



InCoB 2012

# A Linear-time Algorithm for Reconstructing Zero-Recombinant Haplotype Configuration on a Pedigree

June En-Yu Lai

TIGP Bioinformatics Program of Academia Sinica, **Taiwan**

October 3, 2012

# Genotype and Haplotype

For an individual,

		marker loci						
paternal haplotype		A	C	A	T	G	T	C
maternal haplotype	+	A	G	G	C	G	A	G
genotype		AA	CG	AG	TC	GG	TA	CG

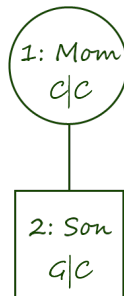
The haplotype structure of a human genome is not available directly from the genotyping and the unordered genotype data does not tell us which allele comes from which parent.

# Pedigree + Genotype $\rightarrow$ Haplotype?

For an individual (son),

		marker loci			
		$\underbrace{\hspace{10em}}$			
paternal haplotype		A	<b>G</b>	A	T
maternal haplotype	+	A	<b>C</b>	G	T
genotype		A A	<b>GC</b>	AG	T T

The parent-child relationship provides the information of inheritance that will help to determine haplotype.



# Zero-Recombinant Haplotype Configuration (1)

- ▶ Input: pedigree + genotype
- ▶ Output: haplotype
- ▶ Assumptions:
  1. The input dataset is free of mutation and recombination.
  2. The input dataset is free of genotyping errors.
  3. All alleles are bi-allelic (denoted by 0 or 1).

## Zero-Recombinant Haplotype Configuration (2)

For a locus,

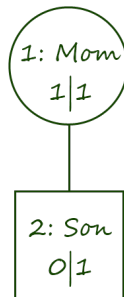
		individuals						
paternal alleles		A	G	A	G	G	A	G
maternal alleles	+	A	A	G	A	G	G	G
genotype		AA	AG	AG	GA	GG	AG	GG
paternal alleles		1	0	1	0	0	1	0
maternal alleles	+	1	1	0	1	0	0	0
genotype		1	2	2	2	0	2	0

## Zero-Recombinant Haplotype Configuration (3)

For a locus,

		individuals						
		<div></div>						
paternal alleles		1	0	1	0	0	1	0
maternal alleles	+	1	1	0	1	0	0	0
genotype		<b>1</b>	<b>2</b>	2	2	<b>0</b>	2	<b>0</b>

If the genotype is homozygous or is the children of the homozygous, it is referred to as **predetermined**.



## A linear system to represent Mendelian law (1)

What is passed = What is received

$$\underline{\mathbf{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j}} = \mathbf{paternal}_j[\text{locus}] + d_{i,j}[\text{locus}]$$

$$\mathbf{paternal}_i[\text{locus}] = \begin{cases} 0 & \text{if paternal allele is 0} \\ 1 & \text{if paternal allele is 1.} \end{cases}$$

$$w_i[\text{locus}] = \begin{cases} 0 & \text{if } \mathbf{genotype}_i[\text{locus}] \text{ is homozygous} \\ 1 & \text{if } \mathbf{genotype}_i[\text{locus}] \text{ is heterozygous.} \end{cases}$$

The variable  $\mathbf{paternal}_i$  and the constant  $w_i$  describe the information of the individual  $i$ .

## A linear system to represent Mendelian law (2)

What is passed = What is received

$$\mathbf{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} = \mathbf{paternal}_j[\text{locus}] + \underline{d_{i,j}[\text{locus}]}$$

$$\mathbf{inheritance}_{i,j} = \begin{cases} 0 & \text{if } i \text{ passes its paternal allele to } j \\ 1 & \text{if } i \text{ passes its maternal allele to } j. \end{cases}$$

$$d_{i,j}[\text{locus}] = \begin{cases} 0 & \text{if } i \text{ is } j\text{'s father} \\ w_j[\text{locus}] & \text{if } i \text{ is } j\text{'s mother.} \end{cases}$$

The variable  $\mathbf{inheritance}_{i,j}$  and the constant  $d_{i,j}$  describe the relationship between the individual  $i$  and the individual  $j$ .



## A linear system to represent Mendelian law (3)

What is passed = What is received

$$\mathbf{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} = \mathbf{paternal}_j[\text{locus}] + d_{i,j}[\text{locus}]$$

An example.

