

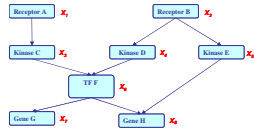
Infinite Mixture and Dirichlet Process

Probabilistic Graphical Models (10-708)

Lecture 20, Nov 28, 2007

Eric Xing

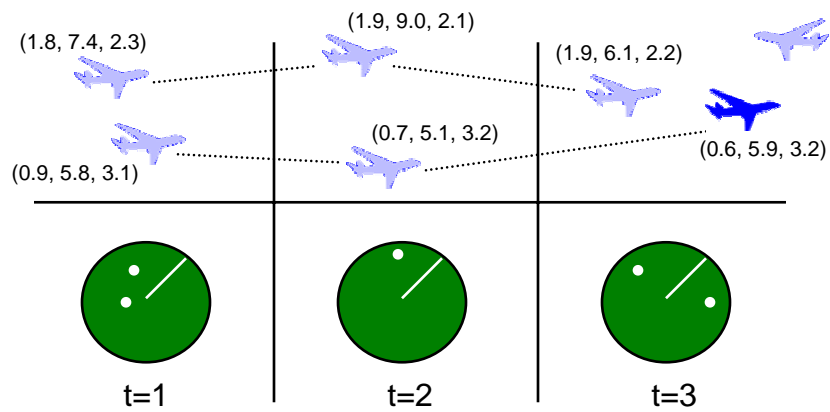
Reading:



Clustering



Object Recognition and Tracking



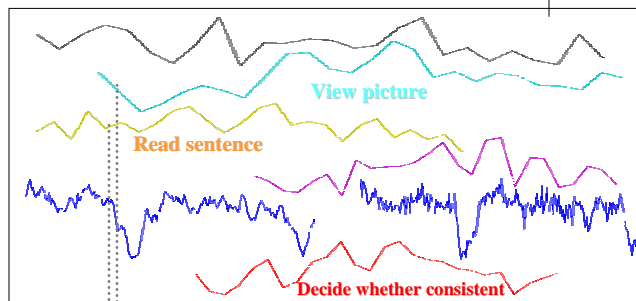
Eric Xing

3

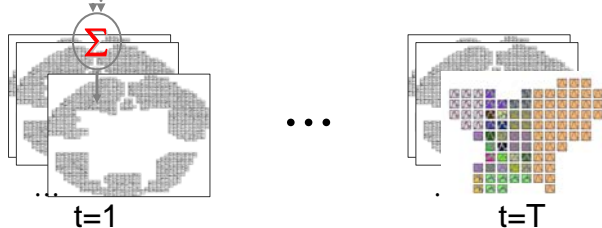
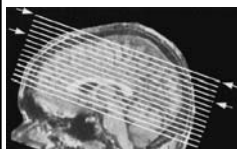
Modeling The Mind ...



Latent brain processes:



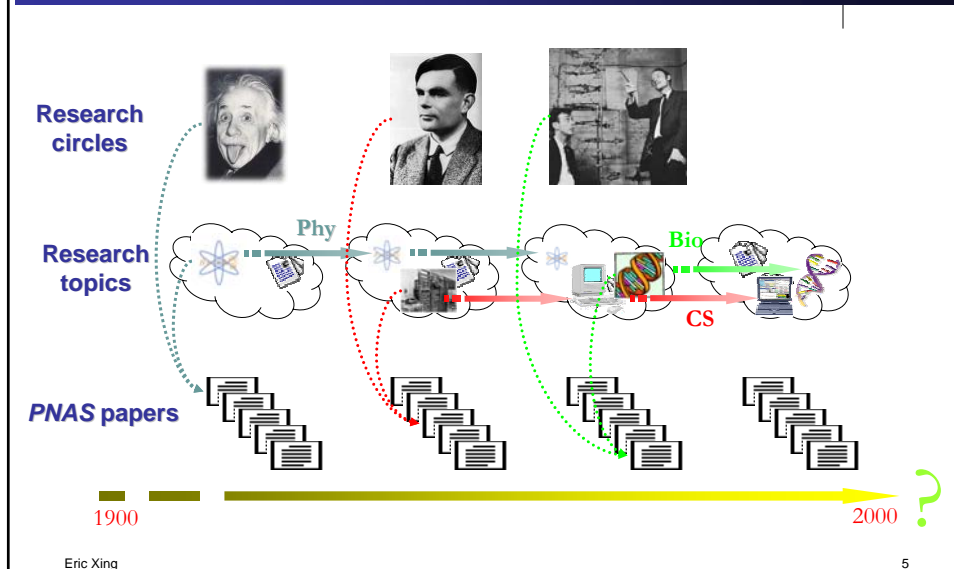
fMRI scan:



