
MCMC for the analysis of genetic data on pedigrees: Tutorial Session 1

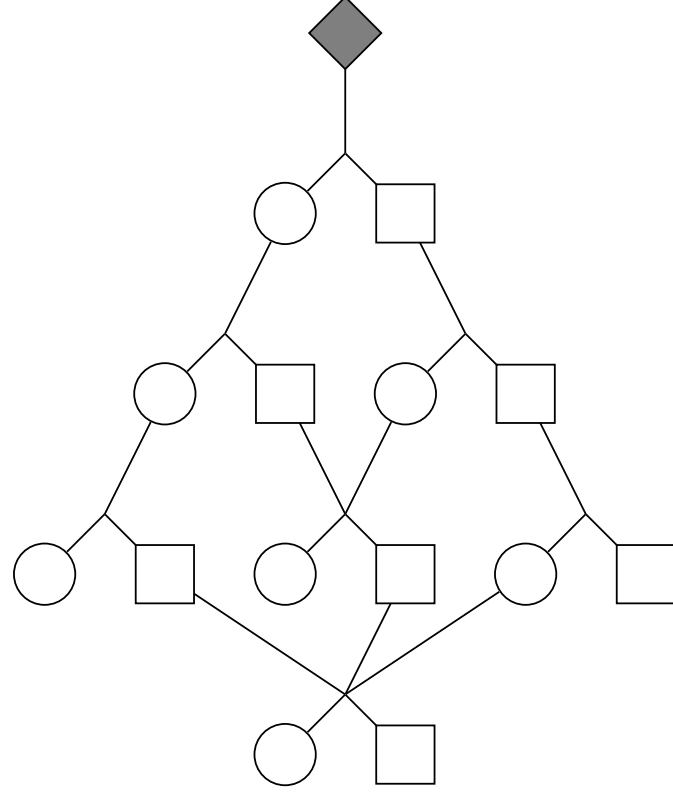
Elizabeth Thompson
University of Washington

- Inheritance and the descent of genes in pedigrees.
- Computation of probabilities on pedigrees.
- MCMC sampling of inheritance patterns in pedigrees.

A PEDIGREE

Individuals have unique identifiers.

To specify pedigree: specify parent identifiers of each individual. **Founders** have unspecified parents: others are **non-founders**.
 Notation: male, female, affected or other data.



GENE IDENTITY BY DESCENT (*ibd*)

- Human individuals are diploid: every cell contains 2 copies of the human genome: one maternal, the other paternal.

• MENDEL'S FIRST LAW (1866)

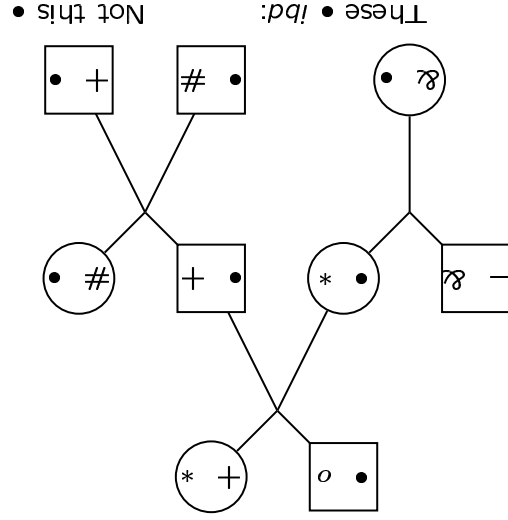
Each parent individual segregates a randomly chosen one of its two genes to each offspring, independently to each offspring. (All meioses are independent.)

- Genes that are copies of the same gene in a recent* common ancestor are said to be identical by descent (*ibd*).
* : *ibd* is defined relative to given pedigree or time point.

- **Simple model:** *ibd* genes are of the same allelic type, non-*ibd* genes are of independent types.

GENE IDENTITY BY DESCENT (*ibd*)

Given I have blood type O, the probability my cousins have blood type O is increased, because with some probability they share genes *ibd* with me at this locus.



