

I Multipoint Likelihoods & Large Families

- A Statistical problems
- B Computational problems

II CIs for trait location (denoted θ)

- A The standard estimator: $\hat{\theta}$
- B The Stewart & Peljto estimator: $\tilde{\theta}$
- C Variance estimators

III A Simulation Study

IV Approximate Importance Sampling

Statistical Problems:

- $\mathcal{H}_0 : \theta = \pm\infty$
- $L(\delta) \neq L(\theta)$
- *iid* no longer applies.
- Confidence Intervals (CIs)? Validity? Efficient?
- How do we know that the *chosen* SNPs are the right SNPs?

Likelihoods on Large Families w/ Dense SNPs

Computational Problems:

- A large latent space
- Sparse SNP panels *increase* error!
- Likelihoods involving dense SNPs are intractable

CI for θ with Dense SNPs: Notation

Let \mathbf{G} denote the obs'd dense SNP genotype data.

Ideally, we want to base inference on $\hat{\theta}(\mathbf{G})$, but this likelihood is intractable. **Why?** Because for certain pairs of SNPs there is linkage disequilibrium (LD), which means that

$$Pr(A_i - B_j) \neq Pr(A_i)Pr(B_j) \quad (1)$$

for any allele A and B of loci i and j , resp.

Note that LE (linkage equilibrium) implies equality in (1), and that likelihoods for sparse subsamples (denoted \mathbf{M}) are tractable provided that the SNPs are all in LE.

Thus, $\mathbf{M} \equiv \mathbf{M}(\mathbf{G}, \mathbf{S})$, where \mathbf{S} denotes a sparse SNP panel.

CI for θ with Dense SNPs: The Stewart-Peljto Estimator

In '10, we proposed $\tilde{\theta} \equiv \mathbf{E}\hat{\theta} \mid \mathbf{G}$, where $\mathbf{E}(\cdot)$ is taken wrt $Pr(\mathbf{S} \mid \mathbf{G})$.

In practice however, we estimate $\tilde{\theta}$ by $\frac{1}{k} \sum \hat{\theta}(\mathbf{G}, \mathbf{S}_j)$, from $\mathbf{S}_1, \dots, \mathbf{S}_k$ realizations of $\mathbf{S} \sim Pr(\mathbf{S} \mid \mathbf{G})$.

Note that (if you can compute it),

$$\begin{aligned} \text{Var}(\tilde{\theta}) &= \text{Var}[\mathbf{E} \hat{\theta} \mid \mathbf{G}] \\ &\leq \text{Var}[\mathbf{E} \hat{\theta} \mid \mathbf{G}] + \mathbf{E} \text{Var}(\hat{\theta} \mid \mathbf{G}) \\ &= \text{Var}[\hat{\theta}(\mathbf{G}, \mathbf{S})] \end{aligned}$$

CI for θ with Dense SNPs: Variance Estimators

For a large number of small families a nonparametric bootstrap approach is quite effective.

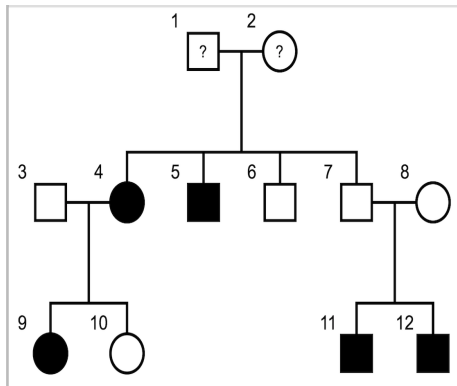
For a small number of large families

- (1) nonparametric bootstrap is no longer applicable,
- (2) a minus 1-LOD unit approach is approximate, at best
- (3) simulation of \mathbf{G} conditional on obs'd trait data is computationally infeasible due to LD.

