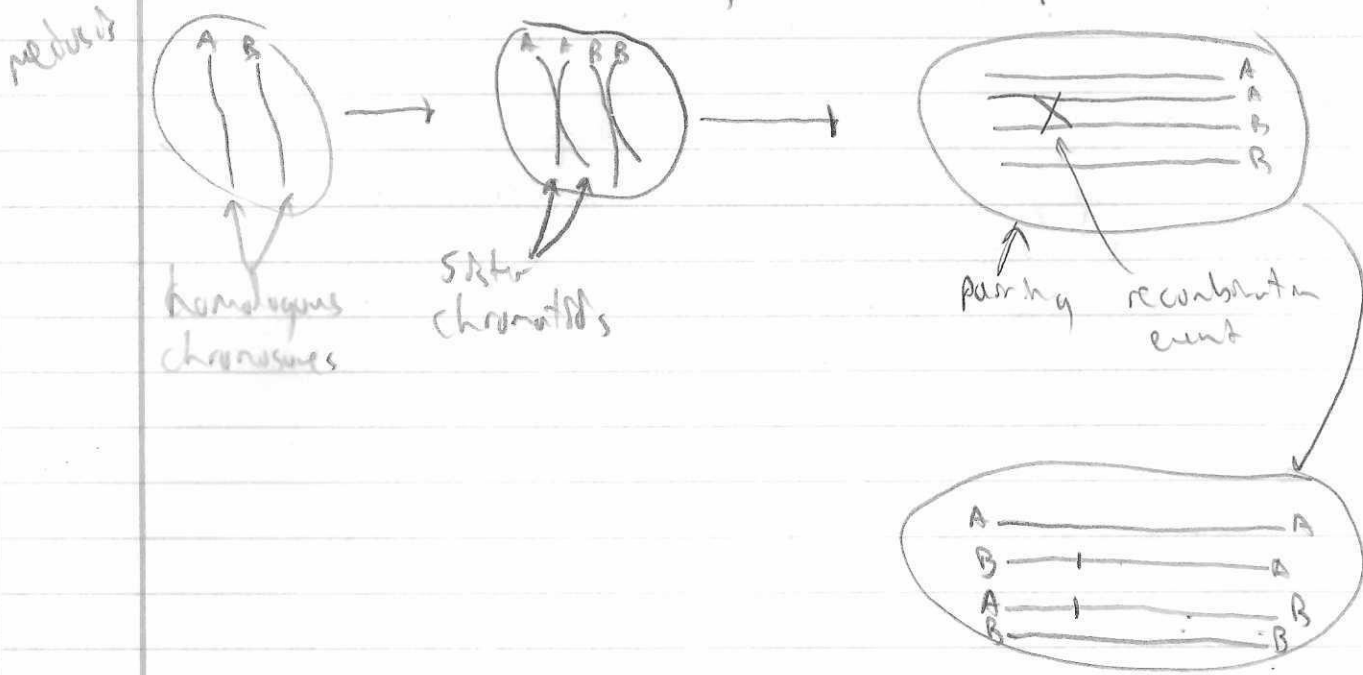


# Li & Stephens (2003)

- Recombination: In meiosis, chromosomes pair and recombine



- Recombination makes it so that loci are not perfectly dependent.
- Dependence between Loci = "LD"
- This article relates LD to recombination  $\rho$  through a model
- Let  $h_n \in \{0, 1\}^S$  be the haplotype for individual  $n$  at  $S$  loci.

Eg, if  $S = 3$ , the possible values of  $h_n$ :

$$h_n \in \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

$$I = A, \quad U = a$$

• We want a likelihood for  $h_1, h_2, \dots, h_n$  given the recombination rate  $\rho$ .

• By definition of conditional probability,

$$P_r(h_1, h_2, \dots, h_n | \rho) = P_r(h_1 | \rho) P_r(h_2 | h_1, \rho) \dots P_r(h_n | h_1, h_2, \dots, h_{n-1}, \rho)$$

• Idea: Find approximation st.

$$P_r(h_1, h_2, \dots, h_n | \rho) \approx \hat{\pi}(h_1 | \rho) \hat{\pi}(h_2 | h_1, \rho) \dots \hat{\pi}(h_n | h_1, h_2, \dots, h_{n-1}, \rho)$$

• Ewen's sampling formulae approximation ( $\theta = 4N\mu$ )

$h_k$  w.p.  $\frac{k}{k+\theta}$  will equal one of  $(h_1, \dots, h_{k-1})$

↑ choose which one uniformly

w.p.  $1 - \frac{k}{k+\theta}$ , will be something completely different

Issues: (1) Does not indicate that # differences should be small  
(2) Does not account for recombination / LD

• Li and Stephen's  $\hat{\pi}$  (Appendix A)

Haplotypes  $h_1, h_2, \dots, h_n$  (so  $n/2$  diploid individuals)

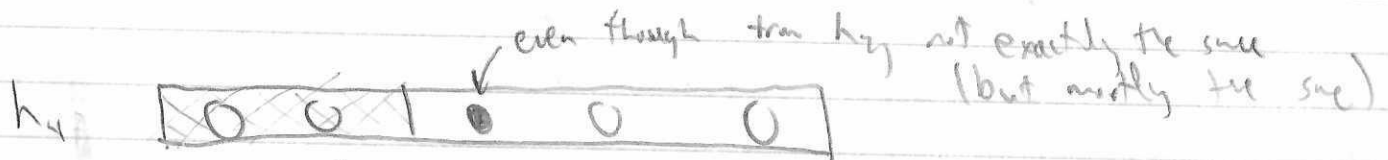
$$\hat{\pi}(h_1) = \frac{1}{2s} \quad (\text{all haplotypes are equally likely})$$

$\pi(h_{k+1} | h_1, \dots, h_k)$

↳ imperfect mosaic of the other  $k$  haplotypes



$h_4$  is sometimes from  $h_1, h_2,$  or  $h_3,$  But it may or may not have copied those values.



