# Infinite Mixture and Dirichlet Process
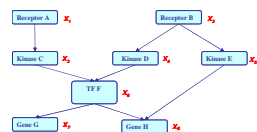
## Probabilistic Graphical Models  (10-708)

**Lecture 20, Nov 28, 2007**
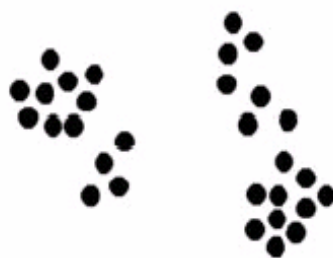
**Eric Xing**

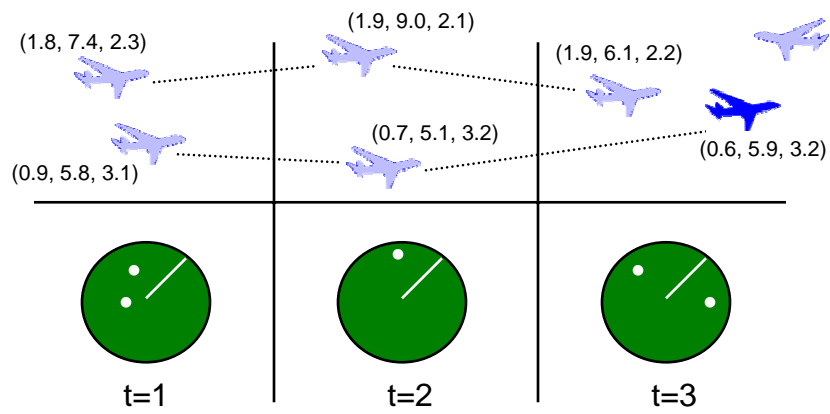**Reading:**

Receptor A $x_1$   Receptor B $x_2$
Kinase C $x_3$   Kinase D $x_4$   Kinase E $x_5$
TF F $x_6$
Gene G $x_7$   Gene H $x_8$

1

---

# Clustering



$K?$

$P(K)$

Eric Xing                                        2

# 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)

(0.9, 5.8, 3.1)

(0.6, 5.9, 3.2)

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

2

# The Evolution of Science



Research circles

Research topics

Phy

Bio

CS

*PNAS* papers

1900 — 2000 ?

---
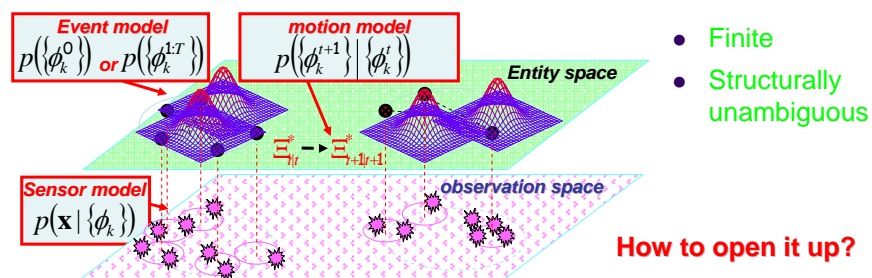
# Partially Observed, Open and Evolving Possible Worlds

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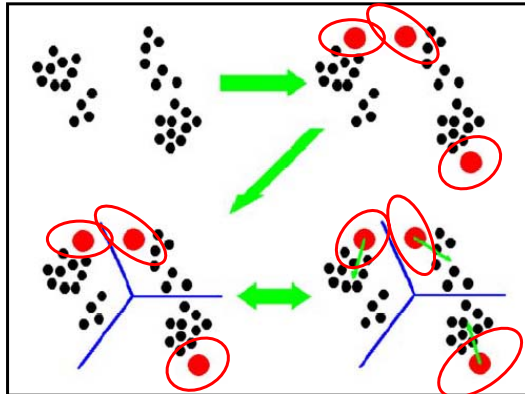- Relational ambiguity

- 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**
$$p\left(\left\{\phi_k^{t+1}\right\}\middle|\left\{\phi_k^t\right\}\right)$$

*Entity space*

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

*observation space*

- Finite
- Structurally unambiguous

**How to open it up?**

3

# A Classical Approach

- Clustering as Mixture Modeling



- Then "model selection"