Eric Xing

4

The Evolution of Science

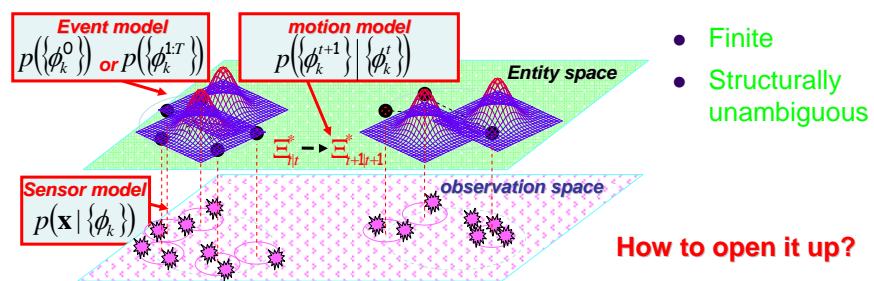


Eric Xing

5

Partially Observed, Open and Evolving Possible Worlds

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

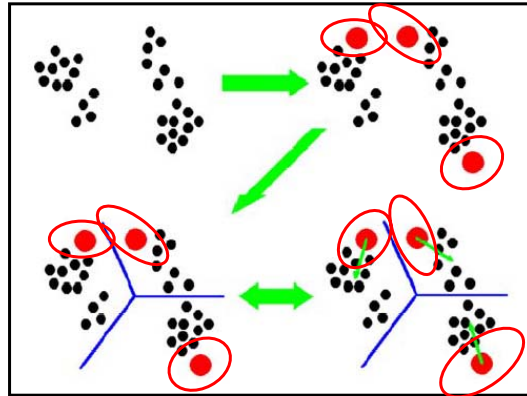


Eric Xing

6

A Classical Approach

- Clustering as Mixture Modeling



- Then "model selection"

Eric Xing

7

Model Selection vs. Posterior Inference

- Model selection
 - "intelligent" guess: ???
 - cross validation: data-hungry ☹
 - information theoretic:
 - AIC
 - TIC
 - MDL : Parsimony, Ockam's Razor
 - Bayes factor: need to compute data likelihood

- Posterior inference:
 - we want to handle uncertainty of model complexity explicitly

$$p(M|D) \propto p(D|M)p(M)$$

$$M \equiv \{\theta, K\}$$

- we favor a distribution that does not constrain M in a "closed" space!

Eric Xing

8

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

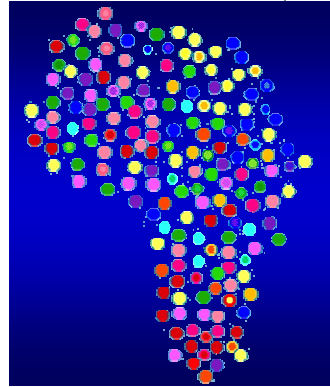
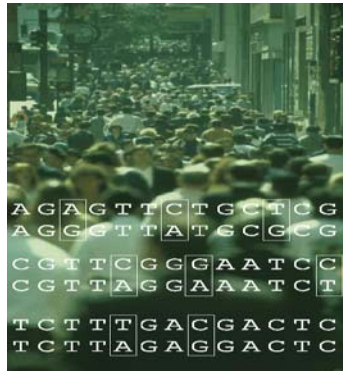


- How to label them ?
- How many clusters ???

