
MCMC for the analysis of genetic data on pedigrees: Tutorial Session 2

Elizabeth Thompson
University of Washington

- Genetic mapping and linkage lod scores
- Monte Carlo likelihood and likelihood ratio estimation
- Monte Carlo estimation of linkage lod scores

GENETIC MARKERS

- Human genome: 3×10^9 bp of DNA.
- DNA variants that can be typed in individuals.
 - **Allele** — type of the DNA at position on chromosome
 - Have been mapped: known locations on the genome.
 - **Locus** — position on a chromosome, or DNA at that position
 - Idea: map genes for traits relative to these markers.
 - Microsatellites; lots of alleles; 350 in a genome scan
 - One every 10^7 bp
 - SNPs: typically only two alleles; lots more exist; 1 per 1000 bp

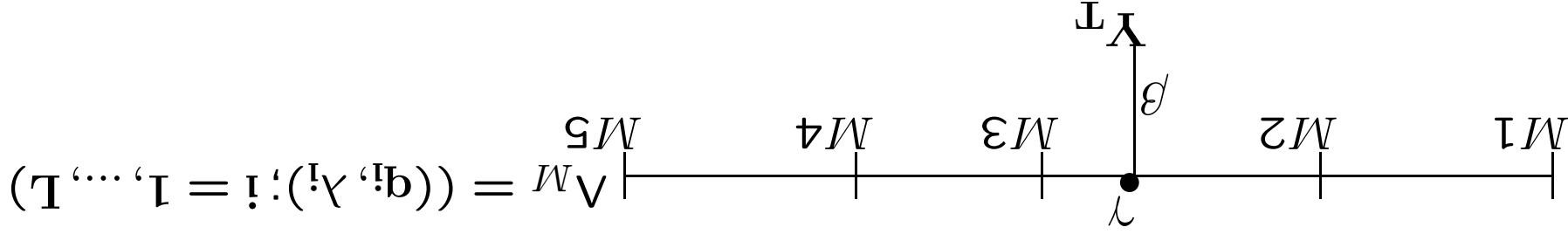
THE STRUCTURE OF A GENETIC MODEL

- Population model: parameters \mathbf{q} , provide probabilities for latent A – allelic types of FGL at each j
- Inheritance model: parameters ρ , 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} .
- $\xi = \xi(\mathbf{q}, \rho, \beta)$.

FROM RECOMBINATION TO LOCATION

- Recall model for $S_{i,\bullet} = (S_{i,1}, \dots, S_{i,T})$:
 $\Pr(S_{i,j} \neq S_{i,j+1}) = d_j$: assumed same $\forall i$ (convenience).
 $S_{i,j}$ assumed Markov in j : no genetic interference.
- Genetic distance d is expected number of crossover events on underlying chromosome: an additive measure.
- Crossovers arise as a Poisson process rate 1 (per Morgan).
- There is a recombination between two loci if there is an **odd number W** of crossovers between them: $W(d) \sim \mathcal{P}(d)$.
- Hence the **Haldane map function**: $d(d) = (1/2)(1 - \exp(-2d))$.
- The key thing is the model: the map function just puts loci onto a linear location map. (See later: MCMC under interference.)

WHAT AND WHY THE LOCATION LOD SCORE



Parameter $\xi = (\beta, \gamma, \Lambda^M)$. Data $\mathbf{Y} = (\mathbf{Y}^M, \mathbf{Y}^T)$

$$\text{lod}(\gamma) = \log_{10} \left(\frac{\text{Pr}(\mathbf{Y}; \Lambda^M, \beta, \gamma)}{\text{Pr}(\mathbf{Y}; \Lambda^M, \beta, \gamma = \infty)} \right)$$