# Model Selection vs. Posterior Inference

- Model selection
    - "intelligent" guess: ???
    - cross validation: data-hungry ☹
    - information theoretic:
        - AIC
        - TIC  $\Big\}$  $\arg\min KL\big(f(\cdot)\mid g(\cdot\mid\hat{\theta}_{ML},K)\big)$
        - MDL :                    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
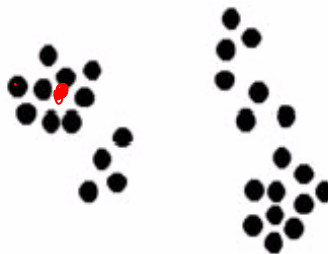
- First order probabilistic languages (FOPLs)
    - Examples: PRM, BLOG …
    - Lift graphical models to "open" world (#rv, relation, index, lifespan …)
    - Focus on complete, consistent, and operating rules to instantiate possible worlds, and formal language of expressing such rules
    - Operational 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

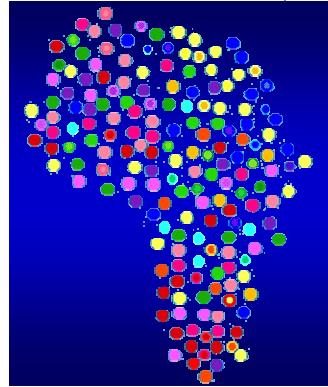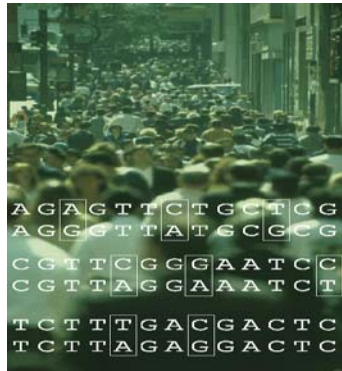Eric Xing 9

---

# Clustering



- How to label them ?

- How many clusters ???

Eric Xing 10

# Genetic Demography





- Are there genetic prototypes among them ?
- What are they ?
- How many ? (how many ancestors do we have ?)

# Genetic Polymorphisms



The ABO Blood System

| Blood Type (genotype) | Type A (AA, AO) | Type B (BB, BO) | Type AB (AB) | Type O (OO) |
|---|---|---|---|---|
| Red Blood Cell Surface Proteins (phenotype) | A agglutinogens only | B agglutinogens only | A and B agglutinogens | No agglutinogens |
| Plasma Antibodies (phenotype) | b agglutinin only | a agglutinin only | No agglutinin | a and b agglutinin |

# Biological Terms

- Genetic polymorphism: a difference in DNA sequence among individuals, groups, or populations

- **Single Nucleotide Polymorphism (SNP):** DNA sequence variation occurring when a single nucleotide - A, T, C, or G - differs between members of the species

  – Each variant is called an "allele"
  – Almost always bi-allelic
  – Account for most of the genetic diversity among different (normal) individuals, e.g. drug response, disease susceptibility
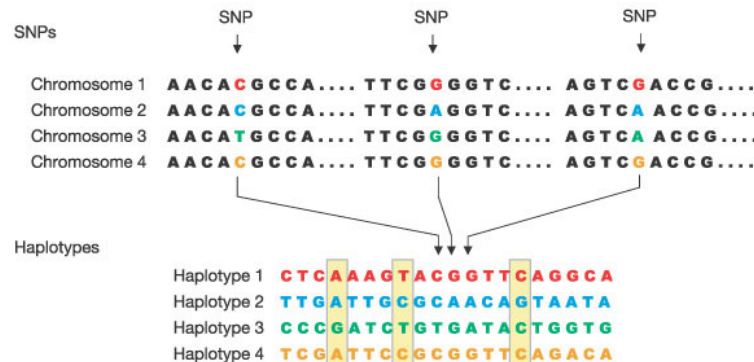


Eric Xing 13

# From SNPs to Haplotypes

- Alleles of adjacent SNPs on a chromosome form **haplotypes**



- Powerful in the study of disease association or genetic evolution

Eric Xing 14

7

# Haplotype and Genotype

- A collection of alleles derived from the same chromosome

**Genotypes**

| | |
|---|---|
| 2 | 13 |
| 1 | 6 |
| 9 | 15 |
| 4 | 17 |
| 1 | 9 |
| 2 | 6 |
| 9 | 17 |
| 2 | 12 |
| 7 | 12 |
| 6 | 14 |
| 1 | 7 |
| 18 | 18 |
| 1 | 4 |
| 10 | 10 |

