

Genetic association models are robust to common population kinship estimation biases

Zhuoran Hou¹, Alejandro Ochoa^{1,2,*}

¹ Department of Biostatistics and Bioinformatics, Duke University, Durham, NC 27705, USA

² Duke Center for Statistical Genetics and Genomics, Duke University, Durham, NC 27705, USA

* Corresponding author: alejandro.ochoa@duke.edu

Abstract

Common genetic association studies for structured populations, including Principal Component Analysis (PCA) and Linear Mixed-effects Models (LMM), model the correlation structure between individuals using population kinship matrices, also known as Genetic Relatedness Matrices or “GRMs”. However, the most common kinship estimators can have severe biases that were only recently characterized. Here we characterize the effect of these kinship biases on genetic association. We employ a large simulated admixed family and genotypes from the 1000 Genomes Project, both with simulated traits, to evaluate a variety of kinship matrices (every bias type has two locus weight types, and their theoretical limits for the simulation). Remarkably, we find nearly equal association statistics and performance for kinship matrices of different bias types (when all other features are matched). These empirical observations lead us to hypothesize that these association tests are invariant to these kinship biases, which using linear algebra we prove holds exactly for LMM and approximately for PCA. Our constructive proof shows that the intercept and relatedness (PCs in PCA, random effect in LMM) coefficients compensate for the kinship bias, so the result extends to generalized linear models as long as those coefficients are present and are nuisance parameters. A corollary of our results is that association testing is also invariant to changing the ancestral population used to determine the kinship matrix. Overall, we find that existing association studies are robust to kinship estimation bias, and our calculations may help improve association methods by taking advantage

of this unexpected robustness, as well as help determine the effects of kinship bias in related problems.

1 Introduction

The goal of genetic association is to detect loci that are related to a specific trait, either causally or by proximity to causal loci. When applied to structured populations with admixed individuals, multiethnic cohorts, or close relatives, controlling for relatedness is crucial to avoid spurious associations and loss of power (Devlin and Roeder, 1999; Voight and Pritchard, 2005; Astle and Balding, 2009; Yao and Ochoa, 2022). The most popular association models for structured populations are Linear Mixed-effects Models (LMM) and Principal Component Analysis (PCA), which are closely related except LMM is capable of modeling high-dimensional structures whereas PCA is strictly a low-dimensional model (Astle and Balding, 2009; Hoffman, 2013; Yao and Ochoa, 2022).

Various association models, including both PCA and LMM, parametrize relatedness using kinship matrices, also known as Genetic Relatedness Matrices or “GRMs”. Kinship coefficients are well suited for this task since they model the covariance structure of genotypes (Malécot, 1948; Jacquard, 1970). Kinship is often encountered in family studies, where they reflect recent relatedness and can be calculated from pedigrees (Wright, 1922; Emik and Terrill, 1949; García-Cortés, 2015). However, as kinship is defined as a probability of identity by descent, it may also capture ancient population relatedness (Malécot, 1948; Astle and Balding, 2009), and common non-parametric kinship estimators from genotypes indeed include population structure in their estimates (Ochoa and Storey, 2021). In LMMs, the kinship matrix is an explicit parameter determining the random effect covariance structure (Xie et al., 1998; Yu et al., 2006; Aulchenko et al., 2007; Astle and Balding, 2009; Kang et al., 2008; Kang et al., 2010; Zhou and Stephens, 2012; Yang et al., 2014; Loh et al., 2015; Sul et al., 2018). In PCA, the principal components (PCs) are in practice the eigenvectors of an empirical genetic covariance matrix that is equivalent to the most common kinship estimator (Price et al., 2006; Astle and Balding, 2009; Hoffman, 2013; Yao and Ochoa, 2022).

Although several kinship estimators have been used with LMMs in the past, work from the last 15 years has converged on what we call the “standard” kinship estimator, which is the same

estimator used in PCA and other related models (Price et al., 2006; Astle and Balding, 2009; Rakovski and Stram, 2009; Thornton and McPeek, 2010; Yang et al., 2010; Yang et al., 2011; Zhou and Stephens, 2012; Speed et al., 2012; Yang et al., 2014; Speed and Balding, 2015; Loh et al., 2015; Wang et al., 2017; Sul et al., 2018). The impetus of our work is the recent characterization of a complex bias for this standard estimator, which varies for every pair of individuals (Weir and Goudet, 2017; Ochoa and Storey, 2021). This recent work also produced two new kinship estimators, which we are interested in characterizing in the context of association. The Weir-Goudet (WG) estimator constitutes a key improvement in that it has a uniformly downward bias (Weir and Goudet, 2017; Ochoa and Storey, 2021). Lastly, the popkin estimator is the only unbiased estimator under arbitrary relatedness (Ochoa and Storey, 2021). To the best of our knowledge, the new WG and popkin estimators have not been used in association studies before, but represent potential improvements over the use of the standard estimator for association.

One potential confounder when comparing the above kinship estimators is that the standard estimator upweights rare variants in a formulation previously called “mean-of-ratios” (MOR), whereas WG and popkin do not, instead following a “ratio-of-means” (ROM) estimation strategy (Bhatia et al., 2013; Ochoa and Storey, 2021). Recent work also formulated a ROM version of the standard estimator, which has a more predictable bias than the widely used MOR version (Ochoa and Storey, 2021). Following a locus weight formulation that allows the standard estimator to weigh loci in both ways (Wang et al., 2017), here we generalized the popkin and WG estimators to have both MOR and ROM versions as well, in order to test for the effect of estimator bias without confounding by locus weighing strategy.

In this work, we originally hypothesized that kinship estimation bias would affect association testing. We perform evaluations using an admixed family simulation (Yao and Ochoa, 2022) as well as real genotypes from the 1000 Genomes project (Consortium, 2010; 1000 Genomes Project Consortium et al., 2012; Fairley et al., 2020), in both cases with simulated traits in order to characterize type I error control and power using robust statistics. Surprisingly, we found that both LMM and PCA association are robust to kinship estimation bias to an extent that most association statistics are invariant to these biases. Lastly, we theoretically characterize the conditions under which

these kinship biases result in invariant association statistics, which encompass changing ancestral population in the kinship matrix too. Overall, we found that long-used association approaches are robust to the most common kinship estimation biases, and developed theoretical results that may help improve association and related approaches such as heritability estimation.

2 Methods

2.1 Genetic model

The following genetic model justifies the use of kinship matrices in association studies, and is the basis of all kinship estimation bias calculations that our theoretical work depends upon.

Suppose there are m biallelic loci and n diploid individuals. The genotype $x_{ij} \in \{0, 1, 2\}$ at a locus i of individual j is encoded as the number of reference alleles, for a preselected but otherwise arbitrary reference allele per locus. These genotypes can be treated as random variables structured according to relatedness. If T is the ancestral population on which allele frequencies are being conditioned, φ_{jk}^T is the kinship coefficient of two individuals j and k , and p_i^T is the ancestral allele frequency at locus i , then under the kinship model (Ochoa and Storey, 2021) the expectation and covariance are given by

$$\mathbb{E} [\mathbf{x}_i | T] = 2p_i^T \mathbf{1}, \quad \text{Cov} (\mathbf{x}_i | T) = 4p_i^T (1 - p_i^T) \boldsymbol{\Phi}^T,$$

where $\mathbf{x}_i = (x_{ij})$ is the length- n column vector of genotypes at locus i , $\boldsymbol{\Phi}^T = (\varphi_{jk}^T)$ is the $n \times n$ kinship matrix, and $\mathbf{1}$ is a length- n column vector of ones. Both $\boldsymbol{\Phi}^T$ and p_i^T are parameters that depend on the choice of ancestral population, for which the Most Recent Common Ancestor (MRCA) population is the most sensible choice (Ochoa and Storey, 2021). However, one of the results of this work is proof that the choice of ancestral population does not affect association testing.

2.2 Kinship estimators

Each subsection below corresponds to a kinship estimator bias type: Popkin is unbiased, while Standard and WG have different bias functions (defined shortly). Each estimator bias type has

two locus-weight versions, some introduced in this work, called *ratio-of-means* (ROM) and *mean-of-ratios* (MOR), a terminology that follows previous convention for these and related estimators (Bhatia et al., 2013; Ochoa and Storey, 2021). Only ROM estimators have closed-form limits. Below $\hat{p}_i^T = \frac{1}{2n} \mathbf{x}_i^\top \mathbf{1}$ is the standard ancestral allele frequency estimator, where the \top superscript denotes matrix transposition (do not confuse with ancestral population superscript T), and $\hat{\Phi}^{T,\text{name}} = (\hat{\varphi}_{jk}^{T,\text{name}})$ relates the scalar and matrix formulas of each named kinship estimator.