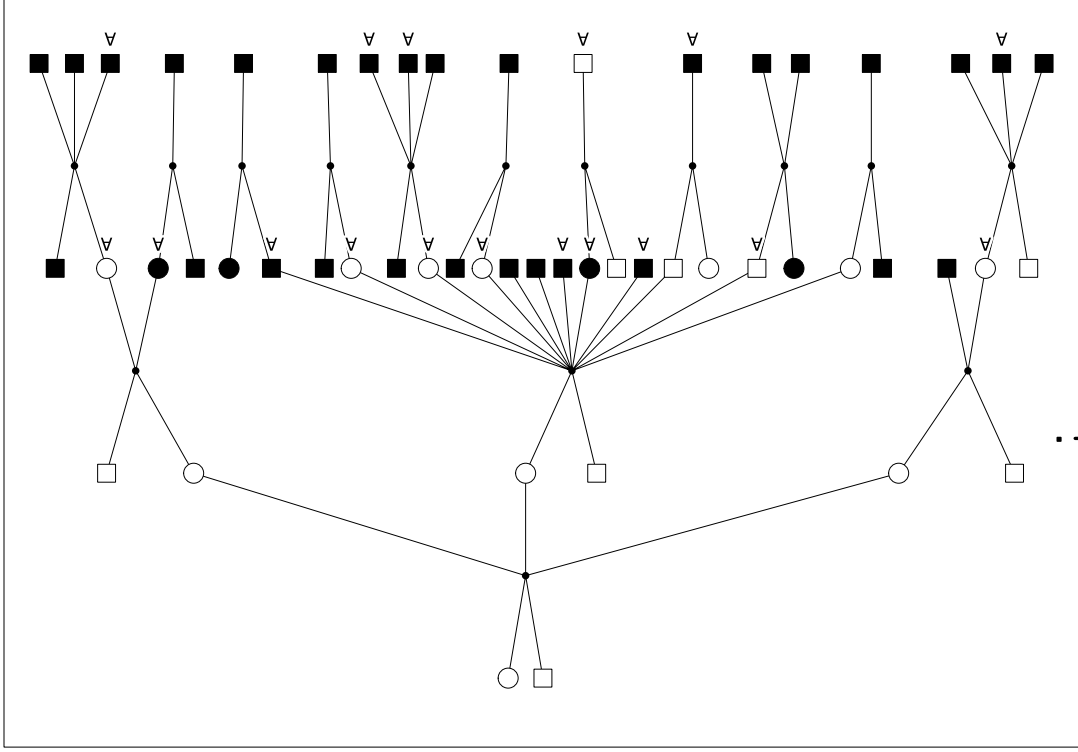
Trait locus location γ is parameter of interest:

$\gamma = \infty$ is no linkage.