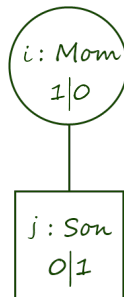
$$\mathbf{paternal}_i[\text{locus}] = 1$$

$$w_i[\text{locus}] = 1$$

$$\mathbf{inheritance}_{i,j} = 0$$

$$\mathbf{paternal}_j[\text{locus}] = 0$$

$$d_{i,j}[\text{locus}] = w_j[\text{locus}] = 1$$



# The concept of a **constraint** (1)

$$\mathbf{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} = \mathbf{paternal}_j[\text{locus}] + d_{i,j}[\text{locus}]$$
$$\mathit{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} = \mathit{paternal}_j[\text{locus}] + d_{i,j}[\text{locus}]$$

Recall that if the genotype is homozygous or is the children of the homozygous, it is referred to as **predetermined** and the value of **paternal** is also known.

## The concept of a **constraint** (2)

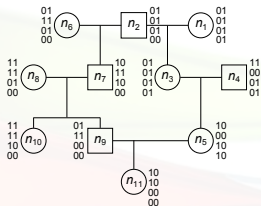
$$\begin{array}{l}
 \text{paternal}_i[\text{locus}] + w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} = \cancel{\text{paternal}_j[\text{locus}]} + d_{i,j}[\text{locus}] \\
 + \cancel{\text{paternal}_j[\text{locus}]} + w_j[\text{locus}] \cdot \mathbf{inheritance}_{j,k} = \text{paternal}_k[\text{locus}] + d_{j,k}[\text{locus}] \\
 \hline
 \text{paternal}_i[\text{locus}] + \text{paternal}_k[\text{locus}] + \sum d = w_i[\text{locus}] \cdot \mathbf{inheritance}_{i,j} + w_j[\text{locus}] \cdot \mathbf{inheritance}_{j,k}
 \end{array}$$



**inheritance** will be absent in the equation if  $w = 0$ . We choose the path along with  $w = 1$  to obtain a constraint in the form  $\sum \mathbf{inheritance} = c$ .

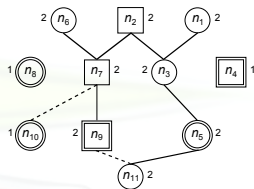
# Pedigree traversal: tree edges & non-tree edges

(a)

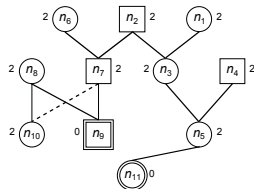


(c)

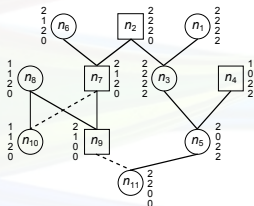
locus 1



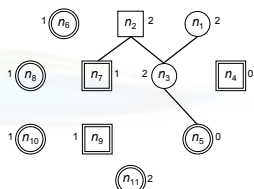
locus 3



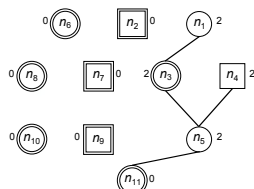
(b)



locus 2



locus 4

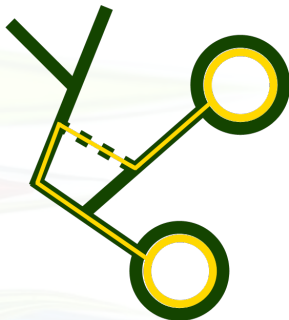


## Tree constraints



$$\sum \text{inheritance} = \text{paternal}_i[\text{locus}] + \text{paternal}_k[\text{locus}] + \sum d$$

## Path constraints



$$\sum \mathbf{inheritance} = \mathit{paternal}_i[\mathit{locus}] + \mathit{paternal}_k[\mathit{locus}] + \sum d$$

## Cycle constraints



$$\sum \text{inheritance} = \cancel{\text{paternal}_j[\text{locus}]} + \cancel{\text{paternal}_j[\text{locus}]} + \sum d$$

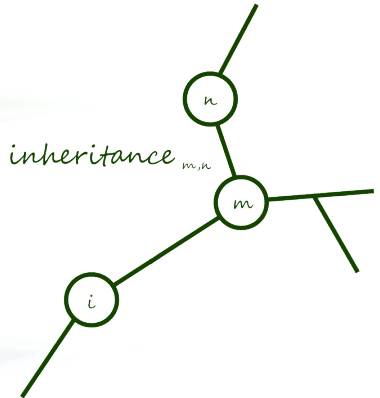
# A good property of tree constraints

There is only one path to connect two specific nodes on a tree.