Eric Xing

10

Genetic Demography



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

Eric Xing

11

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	NONE No agglutinin	a and b agglutinin

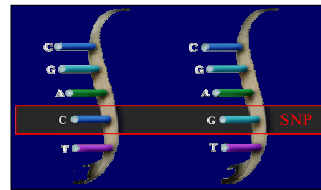
Eric Xing

12

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



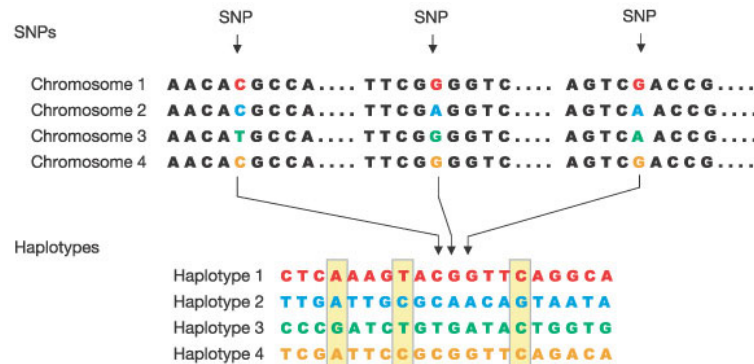
Eric Xing

13

From SNPs to Haplotypes



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



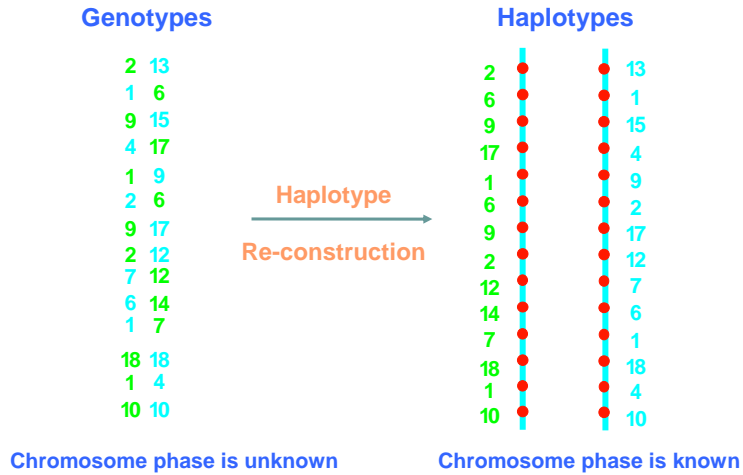
- Powerful in the study of **disease association** or **genetic evolution**

Eric Xing

14

Haplotype and Genotype

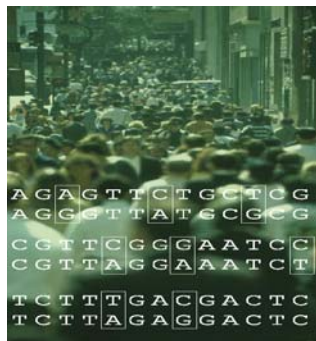
- A collection of alleles derived from the same chromosome



Eric Xing

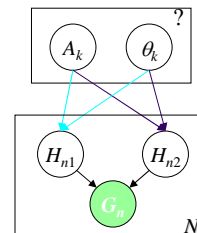
15

Ancestral Inference



```

AGAGTTCTGCTCG
AGGCTTATGCGCG
CGTTCGGGAATCC
CGTTAGGAATCT
TCTTTGACGACTC
TCTTAGAGGACTC
    
```



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

16

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:

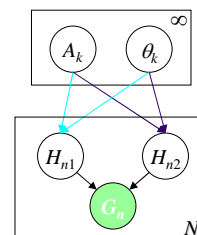
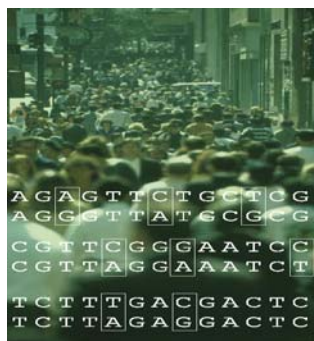
- $|\mathcal{H}| = K \ll 2^J$ fixed-sized population haplotype pool
- $p(h_1, h_2) = p(h_1)p(h_2) = f_1 f_2$ Hardy-Weinberg equilibrium