↑ adjacent values

are more likely from same haplotype (LD/e)

$e \uparrow \Rightarrow$  less likely adjacent pairs are from same haplotype

$e \downarrow \Rightarrow$  more likely

• Let  $X_j =$  which haplotype  $h_{k+1}$  comes from at site  $j$

$$X_j \in \{1, 2, \dots, k\}$$

Es. in example  $X = (3, 3, 2, 2, 2)$

• Model using a Markov chain:

$$P_r(X_1 = x) = \frac{1}{k} \text{ for } x \in \{1, \dots, k\}$$

$$P_r(X_{j+1} = x' \mid X_j = x) = \begin{cases} \exp\{-e d_j / h\} + (1 - \exp\{-e d_j / h\}) \frac{1}{k} & \text{if } x' = x \\ (1 - \exp\{-e d_j / h\}) \frac{1}{k} & \text{otherwise} \end{cases}$$

$d_j =$  distance between adjacent loci (known)

If  $e = 0$  (no recombination)

$$= \begin{cases} 1 & x' = x \\ 0 & \text{otherwise} \end{cases}$$

$\uparrow$  Newer switches.

If  $e = \infty$

$$= \begin{cases} \frac{1}{k} & \text{if } x' = x \\ \frac{1}{k} & \text{otherwise} \end{cases}$$

↑ No dependence between adjacent loci.

• Model Mutation via imperfect copying

Copy is exact w.p.  $\frac{k}{k+\theta}$

Copy is imperfect w.p.  $1 - \frac{k}{k+\theta}$

} via Even's sampling

If imperfect, then 0 w.p.  $\frac{1}{2}$  and 1 w.p.  $\frac{1}{2}$

$$P_r(h_{t+1,j} = a \mid X_j = x, h_{t,j}, h_t) = \begin{cases} \frac{k}{k+\theta} + \frac{1}{2} \frac{\theta}{k+\theta} & \text{if } h_{t,j} = a \\ \frac{1}{2} \frac{\theta}{k+\theta} & \text{if } h_{t,j} \neq a \end{cases}$$

same

different

$$\text{we } \theta = \left( \sum_{m=1}^{M-1} \frac{1}{m} \right)^{-1}$$

HW:

1.) What is length of  $h_1$ ?

(5)

2.) What are possible values of  $h_5$ ?

(0 or 1)

3.)  $n=2$ , what is  $\bar{\theta}$ ?

$$\bar{\theta} = \left( \sum_{n=1}^2 \frac{1}{n} \right)^{-1} = (1)$$

4.)  $n=2$ ,  $\Pr(h_{2j}=0 \mid h_{1j}=0)$ ?

$$= \frac{k}{k+\theta} + \frac{1}{2} \frac{\theta}{1+\theta}$$

$$= \frac{1}{1+1} + \frac{1}{2} \frac{1}{1+1} = \frac{1}{2} + \frac{1}{4} = \left( \frac{3}{4} \right)$$

5.)  $n=2$ ,  $\Pr(h_{2j}=0 \mid h_{1j}=1)$ ?

$$= \frac{1}{2} \frac{\theta}{k+\theta} = \frac{1}{2} \frac{1}{1+1} = \left( \frac{1}{4} \right)$$

6.) Suppose  $\rho = \infty$ , what is  $\Pr(X_{n+1}=1)$ ?

Did this in the notes.  $\left( \frac{1}{k} \right)$

Given observed haplotypes  $h_1, \dots, h_n$ , very easy to compute likelihood and maximize it over  $\theta$

↑ Markov chains make this efficient

- Forward/backward algorithm: for computing forward estimators  $\theta$

- Forward part can calculate likelihood

- Viterbi algorithm can give us hidden states

↑ easy and computationally fast (but 5 lectures to learn them)

• Why use  $\theta = \left( \sum_{m=1}^n \frac{1}{m} \right)^{-1}$ ?

$\theta \sum_{m=1}^n \frac{1}{m}$  is expected # mutations on a tree with  $n$  alleles

↑ (we did this in Gillespie)

Set Expected number to be 1 at each site (since we know it occurred)

$$\Rightarrow \theta \approx \frac{1}{\sum_{m=1}^n \frac{1}{m}}$$

Likelihood depends on order  $h_1, h_2, \dots, h_n$

↳ so average  $\approx 20$  random orderings of haplotypes

• Why that form for transition?

Suppose  $d \in (0, \infty)$ , what is distribution of # jumps?

Does not jump up:  $\exp\{-\rho d/k\}$  ← true for all  $d$

jumps (possibly to same haplotype) up:  $1 - \exp\{-\rho d/k\}$

Let  $D =$  waiting time

$$f(D=d) \propto \exp\left\{-d \frac{\rho}{k}\right\}$$

↑ kernel of exponential

$$E[D] = \frac{k}{\rho}$$

If  $\rho \uparrow$ , then waiting time is small  
If  $\rho \downarrow$ , then waiting time is large

If  $k \uparrow$ , then waiting time is large  
If  $k \downarrow$ , then waiting time is small

If  $k \uparrow$ , then more likely to have a close relative that is more similar in sample, so stay on that relative longer.

Phylogenetics suggests  $E[D] \approx \frac{\rho}{k+\rho} \approx \frac{\rho}{k}$  for small  $\rho$