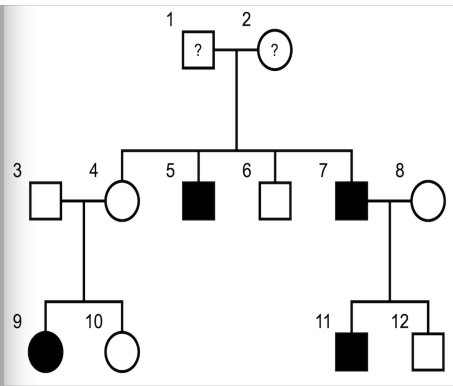
Also note that the $Var(\hat{\theta}) \approx \mathbf{E}[Var(\hat{\theta} | \mathbf{S})]$. Eq. (1)

A Simulation Study: Design

5-DOM pedigrees per replicate

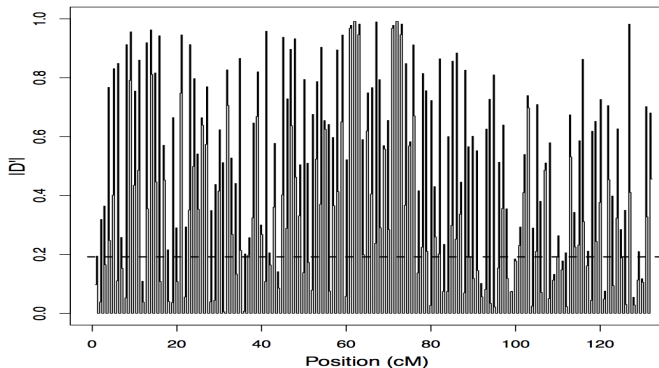


7-REC pedigrees per replicate



A Simulation Study: Design

LD structure of the 132 haplotype blocks:



A Simulation Study: Results

The conditional variance formula holds: $V\hat{\theta} = V\tilde{\theta} + \mathbf{E}V\hat{\theta} | \mathbf{G}$

Trait	$V\hat{\theta}$	$V\tilde{\theta}$	AMCE
DOM	47.27 (46.39)	38.43	8.84
REC	64.25 (63.13)	53.01	11.24

- CI lengths are reduced by 10%.
- $V\hat{\theta}$ in parentheses is computed by $\mathbf{E}V\hat{\theta} | \mathbf{S}$.

Approximate Importance Sampling (AIS)

Recall that the average Monte Carlo error (**AMCE**) is:

$$\begin{aligned} \mathbf{E}V\hat{\theta} \mid \mathbf{G} &= \sum [V\hat{\theta} \mid \mathbf{G}] Pr(\mathbf{G}) \\ &= \sum [V\hat{\theta} \mid \mathbf{G}] \frac{Pr^*(\mathbf{G})}{Pr(\mathbf{G})} Pr(\mathbf{G}) \\ &\approx \sum [V\hat{\theta} \mid \mathbf{G}] \frac{Pr^*(\mathbf{M}')}{Pr(\mathbf{M}')} Pr(\mathbf{G}), \end{aligned}$$

where $\mathbf{M}' \equiv \mathbf{M}(\mathbf{G}, \mathbf{S}')$ and \mathbf{S}' minimizes $V\hat{\theta} \mid \mathbf{S}$.

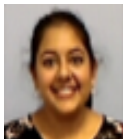
For REC: **AMCE** = 11.24 (cond. sim), and **AMCE** = 11.78 (AIS).

Acknowledgements

I would like to thank the Stewart Lab...

Post-docs: Meng Wang, Jane Cerise

Students: Jessica Fladen, Akanksha Malhotra, Komla Gnona, Valerie Hager, Esther Drill, and Anna Peljto.



For those who are interested in our software:

William.Stewart@nationwidechildrens.org

Thank You!

<http://u.osu.edu/stewart.1212/>

The Kong & Cox Likelihood ('97) at location j

$$\log L(\delta) = c(\mathbf{D}) + \sum_i \log[1 + \delta (f(\mathbf{D}_i) - \mu_i)/\sigma_i],$$

where $f(\mathbf{D}_i) = \mathbf{E}(S_i | \mathbf{D}_i, H_o : \delta = 0)$ and $i = 1, 2, \dots, n$ for n independent families in the data set.

Now define $Z^j \equiv \text{sgn}(\hat{\delta}) \sqrt{2 [\log L(\hat{\delta}) - \log L(0)]}$ for $j = 1, 2, \dots, m$ for m markers in a genome-wide cosegregation scan.

$$Z^j \rightarrow N(0, 1) \text{ as } n \rightarrow \infty$$