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

Eric Xing

17

A Infinite (Mixture of) Allele Model



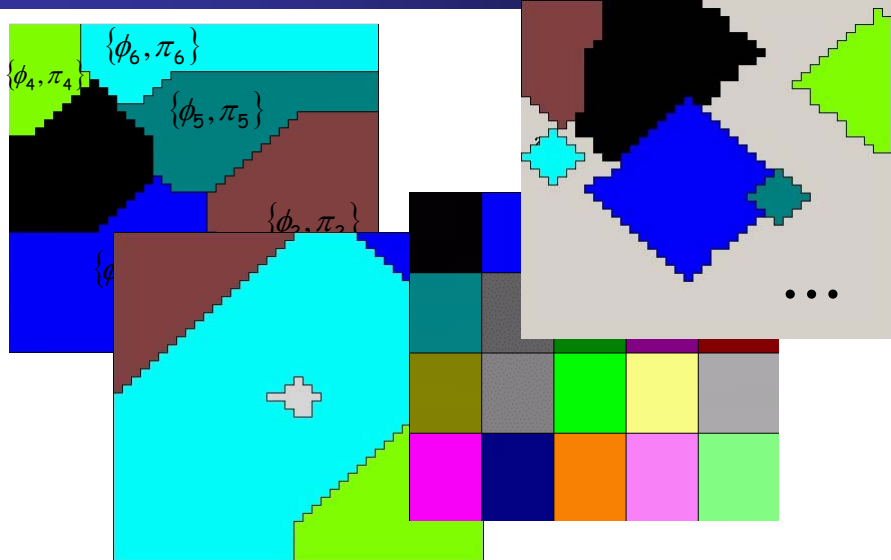
- How?

- Via a nonparametric hierarchical Bayesian formalism !

Eric Xing

18

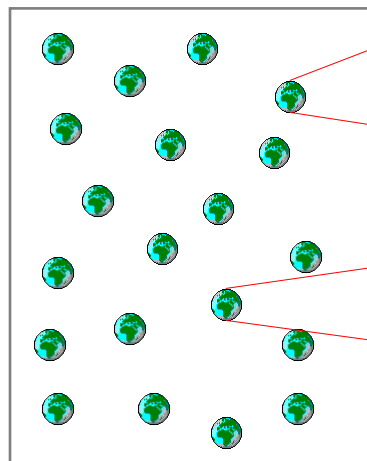
Random Partition of Probability Space



Eric Xing

19

Dirichlet Process



- 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**

Thus a Dirichlet Process G defines a distribution of distribution

Eric Xing

20

Stick-breaking Process



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

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

$\theta_k \sim G_0$

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

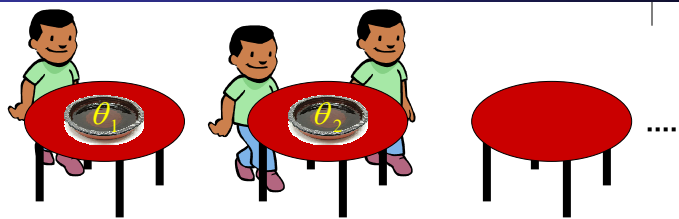
$\pi_k = \beta_k \prod_{j=1}^{k-1} (1 - \beta_j)$

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

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

Eric Xing 21

Chinese Restaurant Process



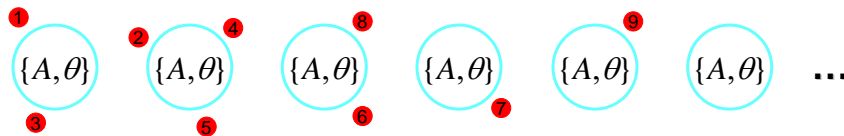
$P(c_i = k \mathbf{c}_{-i}) =$	$\frac{1}{1+\alpha}$	$\frac{0}{1+\alpha}$	$\frac{0}{1+\alpha}$
	$\frac{1}{2+\alpha}$	$\frac{1}{2+\alpha}$	$\frac{\alpha}{2+\alpha}$
	$\frac{1}{3+\alpha}$	$\frac{2}{3+\alpha}$	$\frac{\alpha}{3+\alpha}$
	$\frac{m_1}{i+\alpha-1}$	$\frac{m_2}{i+\alpha-1}$	$\frac{\alpha}{i+\alpha-1}$