2.2.1 Popkin estimator

The popkin (population kinship) estimator (Ochoa and Storey, 2021), generalized here to include locus weights w_i , is given by

$$\hat{\varphi}_{jk}^{T,\text{popkin}} = 1 - \frac{A_{jk}}{\hat{A}_{\min}}, \quad A_{jk} = \frac{1}{m} \sum_{i=1}^m w_i ((x_{ij} - 1)(x_{ik} - 1) - 1), \quad (1)$$

where in this work $\hat{A}_{\min} = \min_{j \neq k} A_{jk}$, and w_i must be positive but need not add to 1. We consider two broad forms for this estimator. The original ROM estimator has $w_i = 1$ and has an unbiased almost sure limit as the number of loci m go to infinity,

$$\hat{\Phi}^{T,\text{popkin-ROM}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \Phi^T,$$

under the assumption that the true minimum kinship is zero. The MOR version, introduced here, upweights rare variants by using $w_i = (\hat{p}_i^T (1 - \hat{p}_i^T))^{-1}$; although it has no closed-form limit, it is approximately unbiased as well (Appendix A) and it is connected to the most common estimator, Standard MOR (Appendix B). The use of locus weights here is inspired by previous calculations relating the standard ROM and MOR estimators (Wang et al., 2017).

2.2.2 Standard estimator

The ROM and MOR versions of the standard kinship estimator are, respectively,

$$\hat{\varphi}_{jk}^{T,\text{std-ROM}} = \frac{\sum_{i=1}^m (x_{ij} - 2\hat{p}_i^T)(x_{ik} - 2\hat{p}_i^T)}{\sum_{i=1}^m 4\hat{p}_i^T(1 - \hat{p}_i^T)}, \quad (2)$$

$$\hat{\varphi}_{jk}^{T,\text{std-MOR}} = \frac{1}{m} \sum_{i=1}^m \frac{(x_{ij} - 2\hat{p}_i^T)(x_{ik} - 2\hat{p}_i^T)}{4\hat{p}_i^T(1 - \hat{p}_i^T)}. \quad (3)$$

The ROM estimator has a biased limit, which is a function of the true kinship matrix (Ochoa and Storey, 2021):

$$\hat{\Phi}^{T,\text{std-ROM}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} F^{\text{std}}(\Phi^T) = \frac{1}{1 - \bar{\varphi}^T} \left(\Phi^T + \bar{\varphi}^T \mathbf{J} - \varphi^T \mathbf{1}^\top - \mathbf{1} (\varphi^T)^\top \right), \quad (4)$$

where $\mathbf{J} = \mathbf{1}\mathbf{1}^\top$ is the $n \times n$ matrix of ones, $\varphi^T = \frac{1}{n}\Phi^T\mathbf{1}$ is a length- n vector of per-row mean kinship values, and $\bar{\varphi}^T = \frac{1}{n^2}\mathbf{1}^\top\Phi^T\mathbf{1}$ is the scalar overall mean kinship. The MOR estimator does not have closed-form limit, but it is well approximated by Eq. (4) in practice, especially when loci with small minor allele frequencies are excluded prior to calculating this estimate. In Appendix B we prove that the two standard estimators are functions of the corresponding popkin estimators, given by the bias function F^{std} :

$$\begin{aligned} \hat{\Phi}^{T,\text{std-ROM}} &= F^{\text{std}}(\hat{\Phi}^{T,\text{popkin-ROM}}), \\ \hat{\Phi}^{T,\text{std-MOR}} &= F^{\text{std}}(\hat{\Phi}^{T,\text{popkin-MOR}}). \end{aligned}$$

2.2.3 Weir-Goudet estimator

The ROM version of the Weir-Goudet (WG) kinship estimator is given by (Weir and Goudet, 2017; Ochoa and Storey, 2021)

$$\hat{\varphi}_{jk}^{T,\text{WG-ROM}} = 1 - \frac{A_{jk}}{\hat{A}_{\text{avg}}}, \quad \hat{A}_{\text{avg}} = \frac{2}{n(n-1)} \sum_{j=2}^n \sum_{k=1}^{j-1} A_{jk}, \quad (5)$$

where A_{jk} is as in Eq. (1). Its biased limit is also a function of the true kinship matrix:

$$\hat{\Phi}^{T,\text{WG-ROM}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} F^{\text{WG}}(\Phi^T) = \frac{1}{1 - \tilde{\varphi}^T} (\Phi^T - \tilde{\varphi}^T \mathbf{J}), \quad (6)$$

where $\tilde{\varphi}^T$ is the mean kinship excluding the matrix diagonal:

$$\tilde{\varphi}^T = \frac{2}{n(n-1)} \sum_{j=2}^n \sum_{k=1}^{j-1} \varphi_{jk}^T. \quad (7)$$

In Appendix C we prove that

$$0 \leq \tilde{\varphi}^T \leq \bar{\varphi}^T \leq \delta^T \leq 1,$$

where $\bar{\delta}^T = \frac{1}{n} \sum_{j=1}^n \varphi_{jj}^T$, and equalities are achieved if and only if all kinship values are equal. Since the WG-ROM estimator closely resembles the popkin estimator in Eq. (1), it follows more straightforwardly that they are related by the bias function F^{WG} , while WG-MOR is introduced here and defined by the below formula:

$$\begin{aligned} \hat{\Phi}^{T,\text{WG-ROM}} &= F^{\text{WG}}(\hat{\Phi}^{T,\text{popkin-ROM}}), \\ \hat{\Phi}^{T,\text{WG-MOR}} &= F^{\text{WG}}(\hat{\Phi}^{T,\text{popkin-MOR}}). \end{aligned}$$

2.3 Association models

LMM and PCA are closely-related association models (Astle and Balding, 2009; Hoffman, 2013; Yao and Ochoa, 2022):

$$\text{LMM: } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{s} + \boldsymbol{\epsilon}, \quad (8)$$

$$\mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma^2 \Phi^T), \quad (9)$$

$$\text{PCA: } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{U}_d\boldsymbol{\gamma}_d + \boldsymbol{\epsilon}, \quad (10)$$

$$\Phi^T = \mathbf{U}\Lambda\mathbf{U}^\top, \quad (11)$$

where \mathbf{y} is a length- n vector of continuous trait values, α is the intercept coefficient, β_i is the genetic effect (association) coefficient of locus i , \mathbf{s} is the (genetic) random effect, σ^2 is the random effect variance factor, \mathbf{U}_d is the $n \times d$ matrix of top- d eigenvectors of Φ^T (often referred to as “principal components” in genetics), γ_d is a length- d vector of coefficients for each eigenvector, $\epsilon \sim \text{Normal}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I})$ are random independent residuals, and \mathbf{I} is the $n \times n$ identity matrix. Furthermore, Eq. (11) is the complete eigendecomposition of Φ^T , where \mathbf{U} is the $n \times n$ matrix of eigenvectors, and Λ is the $n \times n$ diagonal matrix of eigenvalues. As \mathbf{s} and \mathbf{U}_d play analogous roles in modeling the effect of relatedness in LMM and PCA, respectively, we refer to them jointly as “relatedness” effects, and σ and γ_d as their respective coefficients.

2.4 Simulations

2.4.1 Admixed family genotype simulation

An admixed family was simulated following previous work (Yao and Ochoa, 2022), except here only $K = 3$ ancestries were simulated and $F_{ST} = 0.3$ for the admixed individuals, which more closely resembles the parameters of recently-admixed individuals such as Hispanics and African-Americans. Briefly, our admixture model first simulates $n = 1000$ founder individuals with the number of loci $m = 100,000$. Random ancestral allele frequencies p_i^T , subpopulation allele frequencies $p_i^{S_u}$, individual-specific allele frequencies π_{ij} , and genotypes x_{ij} are drawn from this hierarchical model:

$$\begin{aligned} p_i^T &\sim \text{Uniform}(0.01, 0.5), \\ p_i^{S_u} | p_i^T &\sim \text{Beta}\left(p_i^T \left(\frac{1}{f_{S_u}^T} - 1\right), (1 - p_i^T) \left(\frac{1}{f_{S_u}^T} - 1\right)\right), \\ \pi_{ij} &= \sum_{u=1}^K q_{ju} p_i^{S_u}, \\ x_{ij} | \pi_{ij} &\sim \text{Binomial}(2, \pi_{ij}), \end{aligned}$$

where this Beta is the Balding-Nichols distribution (Balding and Nichols, 1995) with mean p_i^T and variance $p_i^T (1 - p_i^T) f_{S_u}^T$. This is implemented in the R package `bnpssd`.

