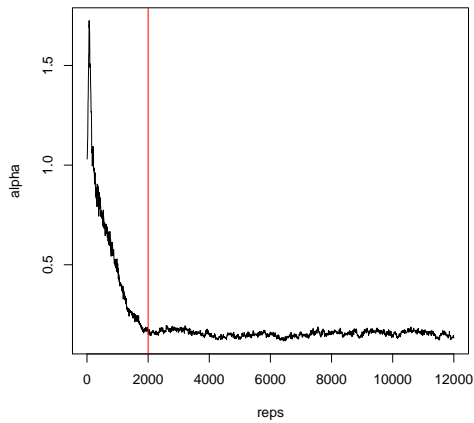


Supplementary Material for  
*Fast Model-Based Estimation of Ancestry in Unrelated  
Individuals*

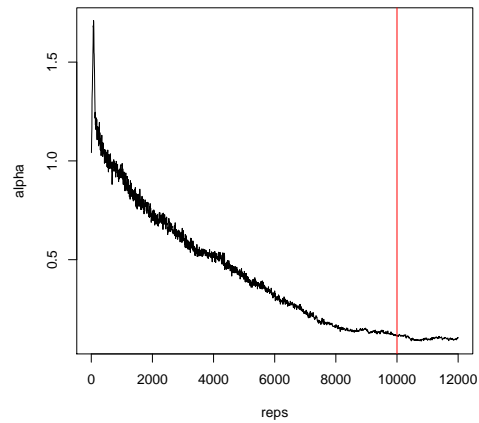
David H. Alexander

John Novembre

Kenneth Lange



(a)  $K = 2$



(b)  $K = 3$

Figure S1: Trajectories of STRUCTURE's estimate of  $\alpha$  for the IBD dataset with (a)  $K = 2$  and (b)  $K = 3$ . It appears that 2,000 burnin iterations are adequate for the IBD dataset with  $K = 2$ . For  $K = 3$  about 10,000 burnin iterations are required.