Exact computation is infeasible

AN EXAMPLE PEDIGREE: APPROXIMATED

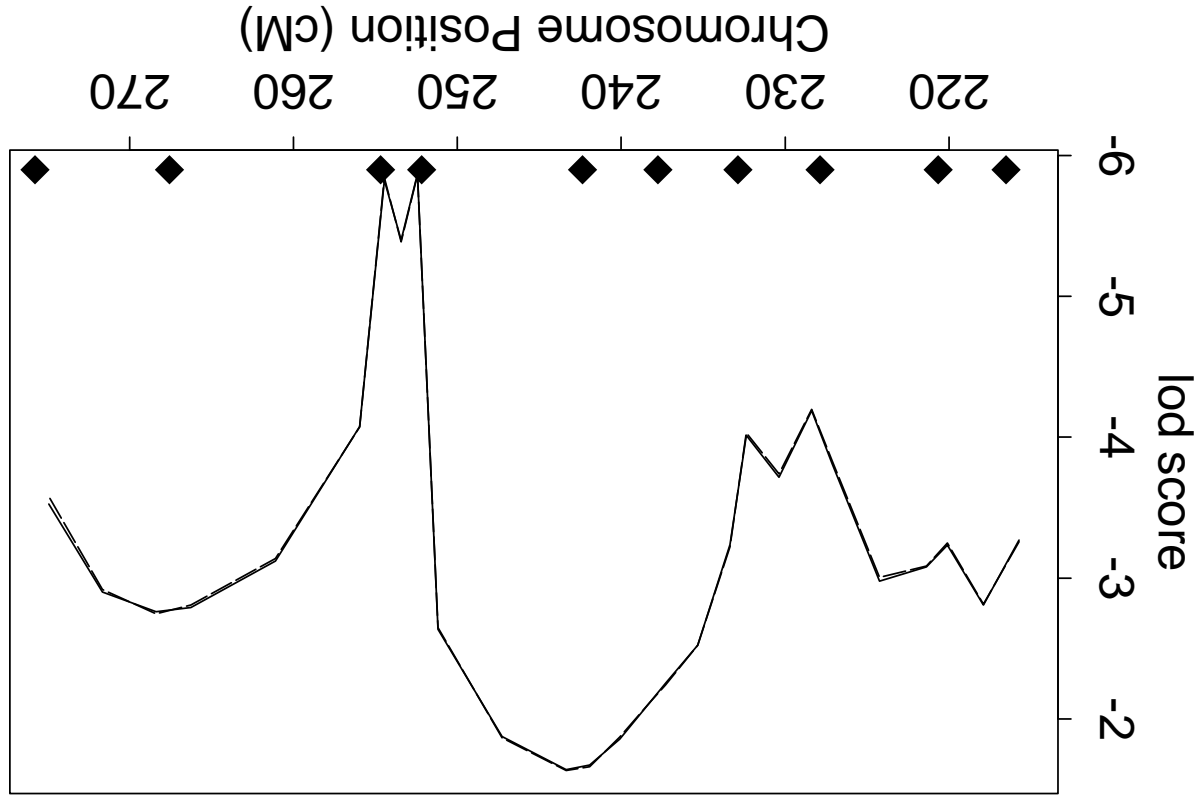
9



SIMPED: disease status and marker availability

Marker data
are SIMULATED
at 10 linked
markers on Chr 1.
Trait is close
to M6

AN EXAMPLE MULTIPPOINT LOD SCORE



MONTE CARLO LIKELIHOODS ON PEDIGREES

- Monte Carlo estimates expectations.

- $L(\xi) = P_\xi(\mathbf{Y}) = \sum_{\mathbf{S}} P_\xi(\mathbf{S}, \mathbf{Y}) = \sum_{\mathbf{S}} P_\xi(\mathbf{Y} | \mathbf{S}) P_\xi(\mathbf{S})$

for parameters ξ and latent variables \mathbf{S} .

- Simple (but not useful) example:

$$L(\xi) = E_\xi(P_\xi(\mathbf{Y} | \mathbf{S}))$$

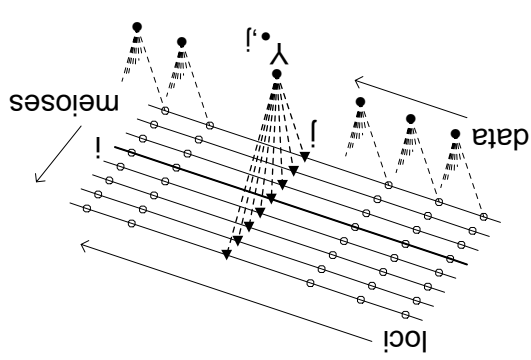
- More generally

$$L(\xi) = \sum_{\mathbf{S}} \left(\frac{P_\xi(\mathbf{S}, \mathbf{Y})}{P_\xi(\mathbf{S}, \mathbf{Y})} \right) P_\xi(\mathbf{S}) = E_{P_\xi} \left(\frac{P_\xi(\mathbf{S}, \mathbf{Y})}{P_\xi(\mathbf{S}, \mathbf{Y})} \right)$$

provided $P_\xi(\mathbf{S}) > 0$ if $P_\xi(\mathbf{S}, \mathbf{Y}) > 0$.