We also include family structure in the simulation. 20 generations are generated iteratively. To

preserve admixture structure mentioned above, individuals in the first generation ($n = 1000$) are ordered by 1D geography, locally unrelated and randomly assigned sex. From the next generation, individuals are paired iteratively: randomly choosing males from the pool and pairing them with the nearest available female with local kinship $< 1/4^3$ until no available males or females. Family sizes are drawn randomly ensuring every family has at least one child. Children are reordered by the average coordinates of their parents, their sex are assigned randomly, and their alleles are drawn from parents independently per locus. The simulation is implemented in the R package **simfam**.

2.4.2 Trait simulation algorithm

Given an $m \times n$ genotype matrix $\mathbf{X} = (\mathbf{x}_i^\top)$, traits are simulated from

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{X}^\top \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{Normal}(\mathbf{0}, (1 - h^2)\mathbf{I}).$$

Given a desired number of causal loci $m_1 = n/10$ and heritability $h^2 = 0.8$, the goal is to choose causal coefficients $\boldsymbol{\beta}$ and the intercept α that result in zero mean and the desired trait heritability. Here, we use the “fixed effect sizes” trait simulation model described in (Yao and Ochoa, 2022). Briefly, first m_1 causal loci are randomly selected, and for these steps only \mathbf{X} is subset to these loci and reindexed. For known p_i^T , causal coefficients are constructed as:

$$\beta_i = \sqrt{\frac{h^2}{2m_1 v_i^T}},$$

where $v_i^T = p_i^T (1 - p_i^T)$; for unknown p_i^T , v_i^T is replaced by the unbiased estimate $\hat{v}_i^T = \hat{p}_i^T (1 - \hat{p}_i^T) / (1 - \bar{\varphi}^T)$, where $\bar{\varphi}^T$ is the mean kinship estimated from **popkin**. Coefficients are made negative randomly with probability 0.5. For known p_i^T , we obtain the desired zero trait mean with $\alpha = -2(\mathbf{p}^T)^\top \boldsymbol{\beta}$, where here $\mathbf{p}^T = (p_i^T)$ contains causal loci only. When p_i^T are unknown, to avoid covariance distortions, the intercept coefficient is constructed as

$$\alpha = -2\hat{p}^T \mathbf{1}_{m_1}^\top \boldsymbol{\beta}, \quad \hat{p}^T = \frac{1}{m_1} \mathbf{1}_{m_1}^\top \hat{\mathbf{p}}^T,$$

where $\mathbf{1}_{m_1}$ is a length- m_1 column vector of ones.

2.5 Real genotype data processing

To evaluate different kinship estimators on a real dataset, we use the high-coverage NYGC version of the 1000 Genomes Project (Fairley et al., 2020), which were processed as before (Yao and Ochoa, 2022). Briefly, using `plink2` (Chang et al., 2015) we kept only autosomal biallelic SNP loci with filter “PASS”, LD-pruned with parameters “`--indep-pairwise 1000kb 0.3`” to remove loci that have a greater than 0.3 squared correlation coefficient with other loci within 1000kb, and lastly remove loci with MAF < 0.01. The resulting data has $m = 1,111,266$ loci and $n = 2,504$ individuals. Traits were simulated for this dataset with $m_1 = n/10 = 250$ causal loci.

2.6 Evaluation of performance

AUC_{PR} and $SRMSD_p$ are used to evaluate approaches as before (Yao and Ochoa, 2022). Briefly, $SRMSD_p$ (Signed Root Mean Square Deviation) is used to measure the difference between the observed null p-value quantiles and the expected uniform quantiles (p-values of continuous test statistics follow a uniform distribution under the null):

$$SRMSD_p = \text{sgn}(u_{\text{median}} - p_{\text{median}}) \sqrt{\frac{1}{m_0} \sum_{i=1}^{m_0} (u_i - p_{(i)})^2},$$

where $m_0 = m - m_1$ is the number of null (non-causal) loci, i indexes null loci only, $p_{(i)}$ is the i th ordered null p-value, $u_i = (i - 0.5)/m_0$ is its expectation, p_{median} is the median observed null p-value, $u_{\text{median}} = \frac{1}{2}$ is its expectation, and sgn is the sign function (1 if $u_{\text{median}} \geq p_{\text{median}}$, -1 otherwise). $SRMSD_p = 0$ corresponds to calibrated p-values, $SRMSD_p > 0$ indicate anti-conservative p-values, and $SRMSD_p < 0$ are conservative p-values.

AUC_{PR} (Area Under the Precision and Recall Curve) is a binary classification measure calculated from the total numbers of true positives (TP), false positives (FP) and false negatives (FN) at some

threshold or parameter t :

$$\text{Precision}(t) = \frac{\text{TP}(t)}{\text{TP}(t) + \text{FP}(t)},$$

$$\text{Recall}(t) = \frac{\text{TP}(t)}{\text{TP}(t) + \text{FN}(t)},$$

followed by calculating the area under the curve traced as t varies recall from zero to one. Higher AUC_{PR} is better, with best performance at $\text{AUC}_{\text{PR}} = 1$ for a perfect classifier, while worst performance at $\text{AUC}_{\text{PR}} = \frac{m_1}{m}$ (overall proportion of causal loci) is for random classifiers.

2.7 Software

Popkin estimates were calculated with the `popkin` R package. Standard kinship estimates were calculated with GCTA (version 1.93.2beta). All other estimators and limits were calculated using the `popkinsuppl` R package. PCs were calculated with the `eigen` function of R.

GCTA was used to run all LMM associations (Yang et al., 2011; Yang et al., 2014). We pass $2\Phi^T$ for all kinship matrices tested (the same scale as its own kinship estimate). PCA association is performed with `plink2` (Chang et al., 2015). We used $r = k - 1 = 2$ for the admixed family simulations, and $r = 10$ for 1000 Genomes.

3 Results

3.1 Empirical analysis using admixed family simulation

To quantify the effect of kinship estimation bias, we simulated genotypes and traits, and calculated association p-values using a factorial design that tests all kinship matrix (3 bias types, times two locus weight versions and one limit) and association model (PCA and LMM) combinations. We first simulated an admixed population with $K = 3$ ancestries, then simulated a 20-generation random pedigree from the admixed population as founders. This high-dimensional admixed family scenario yields a large difference in performance between PCA and LMM (Yao and Ochoa, 2022).

Kinship estimates and available limits on this simulation are shown in Fig. 1. The true kinship