**Haplotype Re-construction** →

**Haplotypes**

| | |
|---|---|
| 2 | 13 |
| 6 | 1 |
| 9 | 15 |
| 17 | 4 |
| 1 | 9 |
| 6 | 2 |
| 9 | 17 |
| 2 | 12 |
| 12 | 7 |
| 14 | 6 |
| 7 | 1 |
| 18 | 18 |
| 1 | 4 |
| 10 | 10 |

**Chromosome phase is unknown**          **Chromosome phase is known**

---

# Ancestral Inference

$p(A, \theta | G)$

$A_k$   $\theta_k$   ?

$H_{n1}$   $H_{n2}$

$G_n$

$N$

**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

# A Finite (Mixture of ) Allele Model

- The probability of a genotype $g$:

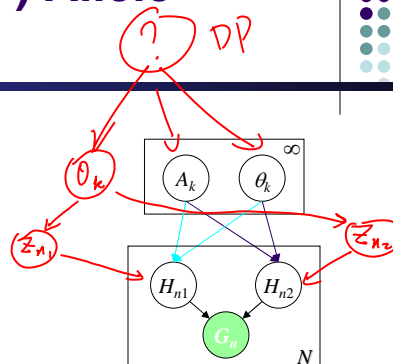$$p(g) = \sum_{h_1, h_2 \in \mathcal{H}} p(h_1, h_2)\, p(g \mid h_1, h_2)$$

| Population haplotype pool | Haplotype model | Genotyping model |

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

---

# A Infinite (Mixture of ) Allele Model



- How?
  - Via a nonparametric hierarchical Bayesian formalism !

9

# Stick-breaking Process

$\beta_k \sim P(\cdot | v_i)$

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

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

$$\theta_k \sim G_0$$

Location

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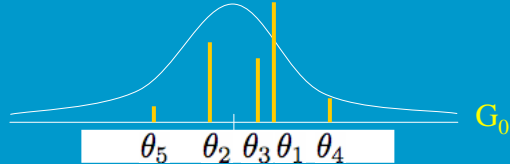
$$\pi_k = \beta_k \prod_{j=1}^{k-1}(1 - \beta_k)$$

Mass

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

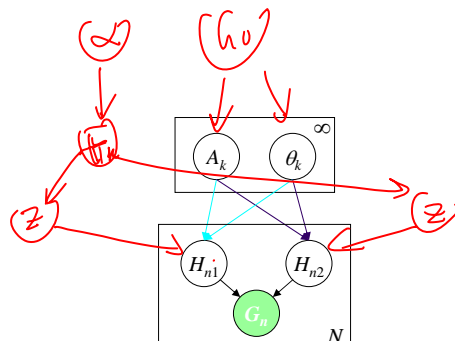| $\prod_{j=1}^{k-1}(1 - \beta_j)$ | $\beta_k$ | $\pi_k$ |
|---|---|---|
| 0 | 0.4 | 0.4 |
| 0.6 | 0.5 | 0.3 |
| 0.3 | 0.8 | 0.24 |

$\theta_5 \quad \theta_2 \quad \theta_3 \theta_1 \quad \theta_4$

$G_0$

Eric Xing
19

---

# Graphical Model



Eric Xing
20

10

# Chinese Restaurant Process



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

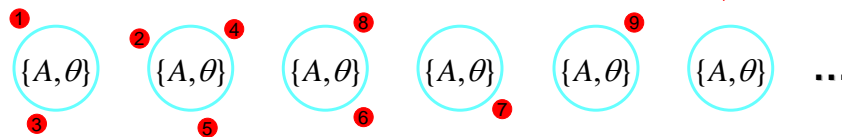| | | |
|---|---|---|
| $1$ | $0$ | $0$ |
| $\dfrac{1}{1+\alpha}$ | $\dfrac{\alpha}{1+\alpha}$ | $0$ |
| $\dfrac{1}{2+\alpha}$ | $\dfrac{1}{2+\alpha}$ | $\dfrac{\alpha}{2+\alpha}$ |
| $\dfrac{1}{3+\alpha}$ | $\dfrac{2}{3+\alpha}$ | $\dfrac{\alpha}{3+\alpha}$ |
| $\dfrac{m_1}{i+\alpha-1}$ | $\dfrac{m_2}{i+\alpha-1}$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

