

*Multipoint Analysis for  
Sibling Pairs*

**Biostatistics 666**

**Lecture 18**

## Previously ...

---

- Linkage analysis with pairs of individuals
  - “Non-parametric” IBS Methods
  - “Maximum Likelihood” IBD Based Method
  - Possible Triangle Constraint

## ASP Methods Covered So Far ...

---

- Increasing degrees of sophistication and complexity
- In each case, only a single marker is evaluated...

# IBS Based Linkage Test

---

$$\chi^2_{2df} = \sum_i \frac{[N_{IBS=i} - E(N_{IBS=i})]^2}{E(N_{IBS=i})}$$

$$LOD = \frac{\chi^2}{2 \ln 10}$$

- Expect counts calculated using:
  - Allele frequencies for marker
  - Relationship for affected individuals

# Likelihood for a Single ASP

---

$$L_i = \sum_{j=0}^2 P(IBD = j | ASP) P(Genotypes | IBD = j) = \sum_{j=0}^2 z_j w_{ij}$$

Risch (1990) defines

$$w_{ij} = P(Genotypes_i | IBD = j)$$

We only need proportionate  $w_{ij}$

# MLS Linkage Test

---

$$L(z_0, z_1, z_2) = \prod_i \sum_j z_j w_{ij}$$

$$LOD = \log_{10} \prod_i \frac{z_0 w_{i0} + z_1 w_{i1} + z_2 w_{i2}}{\frac{1}{4} w_{i0} + \frac{1}{2} w_{i1} + \frac{1}{4} w_{i2}}$$

The MLS statistic is the LOD evaluated at the MLEs of  $z_0, z_1, z_2$

## Possible Triangle Constraint

---

- For any genetic model, we expect ASPs to be more similar than unselected pairs of siblings.
- More precisely, Holmans (1993) showed that for any genetic model
  - $z_2 \geq 1/4$
  - $z_1 \leq 1/2$  and  $z_1 \geq 2 z_0$
  - $z_0 \leq 1/4$

## Further Improvements ...

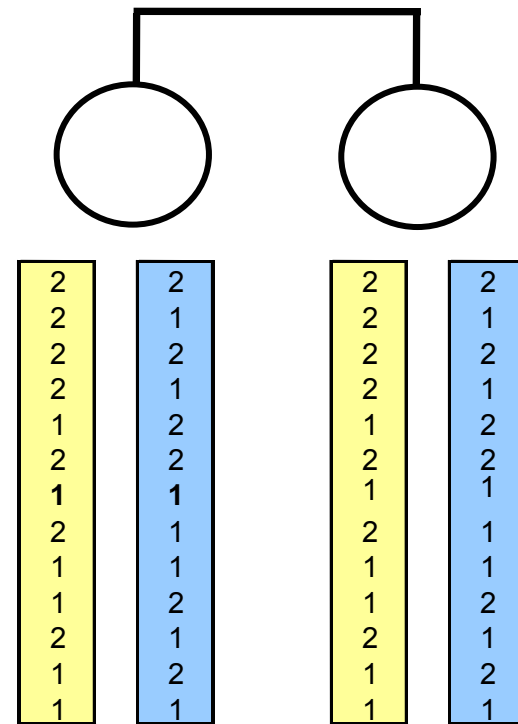
---

- All these methods lose information when a marker is uninformative in a particular family...
- Today, we will see how to use neighboring markers to extract more information about IBD.



# Intuition For Multipoint Analysis

- IBD changes infrequently along the chromosome
- Neighboring markers can help resolve ambiguities about IBD sharing
- In the Risch approach, they might ensure that, effectively, only one  $w$  is non-zero



# Today ...

---

- Framework for multipoint calculations
  - First, likelihood of genotypes for series of markers
  - Discuss application to the MLS linkage test
  - Later, we will use it for useful applications such as error detection and relationship inference
- Refresher on IBD probabilities
- Using a Markov Chain to speed analyses