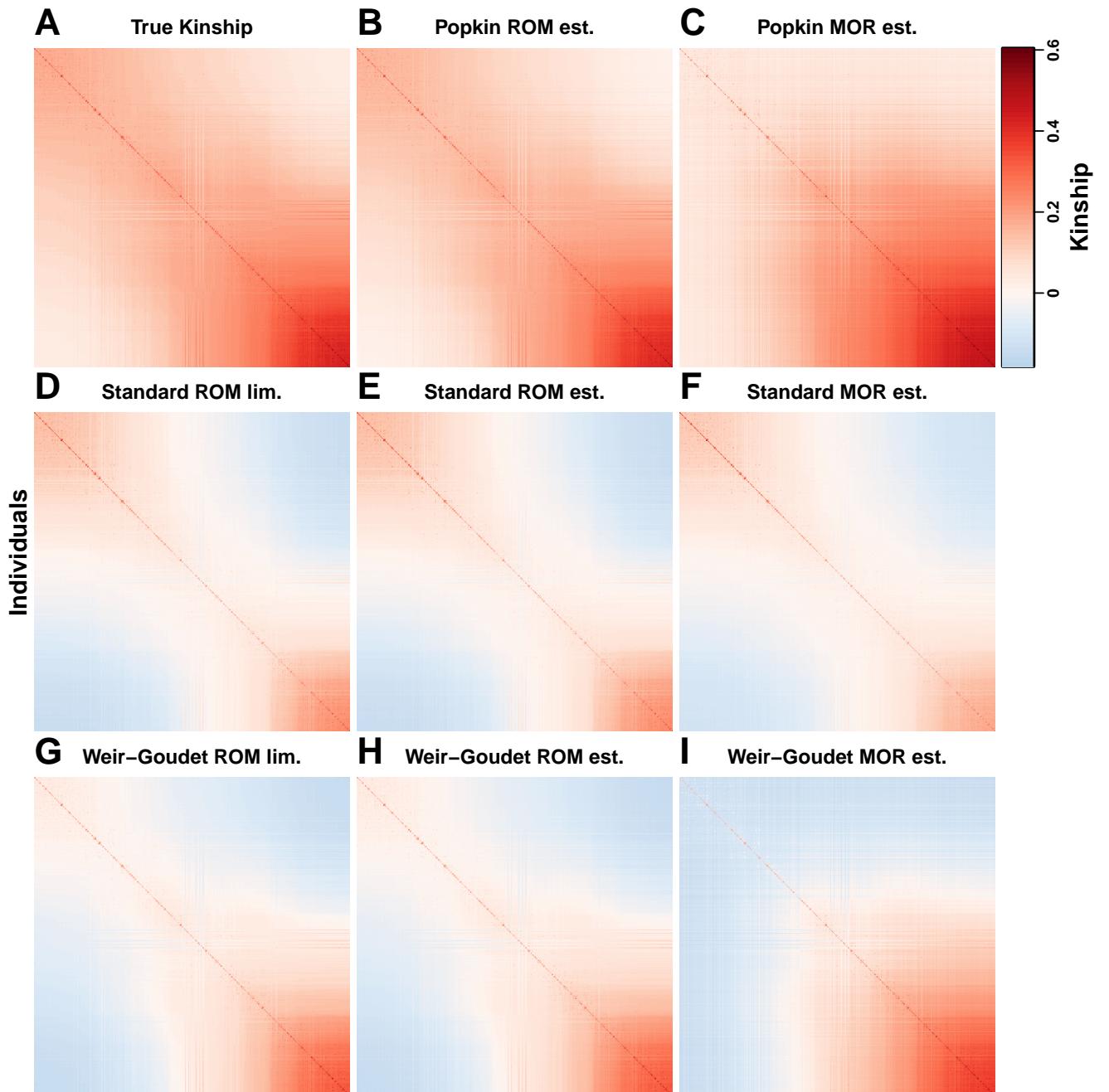


Figure 1: **Kinship estimates and limits on the admixed family simulation.** Each panel shows a kinship matrix as a heatmap, with each of the $n = 1000$ individuals along both x and y axes, color represents kinship: positive estimates in red, negative in blue. Diagonal contains inbreeding estimates. Each estimator bias type (Popkin, Standard, and Weir-Goudet; rows) has three matrices (columns): two locus-weight versions (ROM (ratio of means) and MOR (mean of ratios)) and limit of ROM.

matrix shows the family relatedness as high values concentrated near the diagonal and the ancestry-driven population structure as the broad patterns off-diagonal. Only Popkin ROM is unbiased, while popkin MOR has a slight upward bias that varies across the matrix (Fig. S1A). In contrast, the Standard and Weir-Goudet (WG) estimates have large downward biases overall, resulting in abundant negative values; for Standard these biases vary for every pair of individuals whereas for WG they are uniform.

We performed LMM and PCA association tests to determine how kinship biases affect association performance. Surprisingly, we found that kinship bias type does not have a discernible effect on association performance, as summarized by AUC_{PR} (a robust proxy for power; Fig. 2) and $SRMSD_p$ (measures null statistic calibration; Fig. S2). The largest differences in performance are explained

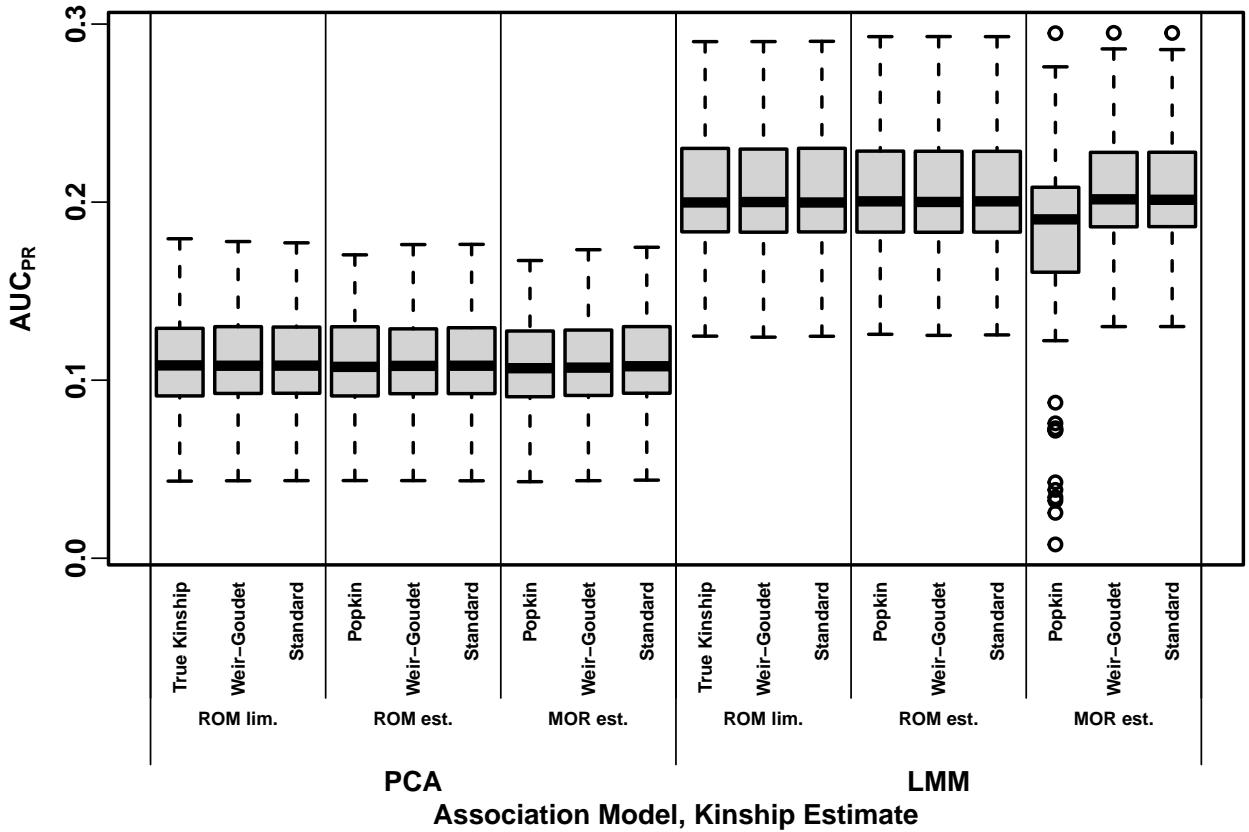


Figure 2: **Distributions of Area Under the Precision-Recall Curve (AUC_{PR}) on the admixed family simulation.** Higher AUC_{PR} is better performance. Results for 100 replicates (each a random genotype matrix and trait vector). Approaches cluster primarily by association model (LMM or PCA), and do not depend much at all on the bias type.

by the association model used (LMM vs PCA), as expected due to our use of a family simulation, where PCA performs poorly. Within association models, there are no clear differences between the performance of any of the kinship matrices, in fact many appear to have identical distributions (both statistics), the only clear exception being LMM popkin MOR, which has a few outlier replicates where performance was exceedingly poor.

