

# A data mishap

## Allele frequencies in sibships

Karl Broman

Biostatistics & Medical Informatics, UW–Madison

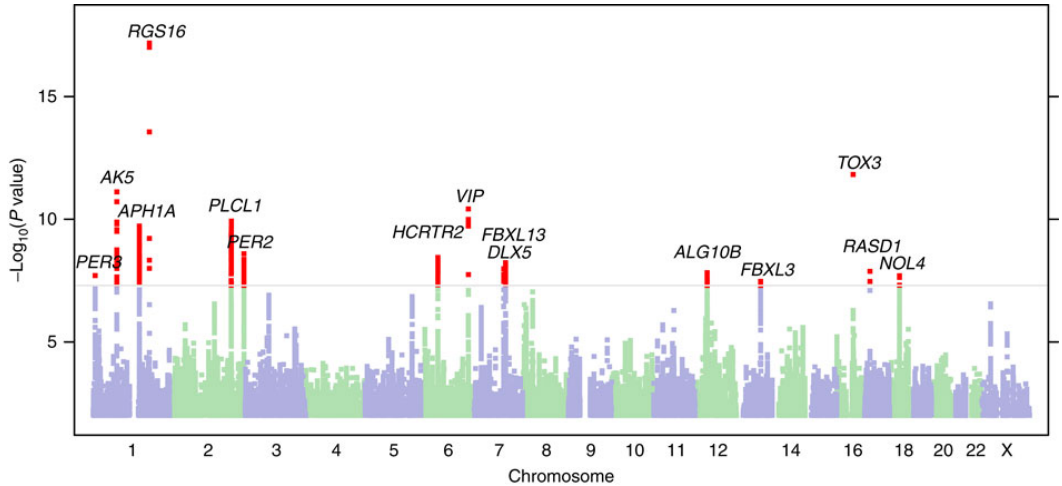
`kbroman.org`

`github.com/kbroman`

`@kwbroman`

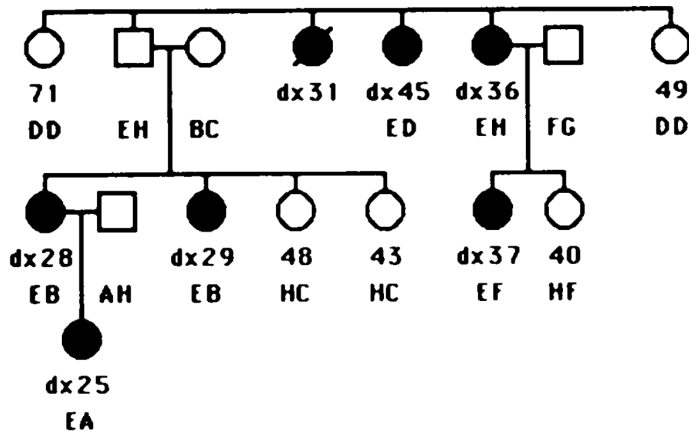
Slides: [kbroman.org/Talk\\_DataMishap](https://kbroman.org/Talk_DataMishap)

# GWAS for “morning person”



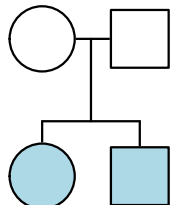
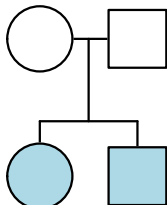
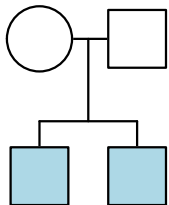
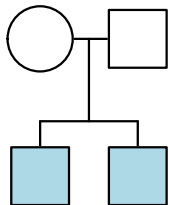
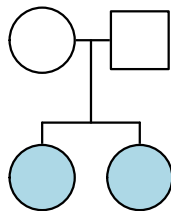
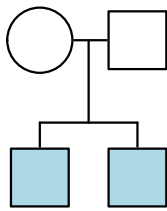
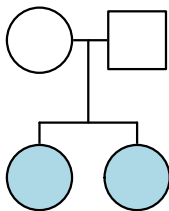
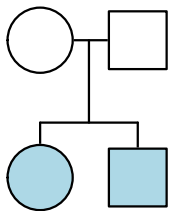
Hu et al (2016) [doi:10.1038/ncomms10448](https://doi.org/10.1038/ncomms10448)