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

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

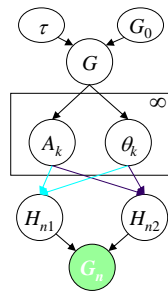


- With $p(h|\{A, \theta\})$ and $p(g|h_1, h_2)$, the CRP yields a posterior distribution on the number of population haplotypes (and on the haplotype configurations and the nucleotide substitution frequencies)

Eric Xing

23

A Hierarchical Bayesian Infinite Allele model



- Assume an individual haplotype h is stochastically derived from a population haplotype a_k with nucleotide-substitution frequency θ_k :

$$h \sim p(h|\{a, \theta\}_k).$$

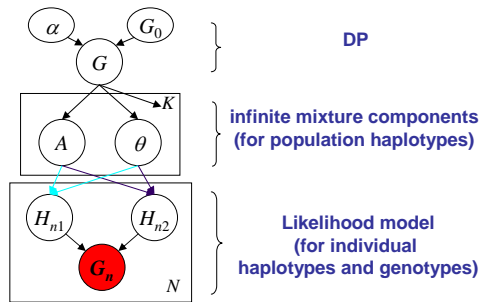
- Not knowing the correspondences between individual and population haplotypes, each individual haplotype is a mixture of population haplotypes.

- The number and identity of the population haplotypes are unknown
 - use a Dirichlet Process to construct a prior distribution G on $\mathcal{H} \times \mathcal{R}^l$.

Eric Xing

24

DP-haplotype



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

Model components



- Choice of base measure:

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

- Nucleotide-substitution model:

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

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

- Noisy genotyping model:

$$p(g_i | h_1, h_2) = \prod_j p(g_{i,j} | h_{1,j}, h_{2,j})$$

$$\text{where } p(g_{i,j} | h_{1,j}, h_{2,j}) = \begin{cases} \gamma & \text{if } h_{1,j} \oplus h_{2,j} = g_{i,j} \\ \frac{1-\gamma}{2} & \text{if } h_{1,j} \oplus h_{2,j} \neq g_{i,j} \end{cases}$$

Gibbs sampling



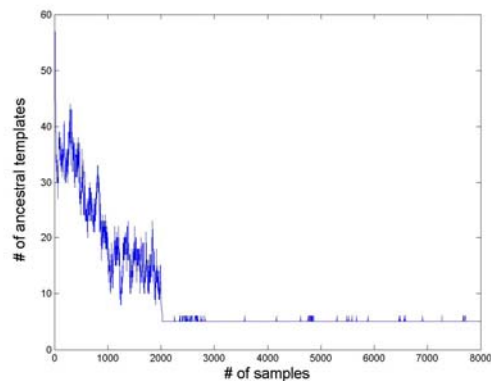
Starting from some initial haplotype reconstruction $H^{(0)}$, pick a first table with an arbitrary $a_j^{(0)}$, and form initial population-hap pool $\mathbf{A}^{(0)} = \{a_j^{(0)}\}$:

- i) Choose an individual i and one of his/her two haplotypes t , uniformly and at random, from all ambiguous individuals;
- ii) Sample $c_i^{(t+1)}$ from $p(c_i^{(t+1)} | c_{-i}^{(t)}, H^{(t)}, \mathbf{A}^{(t)})$, update $c^{(t+1)}$;
- iii) Sample $a_k^{(t+1)}$, where $k = c_i^{(t+1)}$, from $p(a_k^{(t+1)} | \forall h_{-i}^{(t)} \text{ s.t. } c_{i'}^{(t+1)} = k)$; update $\mathbf{A}^{(t+1)}$;
- iii) Sample $h_i^{(t+1)}$ from $p(h_i^{(t+1)} | c_i^{(t+1)}, H_{-i}^{(t)}, \mathbf{A}^{(t+1)})$, update $H^{(t+1)}$.

Eric Xing

27

Convergence of Ancestral Inference



Eric Xing

28

Haplotyping Error



The Gabriel data

