allele frequencies in sibships a data mishap

Karl Broman

Biostatistics & Medical Informatics, UW-Madison

kbroman.org github.com/kbroman @kwbroman Slides: kbroman.org/Talk_DataMishap

These are slides for a 5-min talk about a data mishap, for a community night (https://datamishapsnight.com/) organized by Caitlin Hudon (@beeonaposy) and Laura Ellis (@LittleMissData).



Genome-wide association studies (GWAS) have been a revolution in human genetics. This figure is from a study of 23andMe participants who were asked whether they're a morning person. This binary trait was associated with genotype at markers across the genome, immediately showing genes associated with the trait.





In-between, there was a period where we thought we could find disease genes by gathering a moderate number of affected sibling pairs. You look for regions where affected sibpairs had more similar genotypes than you would expect by chance.



In 1998 I was a postdoc in a genetics lab in Marshfield, WI (2 1/2 hours drive north of Madison). My advisor hooked me up with an affected sibpair study on prostate cancer. I did the initial data cleaning and a basic analysis, hoping to wow the famous people involved with my prowess.



have expected, and on many more chromosomes than I would have expected.



It was so awesome.



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

But as soon as I sent that fax, I was like, "Huh. Those results seem too good to be true."

It turns out that I'd messed up the allele frequencies and so the results were all messed up.



In this prostate cancer study, the affected sibpairs are all old, and there's essentially no data on the parents. In this case, our method for determining sharing is particularly sensitive to the allele frequencies.

It's not obvious how to estimate the allele frequencies, but also the simple approach I took had a bug that really through things off.



The unusually strong results I got were entirely due to a mistake in the code that estimated the allele frequencies. If I use more reasonable estimates, this is what I get. There's maybe evidence for a disease locus on chr 16 and possibly also 15, but the evidence isn't very strong.

And this is sort of what we'd expect given the size of this study. We're hoping to find some evidence of a disease gene, but we're not going to see the whole genome lighting up.



My collaborators were pretty nice about it. And I ended up writing a paper about the problem. That paper also had a major flaw, which is also interesting and instructive, but that's another story.

Slides: kbroman.org/Talk_DataMishap

kbroman.org

github.com/kbroman

@kwbroman