

# A data mishap

## Allele frequencies in sibships

Karl Broman

Biostatistics & Medical Informatics, UW–Madison

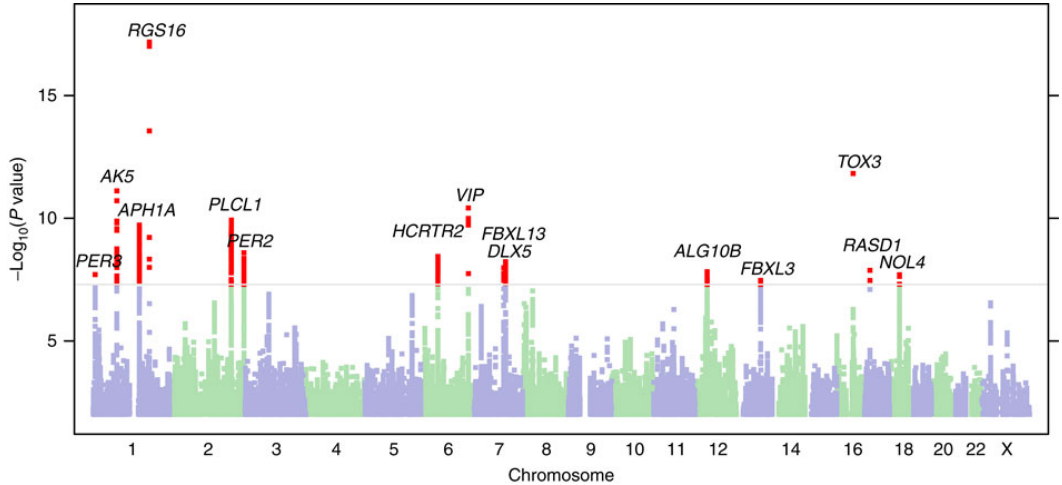
`kbroman.org`

`github.com/kbroman`

`@kwbroman`

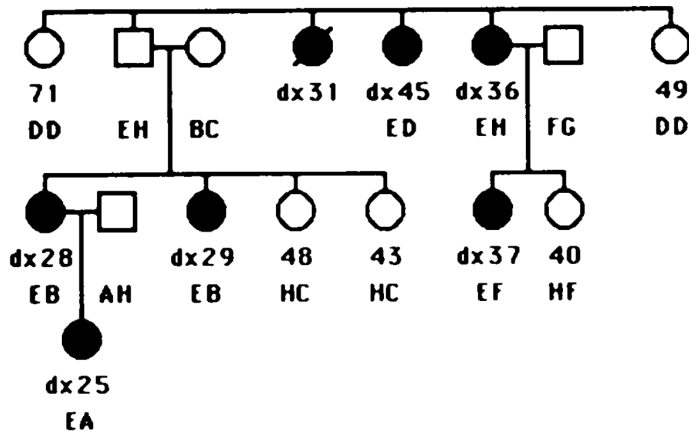
Slides: [kbroman.org/Talk\\_DataMishap](https://kbroman.org/Talk_DataMishap)

# GWAS for “morning person”



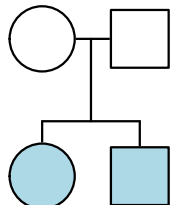
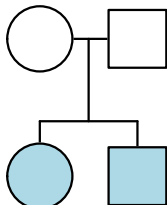
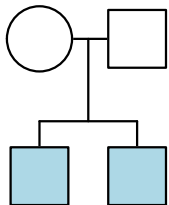
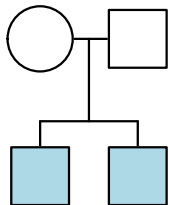
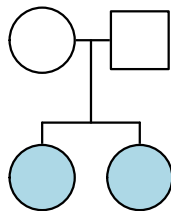
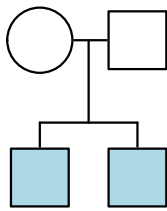
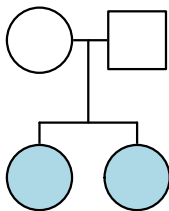
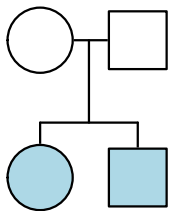
Hu et al (2016) [doi:10.1038/ncomms10448](https://doi.org/10.1038/ncomms10448)