- RELATIVES ARE SIMILAR because they have *ibd* genes.
- MENDELIAN GENETICS APPLIES TO MARKERS.

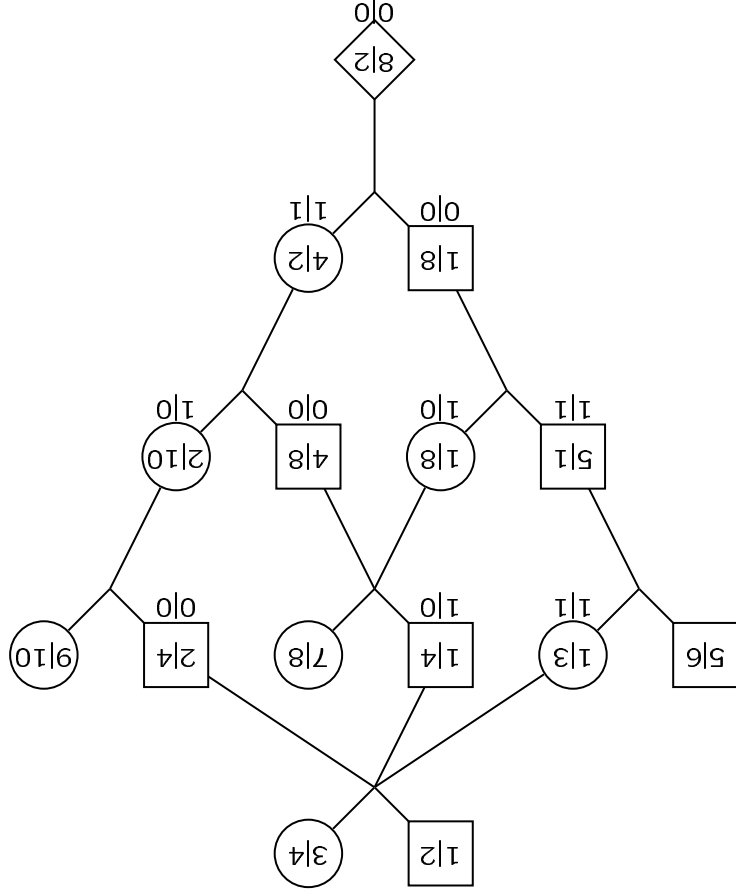
THE INHERITANCE OF GENOME

Label the two haploid genomes of every founder: **Founder genome labels (FGL)**.

Inheritance of FGL:

$$S_{i,j} = 0 \text{ or } 1$$

as in meiosis i at locus j the maternal or paternal gene (respectively) of the parent is transmitted to the offspring.



INHERITANCE AT A SET OF LOCI

... is fully specified by *meiosis indicators*:

For loci $j, j = 1, \dots, T,$

$S_{i,j} = 0$ if gene at meiosis i locus j is parent's maternal gene.
 $S_{i,j} = 1$ if gene at meiosis i locus j is parent's paternal gene.

Notation:

$$S_{i,j}^\bullet = \{S_{i,j';i} = 1, \dots, T\}, \quad j = 1, \dots, m$$

$$S_{i,j}^\bullet = \{S_{i,j';i} = 1, \dots, m\}, \quad j = 1, \dots, T$$

where m is the number of meioses in the pedigree, and T the number of loci along the chromosome.

CHROMOSOMES AND MEIOSIS

Chromosomes duplicate and align and exchange material.