# Ingredients

---

$X_1$

$X_2$

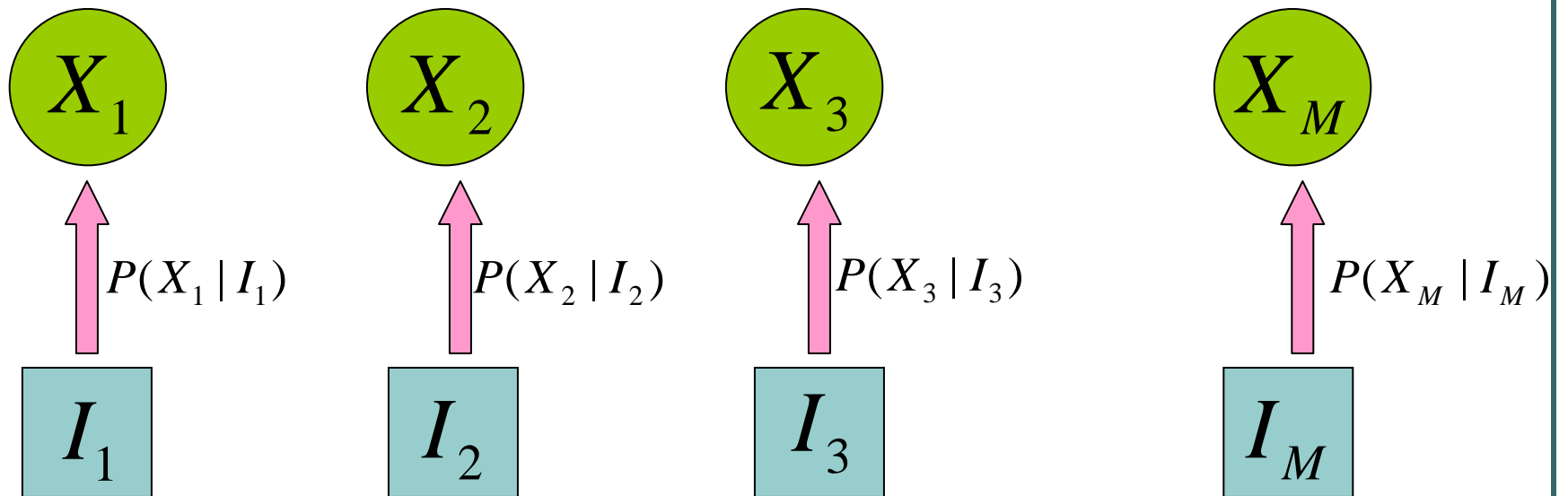
$X_3$

$X_M$

One ingredient will be the observed genotypes at each marker ...