# BRCA pedigree

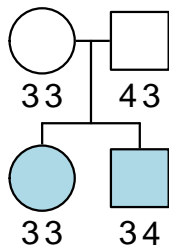
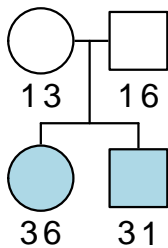
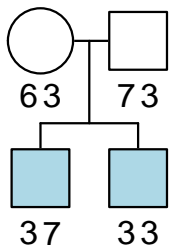
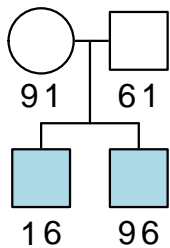
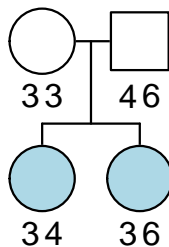
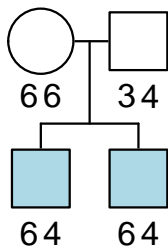
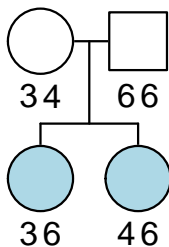
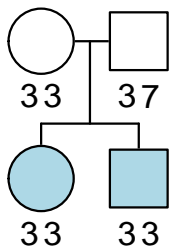


Hall et al (1990) [doi:10.1126/science.2270482](https://doi.org/10.1126/science.2270482)

## Affected sib pairs



## Affected sib pairs



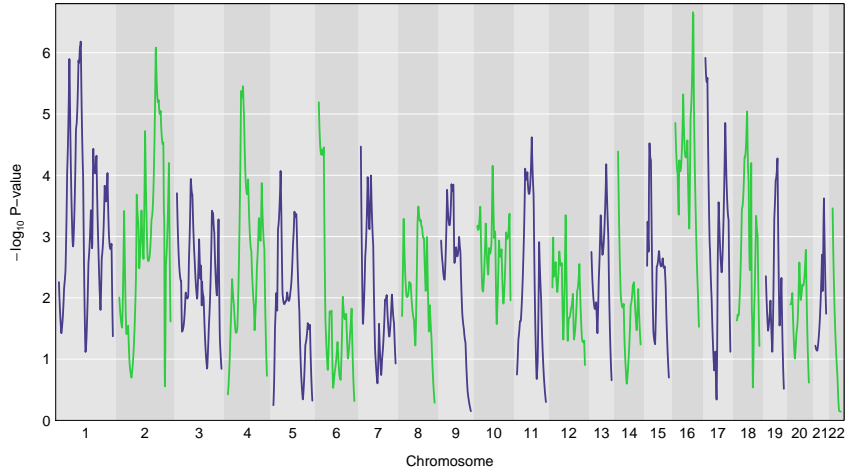
# IBS vs IBD

IBS = identical by **state**  
= same allele number

IBD = identical by **descent**  
= copies of the same ancestral allele

non-inbred sibs are IBD = 0, 1, 2  
with probability = 1/4, 1/2, 1/4

# Prostate cancer genome scan





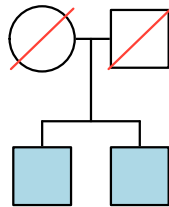
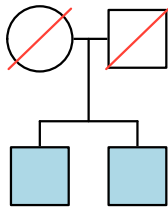
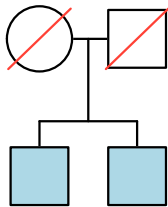
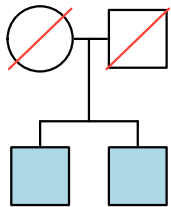
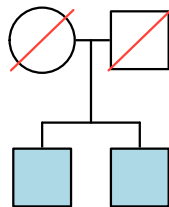
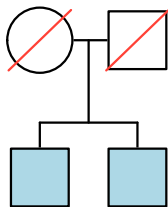
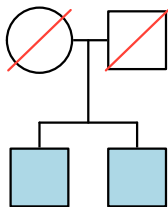
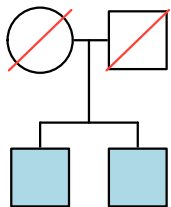
so happy