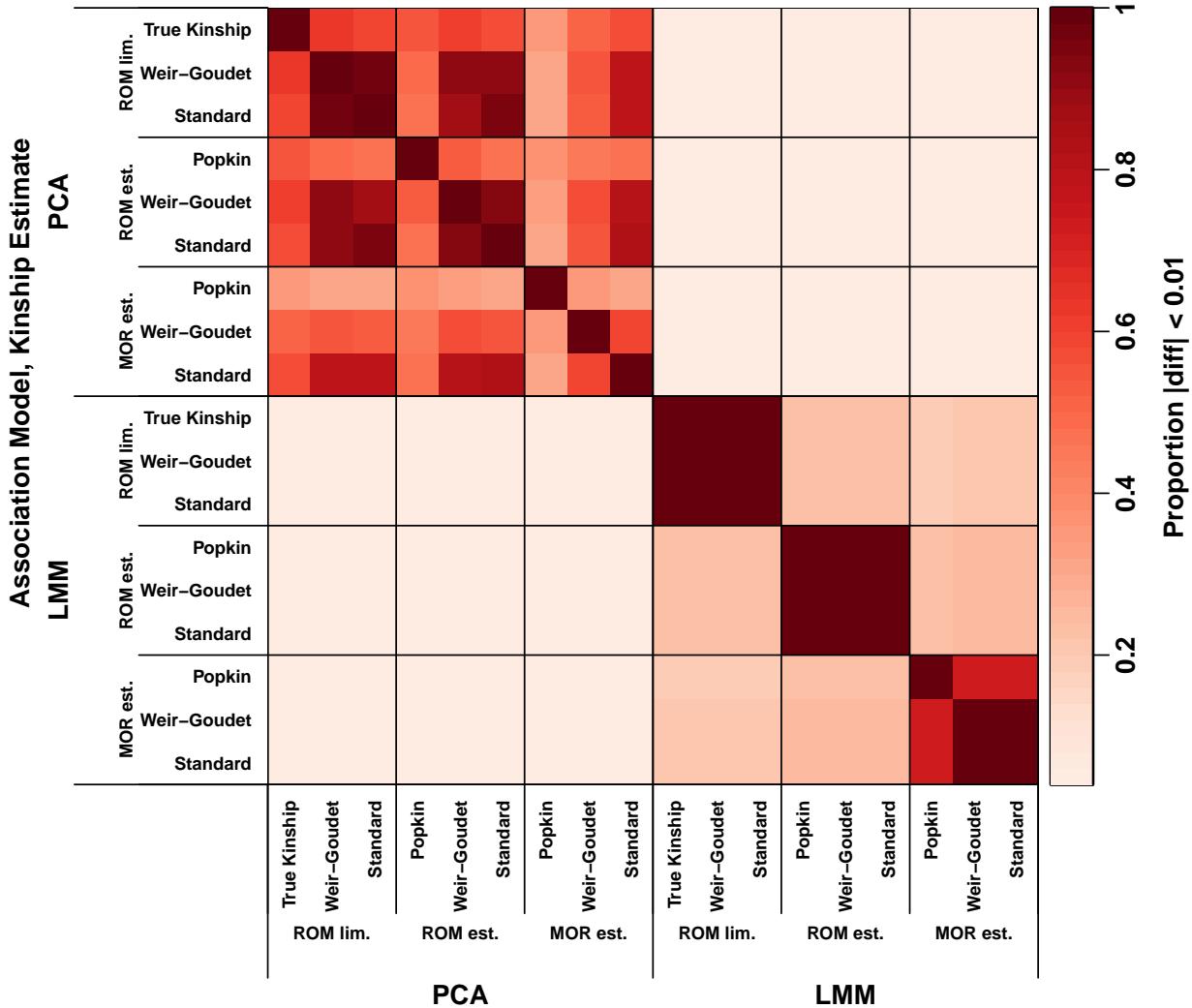


Figure 3: Approximate agreement between p-values on the admixed family simulation. The association p-value vector (one value per tested locus) produced by each combination of association model (LMM vs PCA) and kinship matrix (x and y axes) was used to compute proportions of loci with absolute differences under 0.01 (color). All 100 replicates were used. Methods of all bias types (matched for association model and locus weight type) have large proportions of nearly identical p-values.

To better characterize the nearly-identical performance distribution just observed, we next measured the agreement between individual association p-values. We calculated the proportion of loci between two methods with p-values within 0.01 of each other, which is an approximate measure of agreement, and found a remarkably high agreement between estimators of different bias types after matching association model and locus-weight version or limit (Fig. 3). This is in contrast to the low amounts of agreement across PCA and LMM statistics, and even across LMM statistics with different locus-weight or between those and the ROM limits. Minimum agreements tended to be higher across PCA methods, though here use of the true kinship or either popkin estimates resulted in more disagreements than between Standard and WG estimates or limits. Overall, sets of matched kinship matrices except for different bias types result in nearly identical association statistics.

3.2 Empirical analysis using 1000 Genomes

Now we repeat our analysis using the real genotypes of 1000 Genomes. Kinship estimates are shown in Fig. 4 (note real data has no true kinship or estimator limits). Popkin ROM estimates display an approximate nested block structure that arises from the tree relationships between subpopulations (Fig. 4A; trees were explicitly fit to this data in previous work (Yao and Ochoa, 2022)). However, popkin MOR estimates do not follow the nested blocks tree structure, since kinship between African and non-African populations is higher than kinship within African populations (Fig. 4B). Standard estimates have values closer to zero, and a different bias for each pair of individuals, resulting in higher relative kinship for African compared to non-African populations (Fig. 4C-D). Lastly, WG estimates are uniformly smaller than popkin's and attain large negative values (Fig. 4E-F).

Our association test conclusion are similar to our simulation study: AUC_{PR} and $SRMSD_p$ distributions are nearly identical for estimators of different bias types but same locus-weight version (ROM or MOR) and association model (Figs. S3 and 5). However, unlike the simulation, here the MOR estimates greatly outperform ROM estimates (LMM only), in terms of both AUC_{PR} and $SRMSD_p$. P-values are again nearly identical at a large proportion of loci between approaches with matched association model and locus-weight version (MOR or ROM), regardless of bias type (Fig. S4).

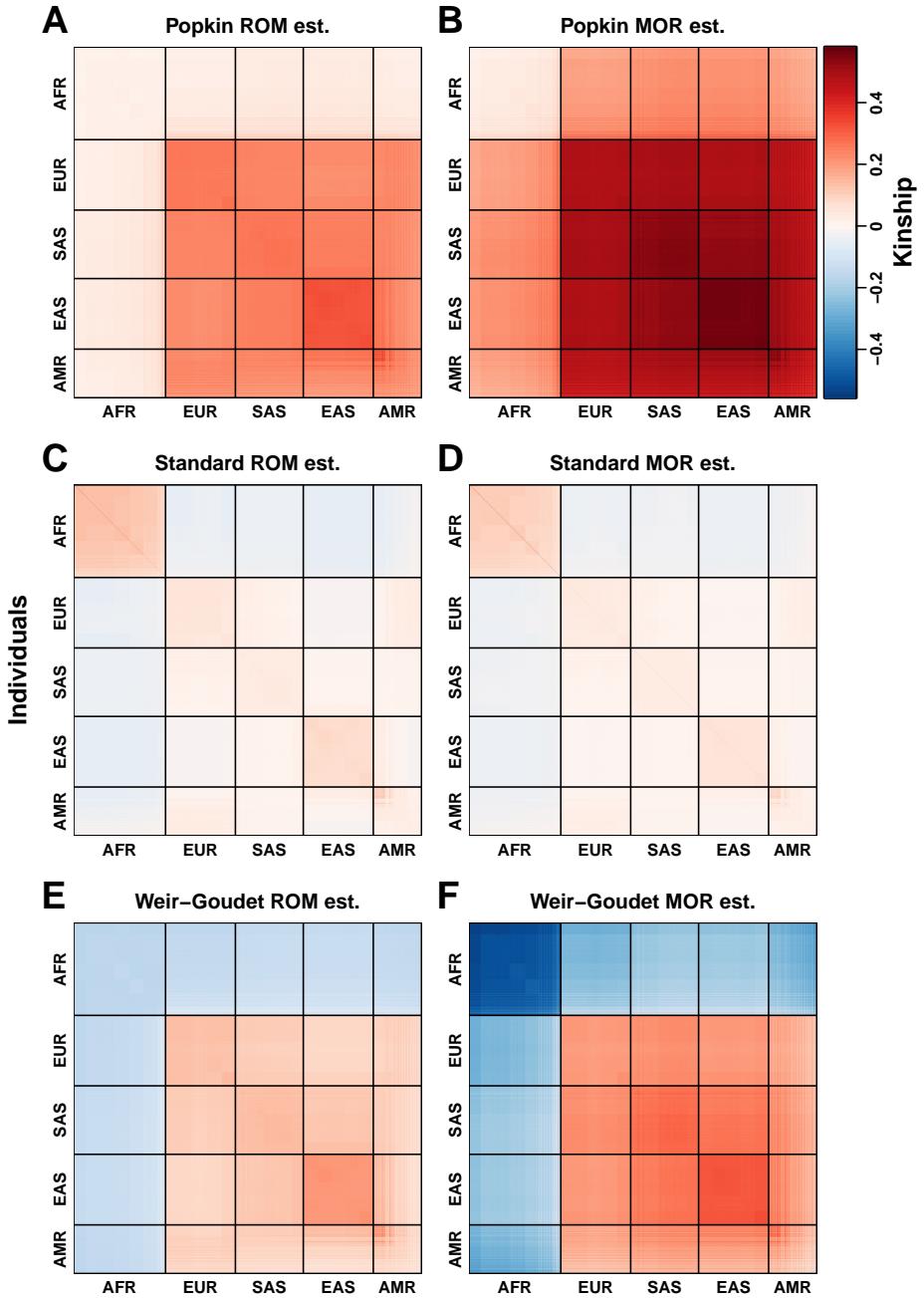


Figure 4: **Kinship estimates on 1000 Genomes.** Each panel represents a kinship matrix as a heatmap, as in Fig. 1. Superpopulation codes: AFR = African, EUR = European, SAS = South Asian, EAS = East Asian, AMR = Admixed Americans (Hispanics). Each estimator bias type (Popkin, Standard, and Weir-Goudet; rows) has two locus-weight versions (columns): ROM (ratio of means) and MOR (mean of ratios). In this visualization the upper range of all panels was capped to the 99 percentile of the diagonal (population inbreeding values) of the popkin MOR estimates.