# BRCA pedigree

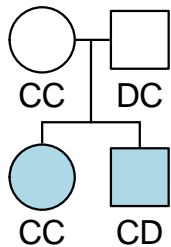
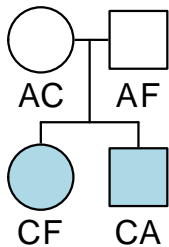
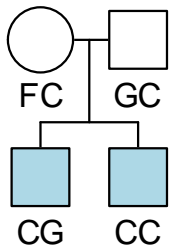
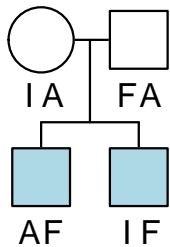
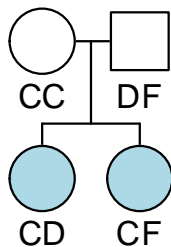
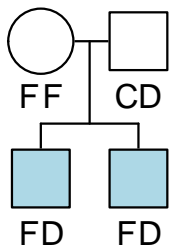
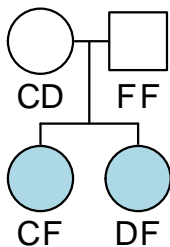
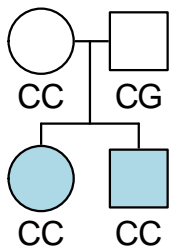


Hall et al (1990) [doi:10.1126/science.2270482](https://doi.org/10.1126/science.2270482)

## Affected sib pairs



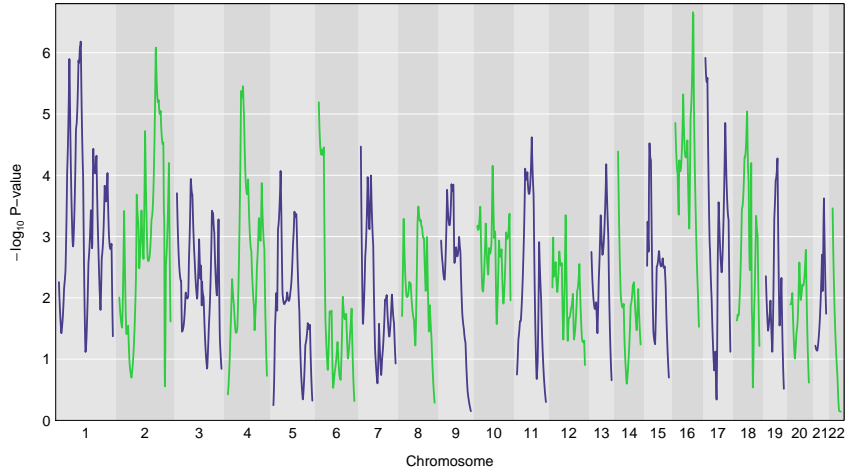
## Affected sib pairs



# Marshfield, Wisconsin



# Prostate cancer genome scan





so happy

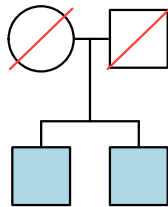
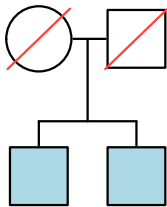
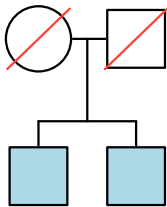
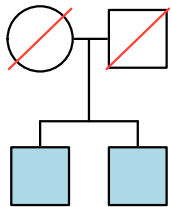
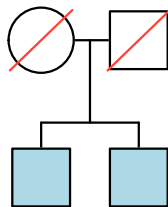
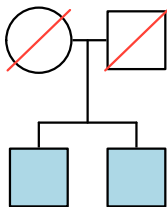
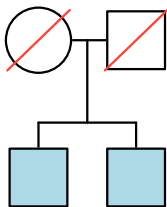
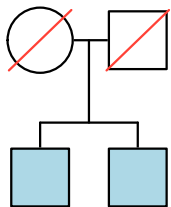




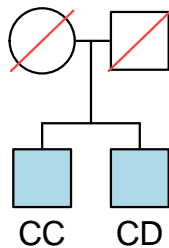
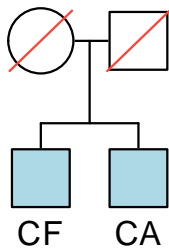
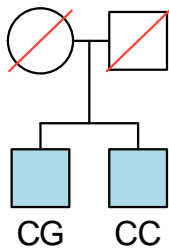
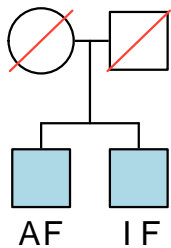
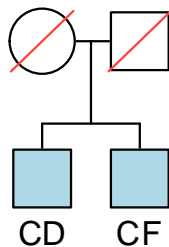
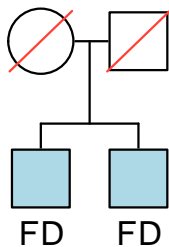
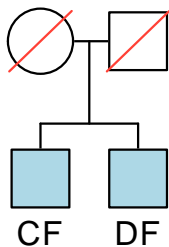
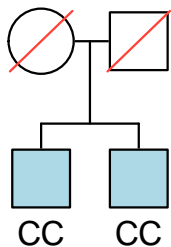
[bit.ly/faxpic](https://bit.ly/faxpic)