SEQUENTIAL IMPUTATION OVER LOCI

Choose the sampling distribution:



$$\begin{aligned}
 P^*(S_{\cdot,j}) &= P^{\xi_0}(S_{\cdot,j} | S^{*(j-1)}, Y^{(j)}) \\
 &= P^{\xi_0}(S_{\cdot,j} | S^*_{\cdot,1}, \dots, S^*_{\cdot,j-1}) \\
 &= P^{\xi_0}(S_{\cdot,j} | S^*_{\cdot,1}, \dots, Y_{\cdot,j-1}, Y_{\cdot,j}) \\
 &= P^{\xi_0}(S_{\cdot,j} | S^*_{\cdot,j-1}, Y_{\cdot,j})
 \end{aligned}$$

Now:

$$\begin{aligned}
 P^{\xi_0}(S_{\cdot,j} | S^{*(j-1)}, Y^{(j)}) &= \frac{P^{\xi_0}(S_{\cdot,j}, Y_{\cdot,j} | S^{*(j-1)}, Y^{(j-1)})}{P^{\xi_0}(Y_{\cdot,j} | S^{*(j-1)}, Y^{(j-1)})} \\
 &= \frac{P^{\xi_0}(S_{\cdot,j}, Y_{\cdot,j} | S^{*(j-1)}, Y^{(j-1)})}{w_j}
 \end{aligned}$$

where, by pedigree-peeling, we can compute

$$w_j = P^{\xi_0}(Y_{\cdot,j} | Y^{(j-1)}, S^{*(j-1)}) = P^{\xi_0}(Y_{\cdot,j} | S^*_{\cdot,j-1}).$$

MONTÉ CARLO LIKELIHOOD ESTIMATE

Thus sequential imputation distribution is

$$P_{\xi_0}^*(\mathbf{S}_*) = \prod_{j=1}^L P_{\xi_0}^*(S_{\bullet,j} | S_{*(j-1)}, Y_{(j)}) = \frac{W^T(\mathbf{S}_*)}{P_{\xi_0}^*(\mathbf{S}_*, \mathbf{Y})}$$

where $W^T(\mathbf{S}_*) = \prod_{j=1}^L w_j$.

Now

$$L(\xi_0) = P_{\xi_0}^*(\mathbf{Y}) = E_{P^*} \left(\frac{P_{\xi}^*(\mathbf{S}, \mathbf{Y})}{P^*(\mathbf{S})} \right) = E_{P^*}(W^T(\mathbf{S}_*))$$

Given N realizations $\mathbf{S}^{(t)}$ the estimate of $L(\xi_0)$ is $N^{-1} \sum_{t=1}^N W^T(\mathbf{S}^{(t)})$.

THE IDEAL SAMPLING DISTRIBUTION

- We want $P^*(\mathbf{S})$ close to proportional to $P^{\xi_0}(\mathbf{Y}, \mathbf{S})$
 - that is $P^*(\mathbf{S}) \approx P^{\xi_0}(\mathbf{S}|\mathbf{Y})$.

- Of course we cannot achieve this, else Monte Carlo would be unnecessary.

- Suppose we use MCMC to sample \mathbf{S} from $P^{\xi_0}(\mathbf{S}|\mathbf{Y})$.

$$\begin{aligned}
 P^{\xi}(\mathbf{Y}) &= \sum_{\mathbf{S}} P^{\xi}(\mathbf{Y}, \mathbf{S}) = \sum_{\mathbf{S}} \frac{P^{\xi}(\mathbf{Y}, \mathbf{S})}{P^{\xi_0}(\mathbf{S}|\mathbf{Y})} P^{\xi_0}(\mathbf{S}|\mathbf{Y}) \\
 &= \mathbb{E}^{\xi_0} \left(\frac{P^{\xi}(\mathbf{Y}, \mathbf{S})}{P^{\xi_0}(\mathbf{S}|\mathbf{Y})} \mid \mathbf{Y} \right) \\
 &= P^{\xi_0}(\mathbf{Y}) \mathbb{E}^{\xi_0} \left(\frac{P^{\xi}(\mathbf{Y}, \mathbf{S})}{P^{\xi_0}(\mathbf{Y}, \mathbf{S})} \mid \mathbf{Y} \right)
 \end{aligned}$$