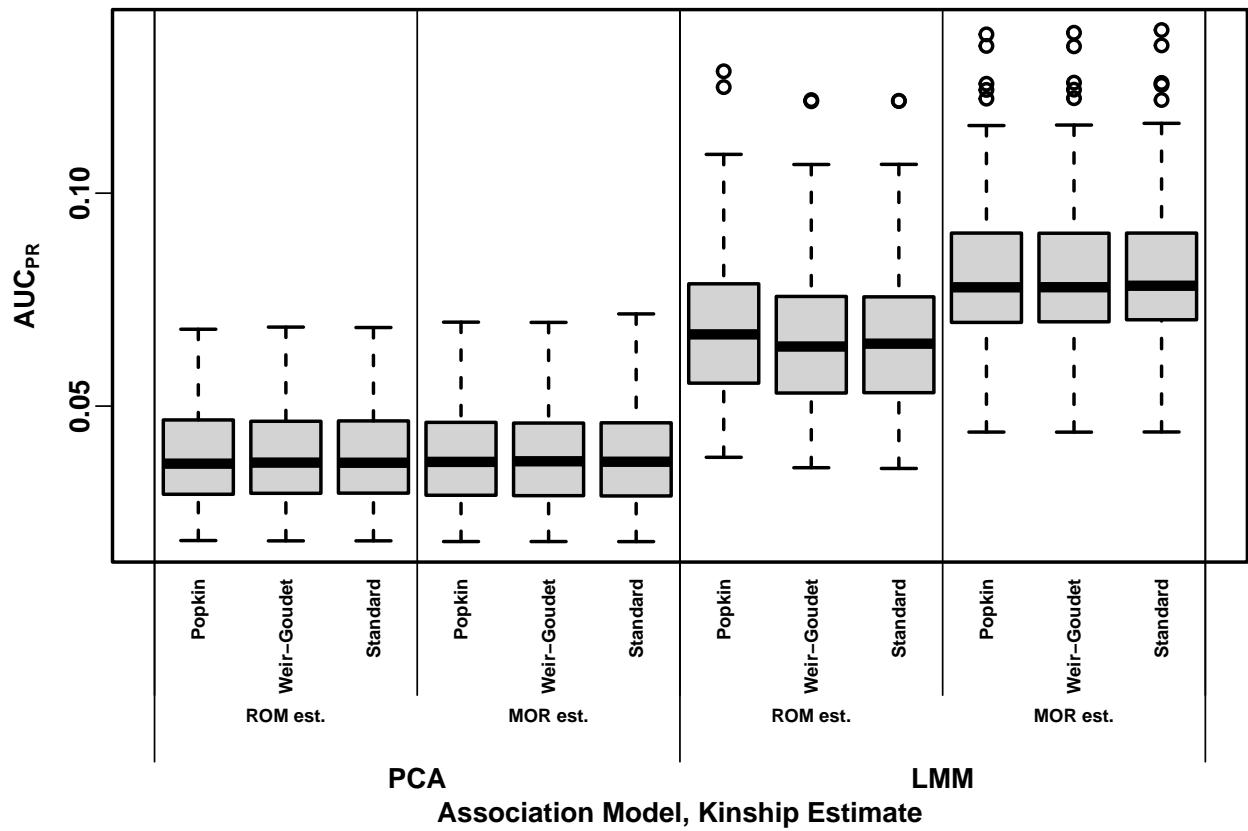


Figure 5: **Distributions of Area Under the Precision-Recall Curve (AUC_{PR}) on 1000 Genomes.** Higher AUC_{PR} is better performance. Results based on 100 simulated trait replicates (real genotype matrix is fixed). Approaches cluster primarily by association model (LMM or PCA) and locus-weight version (ROM or MOR), and do not depend much at all on the bias type.

3.3 Proof of association invariability to common kinship biases

Our empirical observations suggested that replacing a kinship matrix with either the Standard- or WG-biased version does not alter association statistics (with exceptions we attribute to numerical computation artifacts); here prove a general version of these facts mathematically. Our constructive proof shows that only a regression model with relatedness effects as covariates and an intercept is required, whose coefficients adapt to the bias, and no other coefficients change. This is fortunate, as the intercept and relatedness coefficients are nuisance parameters that usually go unreported, while the focal genetic association coefficient and its p-value are unchanged by these biases.

The most general form we identified of the bias function, mapping a kinship matrix to its bias-transformed version, and for which association invariability holds, is

$$\Phi^{T'} = F(\Phi^T) = \frac{1}{c} \mathbf{B} \Phi^T \mathbf{B}^\top, \quad \mathbf{B} = \mathbf{I} - \mathbf{1}\mathbf{b}^\top, \quad (12)$$

where c is any positive scalar and \mathbf{b} is any length- n vector. The key property that the linear operator \mathbf{B} must satisfy is that it shifts the input vector by the same scalar across its values, or

$$\mathbf{B}\mathbf{s} = \mathbf{s} - \mathbf{1}\eta, \quad (13)$$

where \mathbf{s} is any vector and the scalar $\eta = \mathbf{b}^\top \mathbf{s}$ is a function of the input vector. \mathbf{B} in Eq. (12) is the only form that results in Eq. (13).

The Standard bias function $F = F^{\text{std}}$ of Eq. (4) can be written as Eq. (12) with $c = 1 - \bar{\varphi}^T$ and $\mathbf{b} = \frac{1}{n}\mathbf{1}$, in which case \mathbf{B} equals the centering matrix. Further, the generalized standard kinship estimator studied in Ochoa and Storey (2021) instead has \mathbf{b} be the vector of individual weights used in the estimator, whose elements must sum to one, so $\mathbf{b}^\top \mathbf{1} = 1$. In all these cases \mathbf{B} and $\Phi^{T'}$ are singular transformations, since $\mathbf{B}\mathbf{1} = \mathbf{0}$ and $\mathbf{B}^\top \mathbf{b} = \mathbf{0}$.

The WG bias function $F = F^{\text{WG}}$ of Eq. (6) can be written as Eq. (12) with $c = 1 - \tilde{\varphi}^T$ and

$$\mathbf{b} = q \frac{(\boldsymbol{\Phi}^T)^{-1} \mathbf{1}}{\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1}}, \quad (14)$$

$$q = 1 \pm \sqrt{1 - \tilde{\varphi}^T \left(\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1} \right)}. \quad (15)$$

The derivation of this factorization is given in Appendix D. The determinant of the quadratic solution q would be non-negative if $\tilde{\varphi}^T$ satisfied $\tilde{\varphi}^T \leq 1 / (\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1})$. However, the actual $\tilde{\varphi}^T$ does not satisfy this inequality in any of our empirical cases, and in fact $1 / (\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1}) \leq \bar{\varphi}^T$ holds (proven in Appendix E; although $\tilde{\varphi}^T \leq \bar{\varphi}^T$ (Appendix C), in practice those two are very close while $1 / (\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1})$ is much smaller than both), so both values of b above are complex. This is a consequence of WG estimates being non-positive semidefinite, which we elaborate in the following sections. Nevertheless, the PCA approach as well as the GCTA algorithms for fitting variance components and estimating association coefficients work for non-positive semidefinite matrices without invoking complex numbers (following sections and Appendix F).

3.3.1 Proof for LMM case

Consider a random effect \mathbf{s} drawn using $\boldsymbol{\Phi}^T$, as given in Eq. (9). Using the affine transformation property of Multivariate Normal distributions (which despite its name also holds for singular linear transformations) and Eq. (12), it follows that

$$\mathbf{s}' = \mathbf{B}\mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma^{2\prime} \boldsymbol{\Phi}^{T'}) ,$$

where $\sigma^{2\prime} = c\sigma^2$. (This \mathbf{s}' has a degenerate distribution for Standard bias, since $\boldsymbol{\Phi}^{T'}$ is singular, but this is not problematic as long as $\mathbf{s}' + \boldsymbol{\epsilon}$ is non-degenerate, whose covariance $2\sigma^{2\prime} \boldsymbol{\Phi}^{T'} + \sigma_\epsilon^2 \mathbf{I}$ is invertible as long as $\sigma_\epsilon^2 \neq 0$.) Replacing $\mathbf{B}\mathbf{s}$ with the shift form in Eq. (13) shows that

$$\mathbf{s}' = \mathbf{s} - \mathbf{1}\eta$$

are equal in distribution. Therefore, after matching the variance coefficients σ^2 and $\sigma^{2\prime}$ as above, the random effect \mathbf{s}' of the biased kinship matrix differs from the random effect \mathbf{s} of the original kinship only by $\mathbf{1}\eta$, a difference compensated for by adjusting the intercept coefficient in Eq. (8): $\alpha' = \alpha + \eta$. No other regression coefficients, or the total residuals, change when Φ^T is replaced with $\Phi^{T\prime}$, including the association coefficient β_i that is the focus of the test.

The above results required positive semidefinite kinship matrices $\Phi^{T\prime}$. Nevertheless, for the non-positive semidefinite WG bias combined with the generalized least squares association algorithm, which is used by GCTA and other LMMs (Kang et al., 2008; Kang et al., 2010; Yang et al., 2014), we found a stronger result consistent with the above, namely that $\alpha' = \alpha$, or in other words $\eta = 0$ (Appendix F).