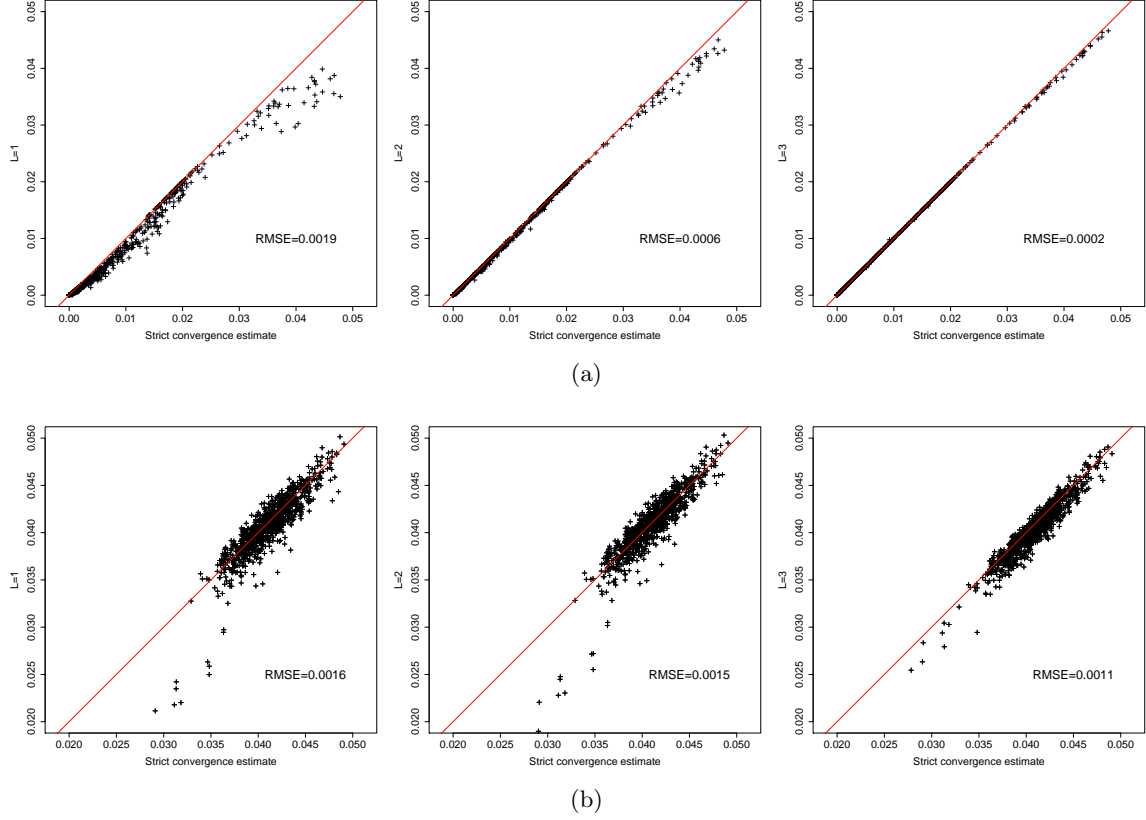


Figure S2: Evidence using the HapMap3 dataset with  $K = 3$  and the IBD dataset with  $K = 2$  suggests that the bootstrap standard errors  $\hat{\sigma}_L$  calculated with just a few accelerated block relaxation iterations per replicate ( $L = 1, 2$ , or  $3$ ) are adequate predictors of the bootstrap standard errors calculated using a stricter convergence requirement ( $\hat{\sigma}_{strict}$ ). This plot represents only the standard errors of the  $Q$  parameters. The diagonal line is  $y = x$ , and RMSE is the square root of the estimated mean square error of  $\hat{\sigma}_L$  as an estimator of  $\hat{\sigma}_{strict}$ .  $L = 3$  is the default in ADMIXTURE.

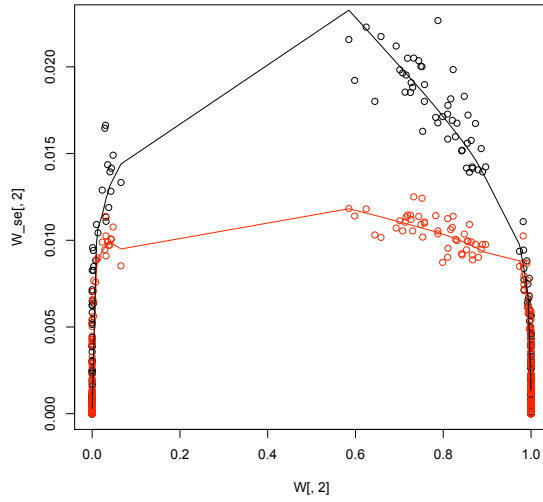


Figure S3: Our block bootstrap procedure produces more conservative standard error estimates than a bootstrap approach that resamples SNPs. The plot above illustrates this effect for the standard error estimates of  $\hat{q}_{i2}$ , the ancestry fraction appearing to correspond to African ancestry among the HapMap3 sample individuals. The points plotted are  $(\hat{q}_{i2}, \hat{\sigma}(\hat{q}_{i2}))$ . The black points correspond to the block bootstrap; the red points correspond to a bootstrap procedure that resamples SNPs (equivalent to a block bootstrap with block size  $h = 1$ ).

## 1 Further Discussion of Simulated Association Studies

Entries in our Table 4 are rounded to four places, giving the potentially confusing impression that correction with estimates from EIGENSTRAT or ADMIXTURE yields more a more powerful association test than correction with the “Ideal” correction using true ancestry values (for example, see the entry for Causal SNPs in Experiment II). In fact, viewing the results in higher precision reveals that ADMIXTURE and EIGENSTRAT are still slightly imperfect correctors, yielding Type I error that is slightly greater than the nominal level of .0001. For example, for the “Differentiated SNPs” in Experiment II, the average proportion of SNPs found significant using the ideal correction was  $.938 \times 10^{-4}$  while for ADMIXTURE and EIGENSTRAT the respective average proportions were  $1.24 \times 10^{-4}$  and  $1.23 \times 10^{-4}$ . The standard error on these average proportions is  $\sim 3 \times 10^{-6}$ , so it appears that EIGENSTRAT- and ADMIXTURE-corrected tests have Type I errors slightly, but significantly, higher than the ideally-corrected test. Thus it is not quite accurate to say that ADMIXTURE and EIGENSTRAT are more powerful correctors than the true ancestry, since the comparison is between tests of different Type I error rates.