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

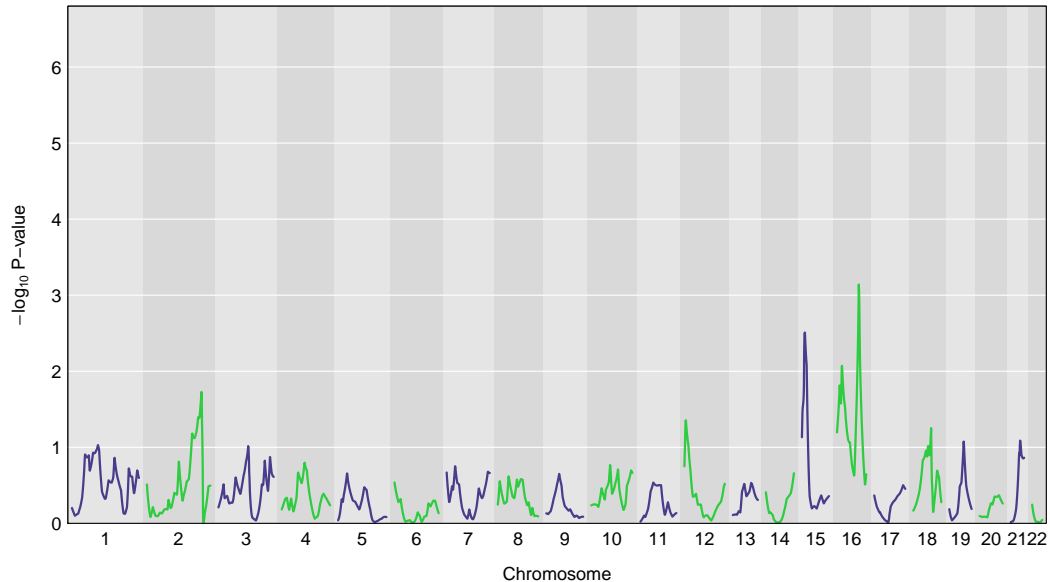
# Prostate cancer pairs



# Prostate cancer pairs



# Prostate cancer genome scan – corrected



# Estimation of Allele Frequencies With Data on Sibships

Karl W. Broman\*

*Department of Biostatistics, Johns Hopkins University, Baltimore, Maryland*

Allele frequencies are generally estimated with data on a set of unrelated individuals. In genetic studies of late-onset diseases, the founding individuals in pedigrees are often not available, and so one is confronted with the problem of estimating allele frequencies with data on related individuals. We focus on sibpairs and sibships, and compare the efficiency of four methods for estimating allele frequencies in this situation: (1) use the data for one individual from each sibship; (2) use the data for all individuals, ignoring their relationships; (3) use the data for all individuals, taking proper account of their relationships, considering a single marker at a time; and (4) use the data for all individuals, taking proper account of their relationships, considering a set of linked markers simultaneously. We derived the variance of estimator 2, and showed that the estimator is unbiased and provides substantial improvement over method 1. We used computer simulation to study the performance of methods 3 and 4, and showed that method 3 provides some improvement over method 2, while method 4 improves little on method 3. *Genet. Epidemiol.* 20:307–315, 2001. © 2001 Wiley-Liss, Inc.

# Estimation of Allele Frequencies With Data on Sibships

Karl W. Broman\*

Department of

Allele frequencies are estimated from data on sibships. The method is based on the fact that the probability of two alleles being identical by descent (IBD) is  $p_i/(1 + p_i)$ , where  $p_i$  is the frequency of allele  $i$ . We derive the maximum likelihood estimates of the allele frequencies and provide a simulation method for testing the hypothesis of Hardy-Weinberg equilibrium.

## Erratum: Broman KW. 2001. Estimation of Allele Frequencies With Data on Sibships. *Genet Epidemiol* 20:307–15.

Karl W. Broman\*

Department of Biostatistics, Johns Hopkins University, Baltimore Maryland

In the April 2001 issue of *Genetic Epidemiology*, in the article “Estimation of Allele Frequencies With Data on Sibships,” by Broman (20:307–15), there is an error on page 310, in the second paragraph under “Method 3: Accounting for Relationships.” The stated probabilities that an allele in the second sibling is not identical by descent (IBD) with one of the first sibling’s alleles, written as  $p_i/(1 + p_i)$ , are incorrect; we had missed two important cases. Let  $(g_{11}, g_{12})$  denote the two alleles of the genotype of the first sibling,  $(g_{21}, g_{22})$  denote the two alleles of the genotype of the second sibling, and  $\mathbf{g} = (g_{11}, g_{12}, g_{21}, g_{22})$ . Further, let  $\mathbf{pg}$  denote the genotypes for the two parents, and  $A$  denote the event “ $g_{21}$  is not IBD with  $g_{11}$  or  $g_{12}$ .” We seek  $\Pr(A|\mathbf{g})$ , which we calculate by conditioning on the parents’ genotypes, as follows:

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

[kbroman.org](http://kbroman.org)

[github.com/kbroman](https://github.com/kbroman)

@kwbroman