That is, if there are two tree constraints,

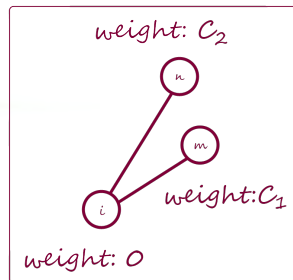
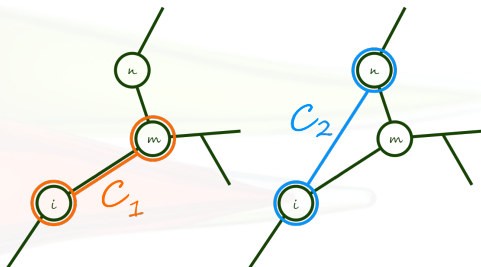
$$\text{tree 1: } \sum_i^m \mathbf{inheritance} = c_1$$

$$\text{tree 2: } \sum_i^n \mathbf{inheritance} = c_2$$





# Constraint graph



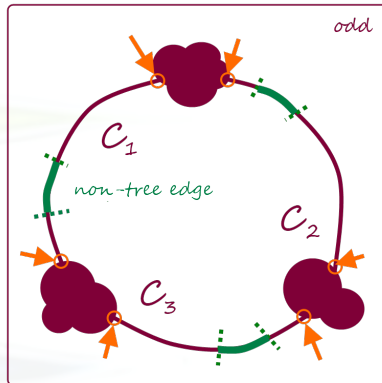
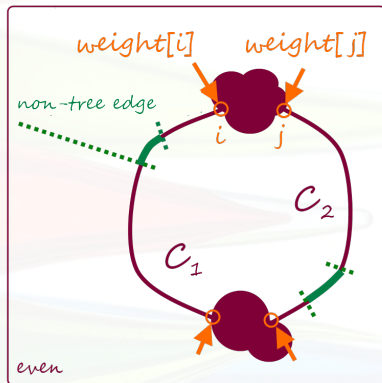
Then we can solve **inheritance** <sub>$m,n$</sub>  by

$$\text{inheritance}_{m,n} = \text{weight}[m] + \text{weight}[n] = c_1 + c_2$$

## Transformation: path constraints into tree constraints

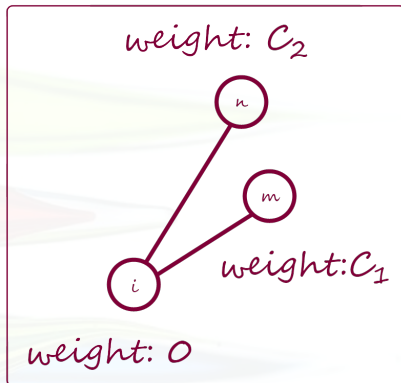


# Synthetic cycle constraints



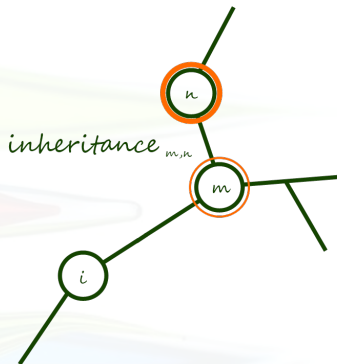
$$\sum \mathbf{inheritance}_{cycle} = \sum weight_{terminal} + \sum C_{path}$$

Use the constraint graph to solve **inheritance**



$$\mathbf{inheritance}_{m,n} = \mathit{weight}[m] + \mathit{weight}[n] = c_1 + c_2$$

## Use *inheritance* to solve **paternal**



$$paternal_n[locus] + w_i[locus] \cdot inheritance_{m,n} = \mathbf{paternal}_m[locus] + d_{m,n}[locus]$$

# Overview

- ▶ Input: genotype + pedigree
  1. Find the predetermined.
  2. Obtain tree, path, and cycle constraints. ( $\sum \mathbf{inheritance} = c$ )
  3. Use tree constraints to build constraint graph.
    - 3.1 Use cycle constraints to transform path constraints into tree constraints.
    - 3.2 Construct a synthetic cycle constraint if the cycle constraint is absent.
  4. Use the constraint graph to solve **inheritance**.
  5. Use *inheritance* to solve **paternal**.
- ▶ Output: haplotype

Thank  
you!