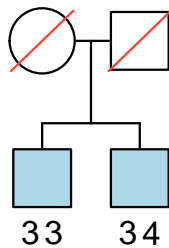
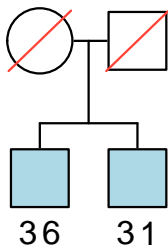
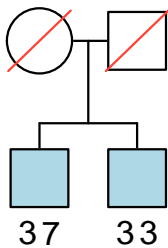
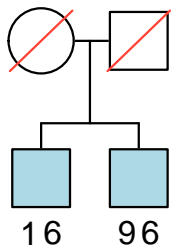
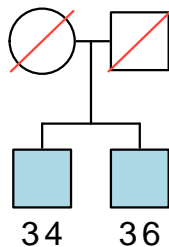
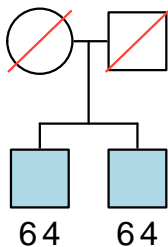
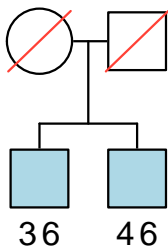
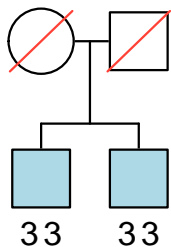
# Lesson

If it seems too good to be true,  
it probably is.

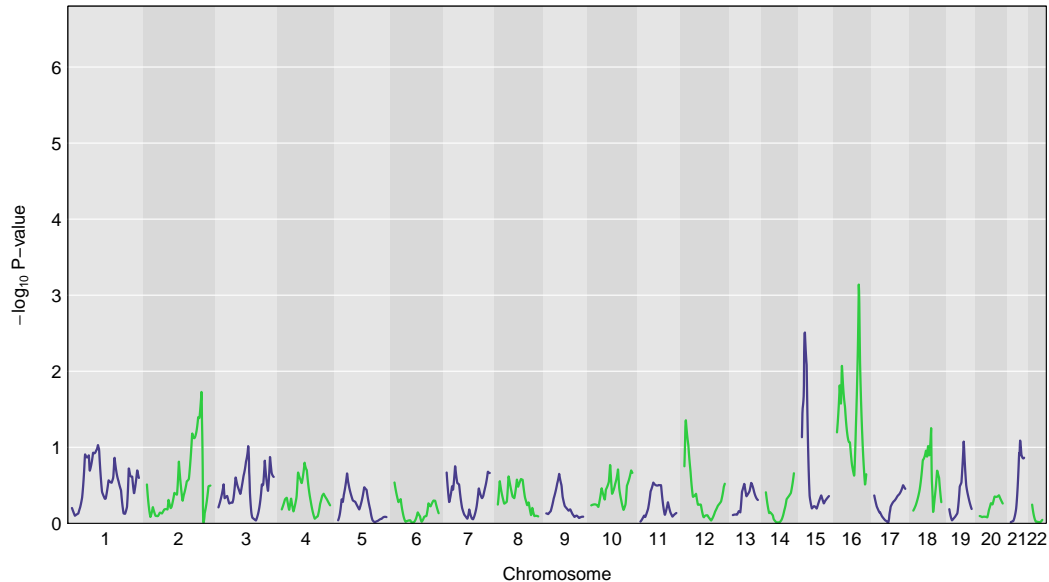
# Prostate cancer pairs



# Prostate cancer pairs



# Prostate cancer genome scan – corrected



# Estimating allele frequencies

Usually, you would use the **founders** in the pedigrees.  
(assumed unrelated)

What if you only have **sibships**?