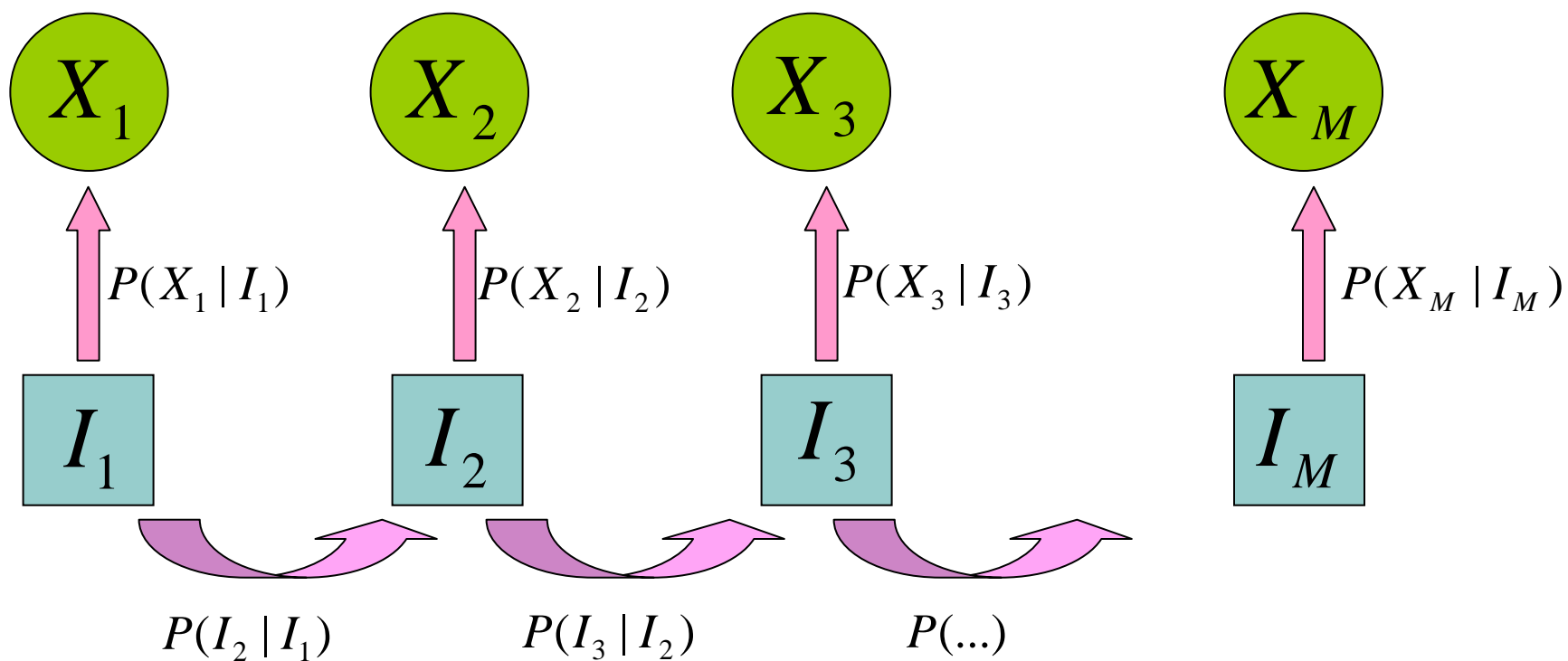
# Ingredients

---



Another ingredient will be the possible IBD states at each marker ...

# Ingredients



**The final ingredient connects IBD states along the chromosome ...**

## The Likelihood of Marker Data

---

$$L = \sum_{I_1} \sum_{I_2} \dots \sum_{I_M} P(I_1) \prod_{i=2}^M P(I_i | I_{i-1}) \prod_{i=1}^M P(X_i | I_i)$$

- General, but slow unless there are only a few markers.
- Combined with Bayes' Theorem can estimate probability of each IBD state at any marker.

## The Ingredients ...

---

- Probability of observed genotypes at each marker conditional on IBD state
- Probability of changes in IBD state along chromosome
- Hidden Markov Model

$$P(I_1)$$

Prior Probability of IBD States



## IBD Probabilities

---

- Number of alleles identical by descent
- For sibling pairs, must be:
  - 0
  - 1
  - 2
- Not always determined by marker data

$$P(X_i | I_i)$$

Probability of Observed  
Genotypes, Given IBD State

$$P(X_m | I_m)$$

Sib	CoSib	IBD		
		0	1	2
(a,b)	(c,d)	$4p_a p_b p_c p_d$	0	0
(a,a)	(b,c)	$2p_a^2 p_b p_c$	0	0
(a,a)	(b,b)	$p_a^2 p_b^2$	0	0
(a,b)	(a,c)	$4p_a^2 p_b p_c$	$p_a p_b p_c$	0
(a,a)	(a,b)	$2p_a^3 p_b$	$p_a^2 p_b$	0
(a,b)	(a,b)	$4p_a^2 p_b^2$	$(p_a p_b^2 + p_a^2 p_b)$	$2p_a p_b$
(a,a)	(a,a)	$p_a^4$	$p_a^3$	$p_a^2$
Prior Probability		$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

Note: Assuming unordered genotypes

Question:

What to do about missing data?

---

- What happens when some genotype data is unavailable?

$$P(I_{i+1} | I_i)$$

Model for Transitions in IBD  
Along Chromosome

$$P(I_{m+1} | I_m)$$

---

- Depends on recombination fraction  $\theta$ 
  - This is a measure of distance between two loci
  - Probability grand-parental origin of alleles changes between loci
- Naturally, leads to probability of change in IBD:

$$\psi = 2\theta(1 - \theta)$$

$$P(I_{m+1} | I_m)$$


---

		IBD State at m + 1		
		0	1	2
IBD state at marker m	0	$(1-\psi)^2$	$2\psi(1-\psi)$	$\psi^2$
	1	$\psi(1-\psi)$	$(1-\psi)^2 + \psi^2$	$\psi(1-\psi)$
	2	$\psi^2$	$2\psi(1-\psi)$	$(1-\psi)^2$

$$\psi = 2\theta(1 - \theta)$$

$$P(I_1)$$
$$P(X_i|I_i)$$
$$P(I_{i+1}|I_i)$$

**All the Ingredients!**



# Example

---

- Consider two loci separated by  $\theta = 0.1$
- Each loci has two alleles, each with frequency .50
- If two siblings have the following genotypes:

	Sib1	Sib2
● Marker A:	1/1	2/2
● Marker B:	1/1	1/1
- What is the probability of IBD=2 at marker B when...
  - You consider marker B alone?
  - You consider both markers simultaneously?

## The Likelihood of Marker Data

---

$$L = \sum_{I_1} \sum_{I_2} \dots \sum_{I_M} P(I_1) \prod_{i=2}^M P(I_i | I_{i-1}) \prod_{i=1}^M P(X_i | I_i)$$

- General, but slow unless there are only a few markers.
- How do we speed things up?