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

Eric Xing 21

---

# The DP Mixture of Ancestral Haplotypes

- The customers around a table 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)
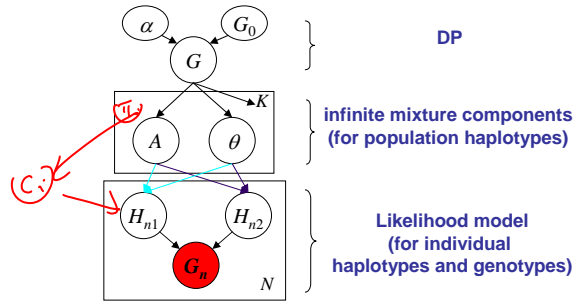
Eric Xing 22

11

# DP-haplotyper



DP

infinite mixture components
(for population haplotypes)

Likelihood model
(for individual
haplotypes and genotypes)

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

---

# Model components

- Choice of base measure:

$$G_0 \sim \mathrm{Unif}(a) \cdot \prod_j \mathrm{Beta}(\theta_j)$$

- Nucleotide-substitution model:

$$p(h_i \mid \{a,\theta\}_k) = \prod_j p(h_{i,j} \mid a_{k,j}, \theta_{k,j})$$

$$\text{where} \quad p(h_{i,j} \mid a_{k,j}, \theta_{k,j}) = \begin{cases} \theta_{k,j} & \text{if} \quad h_{i,j} = a_{k,j} \\ 1 - \theta_{k,j} & \text{if} \quad h_{i,j} = a_{k,j} \end{cases}$$

- Noisy genotyping model:

$$p(g_i \mid h_{i_1}, h_{i_2}) = \prod_j p(g_{i,j} \mid h_{i_1,j}, h_{i_2,j})$$

$$\text{where} \quad p(g_{i,j} \mid h_{i_1,j}, h_{i_2,j}) = \begin{cases} \gamma & \text{if} \quad h_{i_1,j} \oplus h_{i_2,j} = g_{i,j} \\ \dfrac{1-\gamma}{2} & \text{if} \quad h_{i_1,j} \oplus h_{i_2,j} \neq g_{i,j} \end{cases}$$
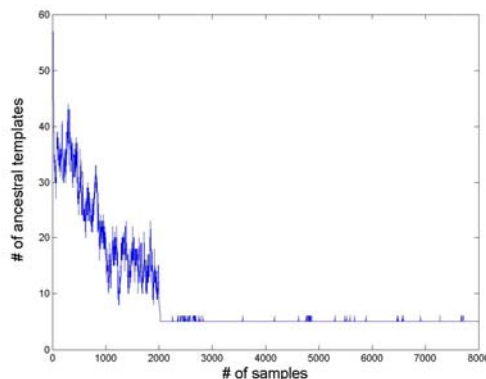
# Gibbs sampling

Starting from some initial haplotype reconstruction $H^{(0)}$, pick a first table with an arbitrary $a_1^{(0)}$, and form initial population-hap pool $\mathbf{A}^{(0)} = \{a_1^{(0)}\}$:

i) Choose an individual $i$ and one of his/her two haplotypes $t$, uniformly and at random, from all ambiguous individuals;

$$/ P_{CRP}(c_i) P(H_i | A_{c_i})$$

ii) Sample $c_{i_t}^{(t+1)}$ from $p(c_{i_t}^{(t+1)} | c_{-i_t}^{(t)}, H^{(t)}, \mathbf{A}^{(t)})$, update $c^{(t+1)}$;

iii) Sample $a_k^{(t+1)}$, where $k = c_{i_t}^{(t+1)}$, from $p(a_k^{(t+1)} | \forall h_{-i'_{t'}}^{(t)} \text{ s.t. } c_{i'_{t'}}^{(t+1)} = k)$; update $\mathbf{A}^{(t+1)}$;

iii) Sample $h_{i_t}^{(t+1)}$ from $p(h_{i_t}^{(t+1)} | c_{i_t}^{(t+1)}, H_{-i_t}^{(t)}, \mathbf{A}^{(t+1)})$, update $H^{(t+1)}$.

---

# Convergence of Ancestral Inference

# Haplotyping Error

## The Gabriel data