Offspring chromosome consists

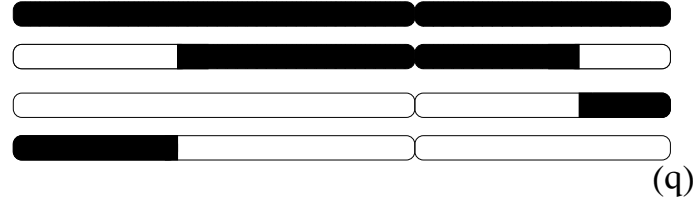
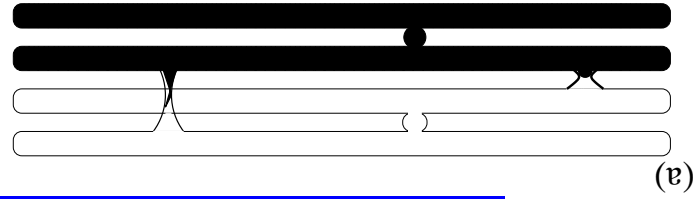
of segments of two parental

chromosomes.

Between two points: recombi-

nation \exists if DNA is from the 2

different parental chromosomes.



$$p(d) = \Pr(\text{recomb.} | d) \nearrow d$$

$$p(0) = 0, p(\infty) = \frac{1}{2}$$

A MODEL FOR LATENT INHERITANCE $S = \{S_{i,j}\}$

- Meioses i are independent: $S_{i,\bullet}$ are independent, a priori.

- Mendel's First Law: $\Pr(S_{i,j} = 0) = \Pr(S_{i,j} = 1) = 1/2$

- Recombination: $\Pr(S_{i,j-1} \neq S_{i,j}) = p_{j-1}$ **same for all i ? – no**

$$\Pr(S_{i,j} | S_{i,j-1}) = p_{j-1}^{R_{j-1}} (1 - p_{j-1})^{m - R_{j-1}}$$

where $R_{j-1} = (\#i : S_{i,j} \neq S_{i,j-1})$

- No genetic interference: $\Pr(S_{i,j} | S_{i,j-1}, S_{i,j-2}, \dots, S_{i,j-1}) = \Pr(S_{i,j} | S_{i,j-1})$

$$\Pr(S) = \prod_{i=1}^2 \Pr(S_{i,\bullet} | S_{i,j-1})$$

A MODEL FOR THE ALLELIC TYPES OF FGL

Allelic types of the FGL are nuisance variables (in most contexts) which we need to marginalize over.

Model: 1. Loci j are independent – good model for $p > 0.005$
 2. At locus j , each FGL g has type k independently with prob $q_{j,k}$.

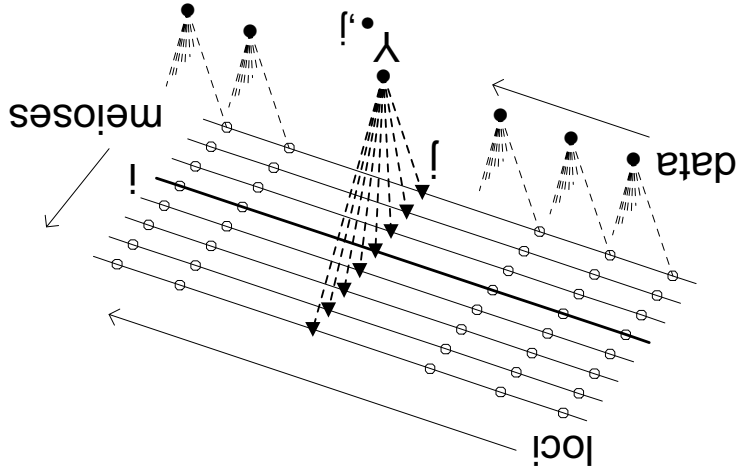
$$\Pr(A_j) = \prod_g q_j(g) = \prod_g q_{j,k}^{n_j(k)}$$

where $n_j(k)$ is number of FGL g with type k at locus j .
 (For (2), better models are possible.)

STRUCTURE OF A GENETIC MODEL

- Population model: parameters \mathbf{q} , provide probabilities for latent A – allelic types of FGL at each j
- Inheritance model: parameters p , provide probabilities for latent S – inheritance of FGL at j , jointly over j .
- Individual genotypes G is deterministic function of (S, A)
- penetrance model, parameters β relates G (and perhaps observable covariates) to observable data \mathbf{Y}
- $\xi = \mathbf{q}, \theta, \beta$.