# A Markov Model

---

- Re-organize the computation slightly, to avoid evaluating nested sum directly
- Three components:
  - Probability considering a single location
  - Probability including left flanking markers
  - Probability including right flanking markers
- Scale of computation increases linearly with number of markers

## A Markov Rearrangement ...

---

$$LEFT_1(j) = P(IBD = j)P(X_1 | I_1 = j)$$

$$LEFT_{i+1}(j) = \sum_{k=0,1,2} LEFT_i(k)P(I_{i+1} = j | I_i = k)P(X_{i+1} | I_{i+1} = j)$$

$$L = \sum_{k=0,1,2} LEFT_{last}(k)$$

- Using this arrangement, we calculate the likelihood by:
  - Evaluating *LEFT* function at the first position
  - Evaluating *LEFT* function along chromosome
    - Each time, re-using results from the previous position only
    - Required effort increases linearly with number of markers
  - Final summation gives overall likelihood

## Improvements ...

---

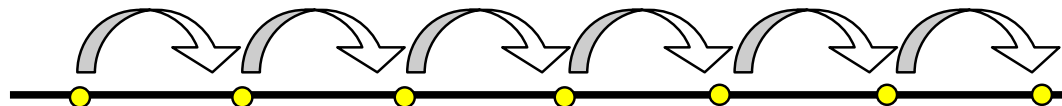
- The previous arrangement, quickly gives the likelihood for any number of markers
- A more flexible arrangement would allow us to quickly calculate conditional IBD probabilities along chromosome...

# A More Flexible Arrangement...

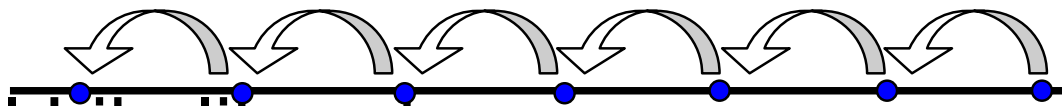
- Single Marker



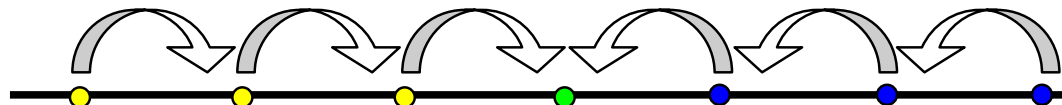
- Left Conditional



- Right Conditional



- Full Likelihood



# The Likelihood of Marker Data

---

$$\begin{aligned} L &= \sum_{I_j} P(I_j) P(X_j | I_j) P(X_1 \dots X_{j-1} | I_j) P(X_{j+1} \dots X_M | I_j) \\ &= \sum_{I_j} P(I_j) P(X_j | I_j) L_j(I_j) R_j(I_j) \end{aligned}$$

- A different arrangement of the same likelihood
- The nested summations are now hidden inside the  $L_j$  and  $R_j$  functions...

## Left-Chain Probabilities

---

$$\begin{aligned} L_m(I_m) &= P(X_1, \dots, X_{m-1} | I_m) \\ &= \sum_{I_{m-1}} L_{m-1}(I_{m-1}) P(X_{m-1} | I_{m-1}) P(I_{m-1} | I_m) \end{aligned}$$

$$L_1(I_1) = 1$$

- Proceed one marker at a time.
- Computation cost increases linearly with number of markers.



## Right-Chain Probabilities

---

$$\begin{aligned} R_m(I_m) &= P(X_{m+1}, \dots, X_M | I_m) \\ &= \sum_{I_{m+1}} R_{m+1}(I_{m+1}) P(X_{m+1} | I_{m+1}) P(I_{m+1} | I_m) \end{aligned}$$

$$R_M(I_M) = 1$$

- Proceed one marker at a time.
- Computation cost increases linearly with number of markers.

## Extending the MLS Method ...

---

$$\begin{aligned}w_j &= P(X_j | I_j)P(X_1 \dots X_{j-1} | I_j)P(X_{j+1} \dots X_M | I_j) \\ &= P(X_j | I_j)L_j(I_j)R_j(I_j)\end{aligned}$$

- We just change the definition for the “weights” given to each configuration!

## Some Extensions We'll Discuss

---

- Modeling error
  - What components might have to change?
- Modeling other types of relatives
  - What components might have to change?
- Modeling larger pedigrees

# Today

---

- Efficient computational framework for multipoint analysis of sibling pairs