The LMM association p-value does not change in several common tests, including the F-test, since it only depends on the residuals and these do not change, as well as the likelihood ratio test, because although covariance determinants change, they cancel out in the ratio. The Wald test used by GCTA (Yang et al., 2014) is also invariant to these kinship biases given our empirical results in Figs. 3 and S4 and proven explicitly for WG bias in Appendix F. Lastly, we confirmed empirically that the Score test for the GCTA model is also invariant to these kinship biases (not shown). These arguments hold whether variance components are fit with maximum likelihood (ML) or restricted maximum likelihood (REML) (Kang et al., 2008; Kang et al., 2010; Yang et al., 2014), since multiplying the estimated genetic variance component σ^2 by c and adjusting the intercept compensates for the bias regardless of how $\sigma^2, \sigma_\epsilon^2$ are estimated.

3.3.2 Proof for PCA case

We present a proof for the PCA case that relies on an approximation that holds well in practice. Based on the PCA model of Eqs. (10) and (11), let \mathbf{U}_d be the top eigenvectors of Φ^T , and \mathbf{U}'_d those of $\Phi^{T\prime}$. The key approximation is that

$$\mathbf{U}'_d \approx \mathbf{B} \mathbf{U}_d, \quad (16)$$

which is not strictly equal (since $\mathbf{B}\mathbf{U}$ is not generally orthogonal, as eigenvectors must be), but we have found it to be a good approximation in practice. In this case the eigenvector coefficients need not change, $\boldsymbol{\gamma}'_d = \boldsymbol{\gamma}_d$, since the difference in scale of the kinship matrices (c in Eq. (12)) is absorbed by the eigenvalues, which are not used in the association model. Applying the shift of Eq. (13) shows that

$$\mathbf{U}'_d \boldsymbol{\gamma}'_d = \mathbf{B}\mathbf{U}_d \boldsymbol{\gamma}_d = \mathbf{U}_d \boldsymbol{\gamma}_d - \mathbf{1}\eta,$$

where $\eta = \mathbf{b}^\top \mathbf{U}_d \boldsymbol{\gamma}_d$ is a scalar. Therefore, the relatedness effects again differ only by $\mathbf{1}\eta$, which is compensated for by adjusting the intercept accordingly, so the association coefficient β_i and the residuals are the same in both cases. Note that the proof works whether there are small numbers of zero or negative eigenvalues in $\Phi^{T'}$ (non-positive semidefinite cases), as those rank last and are simply ignored. The observations from LMMs, for how p-values change depending on the type of test used, also hold for PCA.

We visualize the top PCs of our two datasets in Fig. 6 in order to assess the validity of the approximation in Eq. (16). The approximation is equivalent to each biased PC (Standard or Weir-Goudet) being shifted from the unbiased PC (Popkin), as described in Eq. (13). Fig. 6 indeed shows that PC1 is shifted by noticeable amounts in each of these cases, while PC2 is less shifted. However, a rotation of the PCs is also noticeable, particularly in the simulated data, and other large differences between MOR estimators, as expected since we know the approximation cannot be exact. One last detail worth noting is that PCs can change order upon the bias transformation, which we noticed in the admixed family simulation, where PC2 and PC3 from popkin (and true kinship) actually correspond to PC1 and PC2, respectively, in both Standard and Weir-Goudet estimates (both ROM and MOR), and were plotted as such. No PC reordering occurred in 1000 Genomes. Overall, while the approximation of Eq. (16) can be weakened to merely require that the biased PCs plus intercept span the same subspace of the unbiased PCs plus intercept, the approximate PC shifts better explain intuitively why the result for LMM is also observed for PCA association.

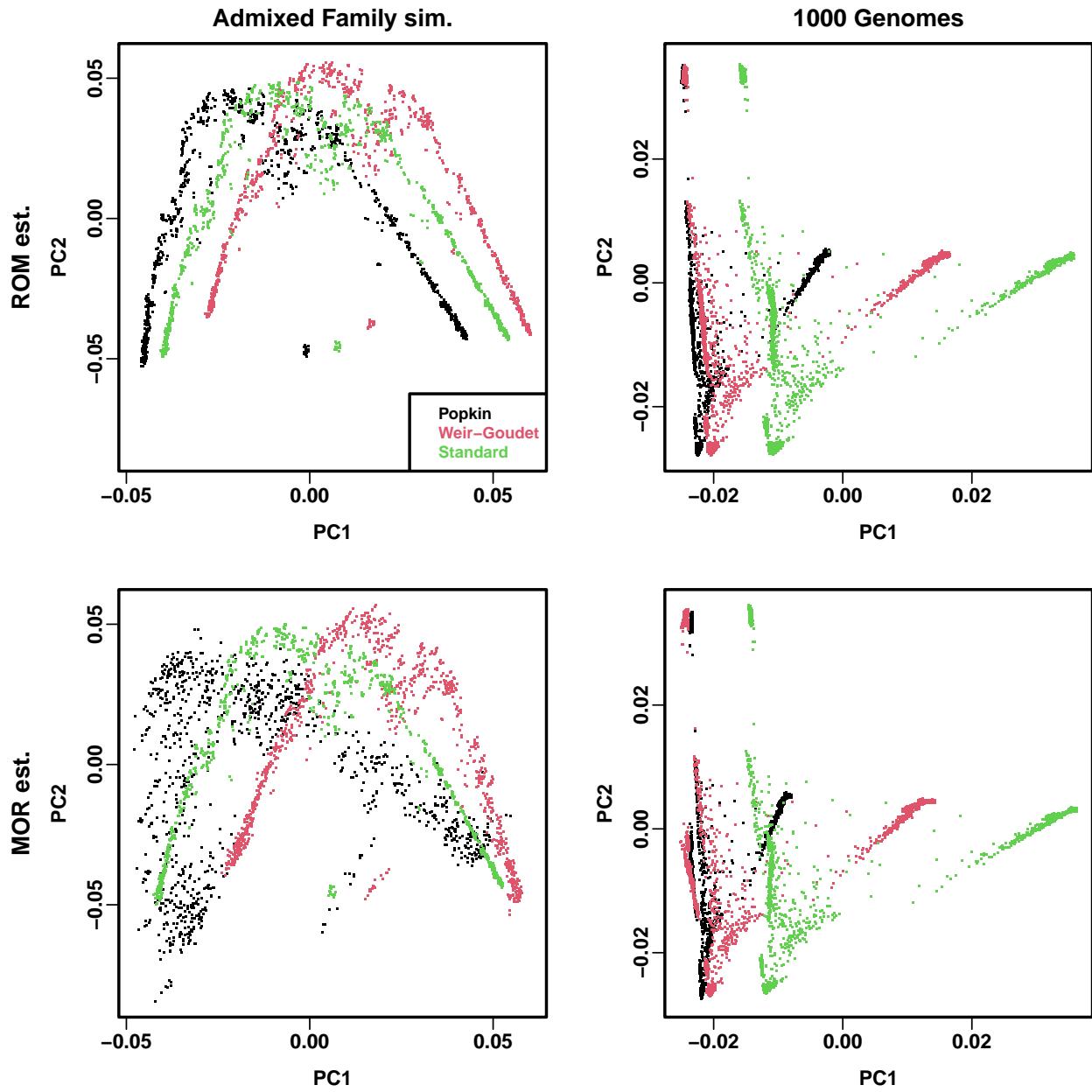


Figure 6: **Visualization of PC shift due to kinship biases.** Each panel shows three estimates (bias types): Popkin, Standard, and Weir-Goudet (all same colors as first panel legend). ROM estimates are shown in first row, MOR in second row. (For simulation, ROM limits are not shown as they were very similar to ROM estimates.) Columns show estimates from each dataset: simulation (1st replicate) and 1000 Genomes. For simulation and popkin only, PC1 and PC2 are replaced with PC2 and PC3 (see text).

3.4 Proof of association invariability to change in ancestral population

The kinship matrices we have used so far have values that depend on the choice of ancestral population T . Here we consider the effect on association of changing ancestral population, and prove that it is compensated for by the relatedness and intercept coefficients, just as it was for common kinship biases.