STATISTICAL VIEW OF A PEDIGREE



Pedigree structure is implicit in the labelling of the meioses.

$Y_{i,j}$: trait or marker data at locus j — trait specific to locus j — can be computed — next we talk about this.

PEELING COMPUTATIONS: THREE CASES

(1) Peeling along a chromosome (HMM): $O(4^m)$

$$\Pr(Y^{(j)}, S^{(j)}) = \sum_{S^{(j-1)}} \Pr(Y^{(j-1)}, S^{(j-1)}) \Pr(S^{(j)} | S^{(j-1)}) \Pr(Y^{(j)} | S^{(j)})$$

Requires sequential summation along the chromosome.
 (2) Pedigree peeling at single locus: $O(K^{3L})$, K large.

$$\Pr(Y^{(j)}) = \sum_{S^{(j)}, A_j} \Pr(Y^{(j)} | S^{(j)}, A_j) \Pr(S^{(j)}) \Pr(A_j)$$

Involves sequential summation over the pedigree structure.
 (3) For both we need

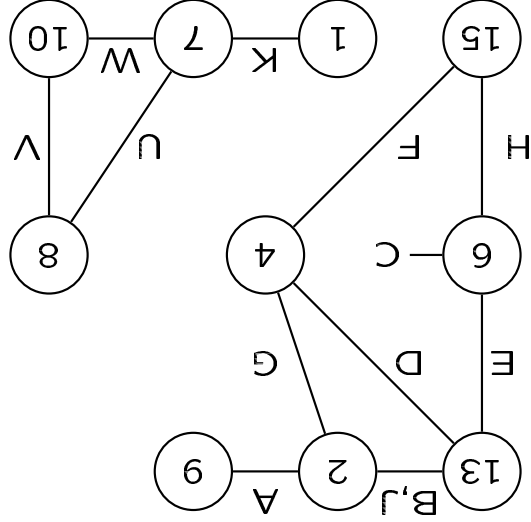
$$\Pr(Y^{(j)} | S^{(j)}) = \sum_{A_j} \Pr(Y^{(j)} | S^{(j)}, A_j) \Pr(A_j)$$

THE FGL-GRAPH STRUCTURE: FOR $\Pr(Y_{\bullet,j} | S_{\bullet,j})$

Only genes in "observed" matter— e.g. FGL 3, 11 ... are not present.
 Only "connected" genes are dependent— lines connect the FGL of observed individuals.

$$\Pr(Y_{\bullet,j} | S_{\bullet,j}) = \sum_{A_j} \Pr(Y_{\bullet,j} | G(S_{\bullet,j}, A_j)) \Pr(A_j) = \sum_{A_j} \Pr(\Pi_n \Pr(Y_{n,j} | G^n(S_{\bullet,j}, A_j))) (\Pi_{g_j} q_j(g))$$

Nodes are FGL
 Edges are observed individuals.