LIKELIHOOD RATIO ESTIMATION

Thus we have

$$\frac{L(\xi)}{L(\xi_0)} = \frac{P_{\xi}(\mathbf{Y})}{P_{\xi_0}(\mathbf{Y})} = E_{\xi_0} \left(\frac{P_{\xi}(\mathbf{Y}, \mathbf{S})}{P_{\xi_0}(\mathbf{Y}, \mathbf{S})} \mid \mathbf{Y} \right)$$

\mathbf{S} is the random variable, \mathbf{Y} is fixed. $\mathbf{S} \sim P_{\xi_0}(\cdot | \mathbf{Y})$.

If $\mathbf{S}^{(\tau)}$, $\tau = 1, \dots, N$, are realized from $P_{\xi_0}(\cdot | \mathbf{Y})$ then the likelihood

ratio can be estimated by

$$\frac{1}{N} \sum_{\tau=1}^{\tau} \left(\frac{P_{\xi}(\mathbf{Y}, \mathbf{S}^{(\tau)})}{P_{\xi_0}(\mathbf{Y}, \mathbf{S}^{(\tau)})} \right)$$

LINKAGE LOCATION LIKELIHOOD RATIO

The form for linkage lod that follows directly from this is

$$\frac{L(\beta, \gamma_1, \Lambda_M)}{L(\beta, \gamma_0, \Lambda_M)} = E_{\xi_0} \left(\frac{P^{\xi_1}(\mathbf{Y}_T, \mathbf{Y}_M, \mathbf{S}_T, \mathbf{S}_M)}{P^{\xi_0}(\mathbf{Y}_T, \mathbf{Y}_M, \mathbf{S}_T, \mathbf{S}_M)} \mid \mathbf{Y}_T, \mathbf{Y}_M \right)$$

for two hypothesized trait locus positions γ_1 and γ_0 .

Now $P_\xi(\mathbf{Y}, \mathbf{S}) = P^\beta(\mathbf{Y}_T | \mathbf{S}_T) P^{\Lambda_M}(\mathbf{Y}_M, \mathbf{S}_M) P^\gamma(\mathbf{S}_T | \mathbf{S}_M)$ so ratio re-

duces to

$$\frac{L(\beta, \gamma_1, \Lambda_M)}{L(\beta, \gamma_0, \Lambda_M)} = E_{\xi_0} \left(\frac{P^{\gamma_1}(\mathbf{S}_T | \mathbf{S}_M)}{P^{\gamma_0}(\mathbf{S}_T | \mathbf{S}_M)} \mid \mathbf{Y}_T, \mathbf{Y}_M \right)$$

LOCAL ESTIMATE IS VERY SIMPLE:
GLOBAL IS HARD

... l T r ...

$$P^{\gamma_1}(\mathbf{S}_T | \mathbf{S}_M) \frac{P^{\gamma_0}(\mathbf{S}_T | \mathbf{S}_M)}{P^{\gamma_1}(\mathbf{S}_T | \mathbf{S}_M)} = \prod_i \left[\begin{matrix} \left(\frac{p_{1r}}{p_{0r}} \right)^{|S_{i,T} - S_{i,r}|} & \left(\frac{p_{0l}}{p_{1l}} \right)^{|S_{i,T} - S_{i,l}|} \\ \left(\frac{1 - p_{1r}}{1 - p_{0r}} \right)^{1 - |S_{i,T} - S_{i,r}|} & \left(\frac{1 - p_{0l}}{1 - p_{1l}} \right)^{1 - |S_{i,T} - S_{i,l}|} \end{matrix} \right]$$

- The above works well only for $\gamma_1 \approx \gamma_0$, and for γ_0, γ_1 with same l and r .
- When likelihoods are not smooth, combining LR estimates does not work well – especially across markers.

