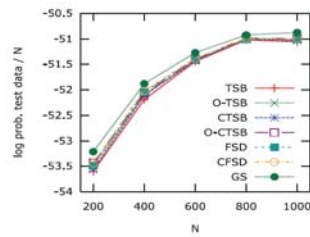
# Experimental results



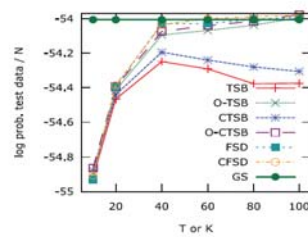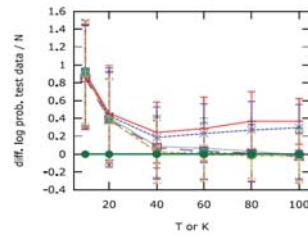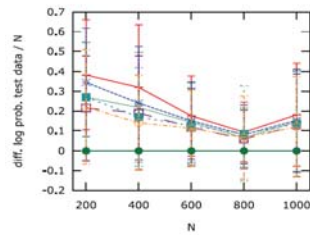Figure 2: Average log probability per data-point for test data as a function of $N$.

Figure 4: Average log probability per data-point for test data as a function of $T$ (for TSB methods) or $K$ (for FSD methods).

# **Summary**

- A non-parametric Bayesian model for Pattern Uncovery

  - Finite mixture model of latent patterns (e.g., image segments, objects)
    - → infinite mixture of propotypes: alternative to model selection
    - → hierarchical infinite mixture
    - → infinite hidden Markov model
    - → temporal infinite mixture model

- Applications in general data-mining …

Eric Xing © Eric Xing @ CMU, 2006-2009 31

16