PEELING A COMPONENT OF THE FGL-GRAPH

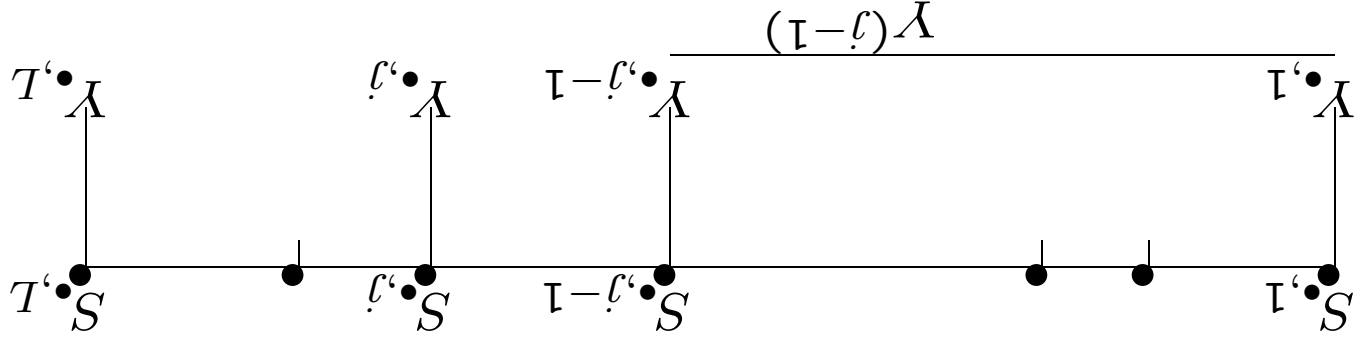
$$\begin{aligned}
 \Pr(Y_{\bullet,j}|S_{\bullet,j}) &= \sum_{\mathbf{g}} \left(\prod_{n=1}^n \Pr(Y_n|g_{n,1}, g_{n,2}) \right) \left(q(g_2)q(g_6)q(g_4) \right. \\
 &\quad \left. \sum_{g_{15}}^{g_{15}} \Pr(Y_H|g_{15}, g_6) \right) \left(\sum_{g_{13}}^{g_{13}} q(g_{13}) \Pr(Y_E|g_6, g_{13}) \Pr(Y_D|g_4, g_{13}) \right) \\
 &\quad \left(\sum_{g_4}^{g_4} q(g_4) \Pr(Y_F|g_{15}, g_4) \right) \left(\sum_{g_2}^{g_2} q(g_2) \Pr(Y_B|g_2, g_{13}) \right) \\
 &\quad \left. \Pr(Y_J|g_2, g_{13}) \Pr(Y_G|g_2, g_4) \right) \left(\sum_{g_9}^{g_9} q(g_9) \Pr(Y_A|g_2, g_9) \right)
 \end{aligned}$$

BAUM ALGORITHM FOR HMM: Lander-Green

For data observations $\mathbf{Y} = (Y_{\bullet,j}, j = 1, \dots, T)$, we want to compute $\Pr(\mathbf{Y})$.

$$\Pr(\mathbf{Y}) = \sum_{\mathbf{S}} \Pr(\mathbf{S}, \mathbf{Y}) = \sum_{\mathbf{S}} \Pr(\mathbf{Y} | \mathbf{S}) \Pr(\mathbf{S})$$

$$= \sum_{\mathbf{S}} \left(\Pr(S_{\bullet,1}) \prod_{j=2}^T \Pr(S_{\bullet,j} | S_{\bullet,j-1}) \prod_{j=1}^T \Pr(Y_{\bullet,j} | S_{\bullet,j}) \right)$$



BAUM ALGORITHM DETAILS

Forwards Baum: Define $R_*^1(s) = \Pr(S_{\bullet,1} = s)$ and $Y^{(j)} = (Y_{\bullet,1}, \dots, Y_{\bullet,j})$, the data up to locus j , so

$$R_*^j(s) = \Pr(Y_{\bullet,k}, k = 1, \dots, j-1, S_{\bullet,j} = s) \\ = \Pr(Y^{(j-1)}, S_{\bullet,j} = s)$$

Now

$$R_*^{j+1}(s) = \sum_{s_*^*} \Pr(S_{\bullet,j+1} = s \mid S_{\bullet,j} = s_*)$$