# Estimating allele frequencies with sibpairs

Method 1: Use a random sibling from each

Method 2: Use everyone, ignoring relationships

# Estimating allele frequencies with sibpairs

Method 1: Use a random sibling from each

$$\text{var}(\hat{p}^{(1)}) = \frac{p(1-p)}{2n}$$

Method 2: Use everyone, ignoring relationships

$$\text{var}(\hat{p}^{(2)}) = ?$$

# Estimating allele frequencies with sibpairs

Method 1: Use a random sibling from each

$$\text{var}(\hat{p}^{(1)}) = \frac{p(1-p)}{2n}$$

Method 2: Use everyone, ignoring relationships

$$\text{var}(\hat{p}^{(2)}) = (3/4) \left( \frac{p(1-p)}{2n} \right)$$



# Estimating allele frequencies with sibpairs

Method 1: Use a random sibling from each

$$\text{var}(\hat{p}^{(1)}) = \frac{p(1-p)}{2n}$$

Method 2: Use everyone, ignoring relationships

$$\text{var}(\hat{p}^{(2)}) = (3/4) \left( \frac{p(1-p)}{2n} \right)$$

relative efficiency =  $4/3 = 1.33$   
(best possible = 1.5)

# My favorite equations

$$E(X) = E[E(X|Z)]$$

$$\text{var}(X) = E[\text{var}(X|Z)] + \text{var}[E(X|Z)]$$

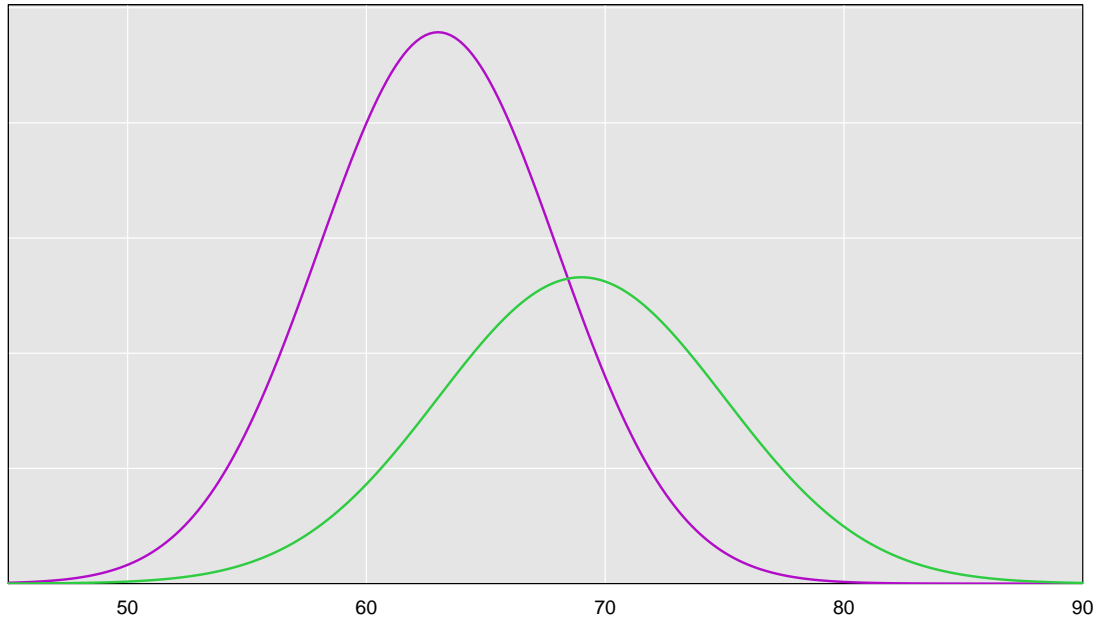
$$\text{cov}(X, Y) = E[\text{cov}(X, Y|Z)] + \text{cov}[E(X|Z), E(Y|Z)]$$