AN MCMC IMPORTANCE SAMPLING ESTIMATE

Lange and Sobel (1996) write the likelihood in the form

$$\begin{aligned}
 L(\beta, \gamma, \Lambda^M) &= P_{\beta, \gamma, \Lambda^M}(\mathbf{Y}^M, \mathbf{Y}^T) \propto P_{\beta, \gamma, \Lambda^M}(\mathbf{Y}^T | \mathbf{Y}^M) \\
 &= \sum_{S^M} P_{\beta, \gamma}(\mathbf{Y}^T | S^M) P_{\Lambda^M}(S^M | \mathbf{Y}^M) \\
 &= E_{\Lambda^M}(P_{\beta, \gamma}(\mathbf{Y}^T | S^M) | \mathbf{Y}^M).
 \end{aligned}$$

- Sample S^M given \mathbf{Y}^M : compute $P(\mathbf{Y}^T | S^M) \forall \beta, \gamma$
- a form of **Rao-Blackwellization** – integrate over S^T .
- Also **importance sampling**: maybe $P(S^M | \mathbf{Y}^M) \approx P(S^M | \mathbf{Y}^M, \mathbf{Y}^T)$
- For “fuzzy” traits it works quite well.

METROPOLIS HASTINGS FOR INTERFERENCE

- Suppose we have interference model $P^{(I)}(\mathbf{S})$ in place of Haldane model $P^{(H)}(\mathbf{S})$ we have used so far.

- Use block-Gibbs update of meiosis i (S_i^{\bullet}) to propose S_i^{\dagger} .

- Hastings ratio is for current \mathbf{S} and proposed S_i^{\dagger} is

$$\begin{aligned}
 & \frac{P^{(I)}(\mathbf{S}^{\dagger}, \mathbf{Y}) P^{(H)}(S_i^{\bullet} | S_i^{\dagger}, \mathbf{S}^{\dagger}, k, \mathbf{Y})}{P^{(I)}(\mathbf{S}, \mathbf{Y}) P^{(H)}(S_i^{\dagger} | S_i^{\bullet}, \mathbf{S}, k, \mathbf{Y})} = \frac{P^{(I)}(\mathbf{S}^{\dagger}, \mathbf{Y}) P^{(I)}(\mathbf{S}, \mathbf{Y})}{P^{(I)}(\mathbf{S}, \mathbf{Y}) P^{(I)}(\mathbf{S}^{\dagger}, \mathbf{Y})} \\
 & = \frac{P^{(I)}(\mathbf{S}^{\dagger}, \mathbf{Y}) P^{(I)}(\mathbf{S}, \mathbf{Y})}{P^{(I)}(\mathbf{S}, \mathbf{Y}) P^{(I)}(\mathbf{S}^{\dagger}, \mathbf{Y})} = \frac{P(\mathbf{S}^{\dagger} | \mathbf{S})}{P(\mathbf{S} | \mathbf{S}^{\dagger})}
 \end{aligned}$$

INTERFERENCE ctd.

$$h(\mathbf{S}^\dagger; \mathbf{S}) = \prod_{k=1}^m \frac{P(I)(S_k^\dagger, \bullet) P(H)(S_k^\dagger, \bullet)}{P(I)(S_k, \bullet) P(H)(S_k, \bullet)} = \frac{P(I)(S_i^\dagger, \bullet) P(H)(S_i^\dagger, \bullet)}{P(I)(S_i, \bullet) P(H)(S_i, \bullet)}$$

- $\Pr(\mathbf{S}^* = \mathbf{S}^\dagger) = a = \min(1, h)$. $\Pr(\mathbf{S}^* = \mathbf{S}) = 1 - a$.

- Question: better to sample under H and reweight, or use M-H to sample under model I ?