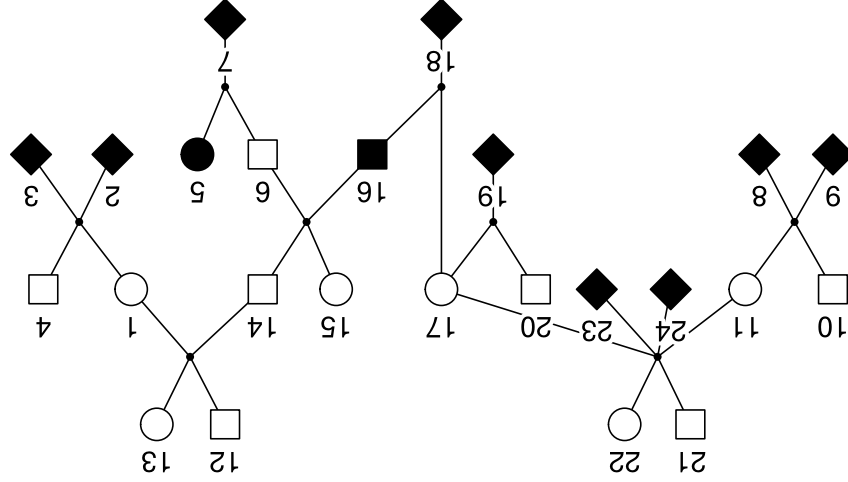
$$\Pr(Y_{\bullet,j} \mid S_{\bullet,j} = s_*) R_*^j(s_*) \Big]$$

for $j = 1, 2, \dots, L-1$, with

$$L = \Pr(Y) = \sum_{s_*^*} \Pr(Y_{\bullet,T} \mid S_{\bullet,T} = s_*) R_*^T(s_*)$$

$S_{\bullet,j}$ can take 2^m values, where m is number of meioses. Computation is limited to small pedigrees.

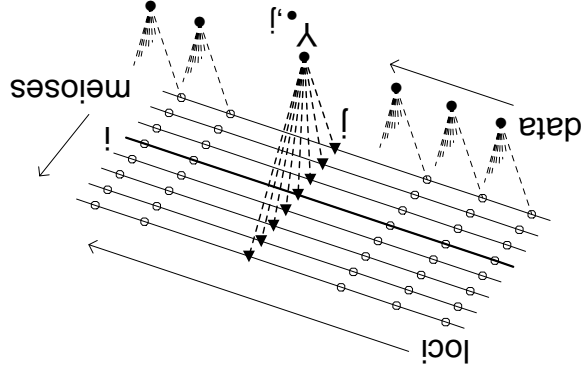
PEDIGREE PEELING: Eiston-Stewart



Conditional on genotypes of parents, grandparent couples and all offspring are all mutually independent.
 Accumulate probabilities over pedigree, using genotypes of (cut-set of) individual(s) as latent state space.
 Linear in pedigree size. Exponential in number of loci.
 (Do with S not genotypes; do by "ordered genotype" not genotype.)

MCMC FOR LARGE PEDIGREES WITH MULTIPLE LOCI

- If pedigree (m) and number of loci (L) are large, we cannot chromosome-peel or pedigree-peel.
- We can still compute $\Pr(Y_{i,j} \bullet | S_{i,j} \bullet)$ by FGL-peeling, quickly and easily (relatively).



The HMM structure lends itself to several **block-Gibbs** schemes, each updating a subset S_u of $\{S_{i,j}\}$ conditional on rest (S_f) and \mathbf{Y} .

BLOCK GIBBS SAMPLERS FOR $S \sim \Pr(\cdot|Y)$

L-sampler: requires peeling over the pedigree to resample $S_{\cdot,j}$.

$$\Pr(Y_{\cdot,j} | S_{\cdot,j-1}, S_{\cdot,j+1}) = \sum_{S_{\cdot,j}} \Pr(Y_{\cdot,j} | S_{\cdot,j}) \Pr(S_{\cdot,j} | S_{\cdot,j-1}, S_{\cdot,j+1})$$

M-sampler: requires peeling along the chromosome (Baum algorithm) using $\Pr(S_{\cdot,j} | Y_{\cdot,1}, \dots, Y_{\cdot,j})$.

$$\Pr(\{S_{i,\cdot}, i \in I_*\} | Y, \{S_{i',\cdot}, i' \notin I_*\})$$

L-sampler mixes poorly for tight linkage

M-sampler mixes poorly on extended pedigrees.

L-sampler is irreducible (theoretically).

Together (**LM-sampler**) they can do better.