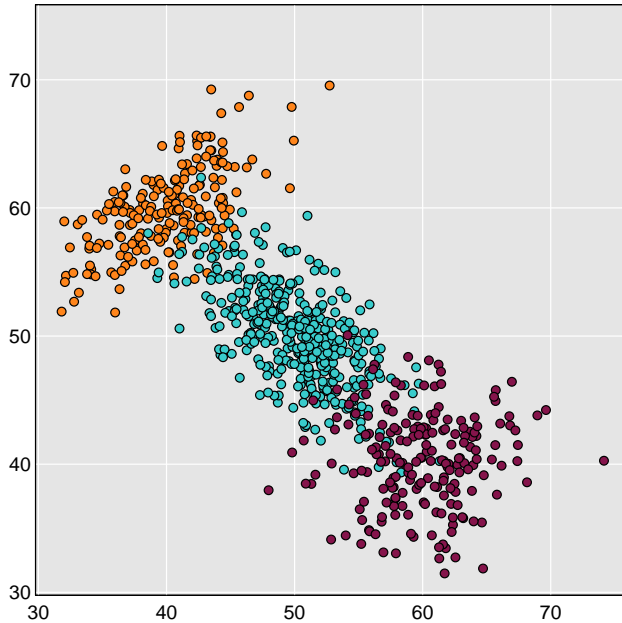
Everything is a mixture

## Another fave

$$\text{cov}(X, aY + bZ) = a \text{cov}(X, Y) + b \text{cov}(X, Z)$$

$$\begin{aligned}\text{thus } \text{var}(X + Y) &= \text{cov}(X + Y, X + Y) \\ &= \text{cov}(X + Y, X) + \text{cov}(X + Y, Y) \\ &= \text{cov}(X, X) + \text{cov}(X, Y) + \text{cov}(Y, X) + \text{cov}(Y, Y) \\ &= \text{var}(X) + \text{var}(Y) + 2 \text{cov}(X, Y)\end{aligned}$$





## Back to that SE

Let  $X_i$  = number of 1 alleles in sib  $i$ .

We want  $\text{var}(X_1 + X_2)$

$$\begin{aligned}\text{And so really } \text{cov}(X_1, X_2) &= E[\text{cov}(X_1, X_2 | \text{IBD})] + \text{cov}[E(X_1 | \text{IBD}), E(X_2 | \text{IBD})] \\ &= E[\text{cov}(X_1, X_2 | \text{IBD})] \\ &= \sum_{k=0}^2 \text{cov}(X_1, X_2 | \text{IBD} = k) \Pr(\text{IBD} = k)\end{aligned}$$

Also

$$\text{cov}(X, Y) = E(XY) - E(X)E(Y)$$

$$\text{cov}(X, Y|Z) = E(XY|Z) - E(X|Z)E(Y|Z)$$

**TABLE IV. Joint Distribution of the Numbers of 1 Alleles Carried by Two Individuals, Given the Number of Alleles They Share IBD**

IBD	$X_1, X_2$	$\Pr(X_1, X_2 \mid \text{IBD})$
0	0,0	$(1 - p)^4$
	0,1	$2p(1 - p)^3$
	1,0	$2p(1 - p)^3$
	1,1	$4p^2(1 - p)^2$
	0,2	$p^2(1 - p)^2$
	2,0	$p^2(1 - p)^2$
	1,2	$2p^3(1 - p)$
	2,1	$2p^3(1 - p)$
	2,2	$p^4$
1	0,0	$(1 - p)^3$
	0,1	$p(1 - p)^2$
	1,0	$p(1 - p)^2$
	1,1	$p(1 - p)$
	1,2	$p^2(1 - p)$
	2,1	$p^2(1 - p)$
	2,2	$p^3$
2	0,0	$(1 - p)^2$
	1,1	$2p(1 - p)$
	2,2	$p^2$



# Lessons

- ▶ Omitting data is usually bad
- ▶ Crudely ignoring correlations can be good  
You might even be able to figure out the SE

# Method 3

Account for relationships in the estimate

Missing data = IBD status for a sib pair at a marker

Use EM algorithm:

E step: estimate IBD status given allele frequencies

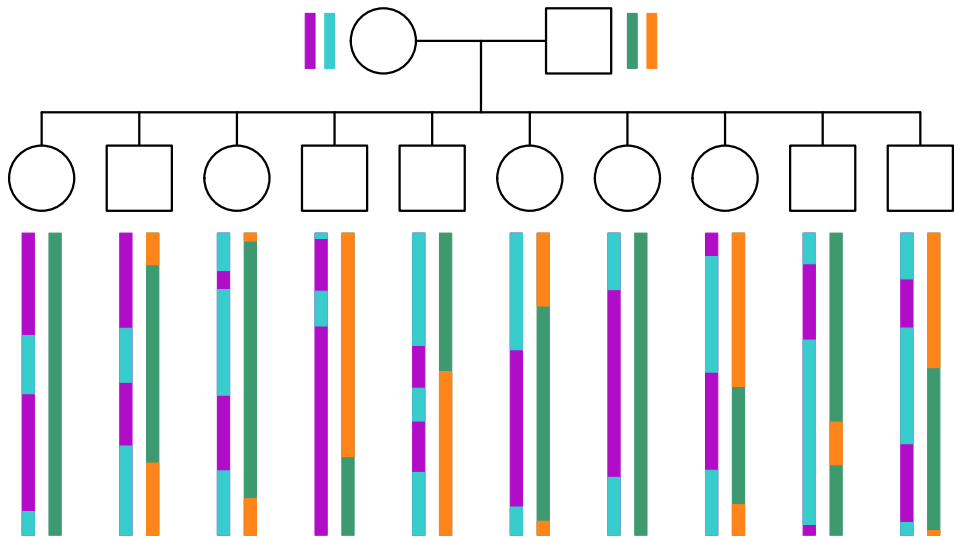
M step: estimate allele frequencies given IBD status

# Method 4

Make use of the multiple markers on a chromosome

- ▶ Markers along chromosome give improved info about IBD status
- ▶ Again, an EM algorithm:
  - Estimate IBD along chromosome given allele frequencies
  - Re-estimate allele frequencies using IBD information

# Siblings' chromosomes



## Average relative efficiency

Method	Allele frequency			
	0.05	0.10	0.15	0.20
1		1.00		
2		1.33		
3	1.46	1.45	1.44	1.43
4	1.48	1.46	1.45	1.44

## Method 3

<b>het</b>	<b>Allele frequency</b>			
	0.05	0.10	0.15	0.20
0.7	1.45	1.44	1.43	1.42
0.8	1.46	1.45	1.44	1.43
0.9	1.48	1.47	1.47	1.45

## Method 4

<b>d (cM)</b>	<b>Allele frequency</b>			
	0.05	0.10	0.15	0.20
0.1	1.50	1.49	1.48	1.48
1	1.49	1.48	1.47	1.46
5	1.48	1.46	1.44	1.44
10	1.47	1.45	1.43	1.42
(method 3)	1.46	1.45	1.44	1.43

# Summary

Method	Progr. time	CPU time	Rel. Eff.
1	2 min	1 msec	1.00
2	2 min	1 msec	1.33
3	1 morning	2 msec	1.45
4	1 afternoon	2.5 sec	1.46



# One last thing

- ▶ Turns out, I made a **mistake** in Method 3  
Mary Sara McPeck (U Chicago) spotted it
- ▶ Fixed problem, re-ran simulations, and...

# One last thing

- ▶ Turns out, I made a **mistake** in Method 3  
Mary Sara McPeck (U Chicago) spotted it
- ▶ Fixed problem, re-ran simulations, and...

the correct MLE was **worse** than my mistaken estimate