Suppose we had started with a kinship matrix Φ^S in terms of ancestral population S , and T is a population ancestral to S . If the inbreeding coefficient of S when T is the reference ancestral population is f_S^T , then the kinship matrix Φ^T in terms of T is given by (Ochoa and Storey, 2021)

$$(\mathbf{J} - \Phi^T) = (\mathbf{J} - \Phi^S) (1 - f_S^T).$$

Solving for Φ^T and simplifying results in

$$\Phi^T = (1 - f_S^T) \Phi^S + f_S^T \mathbf{J}.$$

Notice that this resembles the WG bias function but in reverse: whereas WG bias reduces and rescales all kinship values by $\tilde{\varphi}^T$, changing to a more ancestral population rescales and increases all kinship values by f_S^T . Indeed, excluding the trivial degenerate case $f_S^T = 1$, this transformation can be written as Eq. (12) with $c = (1 - f_S^T)^{-1}$ and

$$\begin{aligned} \mathbf{b} &= q \frac{(\Phi^S)^{-1} \mathbf{1}}{\mathbf{1}^\top (\Phi^S)^{-1} \mathbf{1}}, \\ q &= 1 \pm \sqrt{1 + \frac{f_S^T}{1 - f_S^T} (\mathbf{1}^\top (\Phi^S)^{-1} \mathbf{1})}. \end{aligned}$$

The determinant of q is strictly positive, since $\mathbf{1}^\top (\Phi^S)^{-1} \mathbf{1} > 0$ (since Φ^S is positive definite, its inverse is too) and $0 \leq f_S^T < 1$. Thus, the results of the previous section also apply to this transformation: ancestor change is also compensated for by the relatedness and intercept coefficients, which are the only coefficients that depend on the ancestor population, so the association statistics are invariant to this transformation.

3.5 Characterization of non-positive semidefinite and singular kinship and trait covariance estimators

While attempting to validate and characterize the earlier factorization of the WG bias function (Eqs. (12) to (15)), we discovered that this estimator does not produce positive semidefinite matrices, which covariance matrices are required to be. To characterize this problem more generally, we calculated the eigenvalues of all of the kinship matrices Φ^T we produced in all evaluations as well as the trait covariance matrices $\mathbf{V} = 2\sigma^2\Phi^T + \sigma_\epsilon^2\mathbf{I}$ used by LMMs, which we calculated using GCTA's estimates of σ^2 and σ_ϵ^2 .

We found that all WG matrices have very large negative minimum eigenvalues, and popkin MOR estimates also have smaller negative minimum eigenvalues (Fig. S5). Since small negative eigenvalues may be obscured in the previous evaluation, we also counted the proportion of matrices across replicates that had negative eigenvalues. We discovered that, besides all WG matrices and most popkin MOR estimates, Standard matrices were also often non-positive semidefinite but only in 1000 Genomes (Fig. S6). The likely explanation is that, while Standard estimates are positive-semidefinite by construction when there is no missing data, this is no longer true under missingness. Every non-positive semidefinite matrix only had one negative eigenvalue. The eigenvector with the negative eigenvalue for WG bias is partially characterized in Appendix G. Notably, unlike popkin MOR, all popkin ROM estimates are positive definite in every evaluation, including under missingness in 1000 Genomes.

In order to quantify matrix singularity, as well as numerical accuracy problems caused by multiplying by inverses of nearly-singular matrices, we calculated condition numbers, which equal the maximum absolute eigenvalue divided by the minimum absolute eigenvalue of our covariance matrices. As expected, we saw that standard kinship matrices are singular on our admixed family simulation (which lacks missingness), as reflected by extremely high condition numbers, but their trait covariances are well conditioned (have small condition numbers; Fig. S7). Although no other matrices are singular, we do see that popkin MOR estimates in the admixed family simulation have relatively high condition numbers for both the kinship and the trait covariance.

Let us consider the theoretical connection between the eigenvalues of the kinship matrix and

those of \mathbf{V} . The eigendecomposition trick widely used to fit variance components in LMMs (Kang et al., 2008; Lippert et al., 2011; Svishcheva et al., 2012; Zhou and Stephens, 2012; Sul et al., 2018) yields

$$\mathbf{V} = \mathbf{U} (2\sigma^2 \boldsymbol{\Lambda} + \sigma_\epsilon^2 \mathbf{I}) \mathbf{U}^\top,$$

where \mathbf{U} and $\boldsymbol{\Lambda}$ are the eigenvectors and eigenvalues of $\boldsymbol{\Phi}^T$, respectively (Eq. (11)), so the eigenvectors of \mathbf{V} are also \mathbf{U} and its eigenvalues are $2\sigma^2 \boldsymbol{\Lambda} + \sigma_\epsilon^2 \mathbf{I}$. It follows that if $\boldsymbol{\Phi}^T$ is positive definite (all of its eigenvalues are positive) then so is \mathbf{V} (its eigenvalues are also positive), and for these cases the condition number of \mathbf{V} is always smaller (better) than that of $\boldsymbol{\Phi}^T$. A negative kinship eigenvalue λ_i may become positive for \mathbf{V} only if $\lambda_i > -\sigma_\epsilon^2/(2\sigma^2) = -(1-h^2)/(2h^2)$, so very large negative λ_i values as observed for WG do not become positive in \mathbf{V} , in fact often they became more negative (Fig. S5). This equation also shows that \mathbf{V} is always invertible and well-conditioned even when $\boldsymbol{\Phi}^T$ is singular positive semidefinite, as the Standard estimator is under no missingness, since a kinship eigenvalue of zero becomes σ_ϵ^2 for \mathbf{V} . Conversely, the above equation explains why some non-positive semidefinite kinship matrices are particularly problematic: small negative eigenvectors with values near $-\sigma_\epsilon^2/(2\sigma^2)$ can result in ill-conditioned \mathbf{V} . We see that popkin MOR estimates are non-positive semidefinite (Fig. S5) in such a way that some of their \mathbf{V} are ill-conditioned (Fig. S7), and we find that this explains its poorer performance in the admixed family evaluations (Figs. 2 and S2), as shown in the next subsection.

3.6 Further empirical validation of theoretical predictions

Seeing that WG is always non-positive semidefinite, and to query other instances where predictions were not fully met, here we analyze estimation accuracy of various parameters to better understand theoretically and empirically how this broken assumption affects them.

With PCA, no deviations from expectation of AUC_{PR} and $SRMSD_p$ were observed for WG (or any other kinship matrices). This is as expected since PCA simply ignores eigenvectors with negative or zero eigenvalues (which are ranked last). Therefore, our analysis focused on LMM, where deviations were observed and clarification regarding WG is needed.

LMMs such as GCTA perform association testing in two steps. First is the restricted maximum

likelihood (REML) step used to fit variance components. Although the eigendecomposition approaches (Kang et al., 2008; Lippert et al., 2011; Svishcheva et al., 2012; Zhou and Stephens, 2012; Sul et al., 2018) require positive definite \mathbf{V} (lest the determinant of \mathbf{V} be negative), surprisingly the GCTA average information algorithm only requires in practice that \mathbf{V} be invertible (Yang et al., 2011). Thus, we find that variance components estimated from WG are generally close to those of Standard (Fig. S8). Furthermore, the relationship between WG, Standard, and True or Popkin variance components are largely as expected from our theory, with the exception of popkin ROM on 1000 Genomes only, whose genetic variance estimates tended to be a bit smaller than expected (Fig. S9). We were not able to determine a reason that these popkin ROM estimates were worse than expected on 1000 Genomes.

Next we determine the effect of WG bias on coefficient estimates. In this second step of LMM association testing, once \mathbf{V} is determined, GCTA and other LMMs use generalized least squares to estimate fixed effects coefficients (Kang et al., 2008; Kang et al., 2010; Yang et al., 2014). Using the first replicate of the admixture family simulation and the true kinship matrix and the Standard and WG limits only, we recalculated the genetic effect β_i and intercept coefficients α in R for all loci, and confirmed that we recovered the GCTA estimates for β_i to the given precision. We then compare intercept coefficients, which are not given by GCTA, and confirmed our theoretical prediction (Appendix F) that they are identical whether the True or WG ROM limit kinship matrices are used (the mean absolute difference was below 10^{-7}). In contrast, intercepts fit using the Standard ROM limit kinship matrix are different than those of True, which agrees with our theoretical prediction that the intercept compensates for the kinship matrix bias (Fig. S10). However, if the random effect were truly random, then since $\eta = \mathbf{b}^\top \mathbf{s}$ we would expect $\eta \sim \text{Normal}(0, \mathbf{b}^\top \boldsymbol{\Phi}^T \mathbf{b})$, which does not agree with our empirical observations that the mean of η is non-zero, and which has some dependence on α . This disagreement is interesting but not unexpected, since maximum likelihood estimates of α are bound to differ from their theoretical values.

Lastly, we want to explain the largest deviations from our predictions at the AUC_{PR} and $SRMSD_p$ level, which are limited to popkin ROM in the admixture simulation and popkin MOR for 1000 Genomes. Here we use the WG estimates as reference since WG \mathbf{V} matrices generally

had the lowest condition numbers, so they should yield the most numerically accurate estimates (Fig. S7). We find that errors in the genetic variance component estimation (σ^2) strongly drive the observed popkin ROM errors in 1000 Genomes in SRMSD_p expectation, and to a lesser extent errors in AUC_{PR} as well (Fig. S11). However, for the popkin MOR errors in the admixed family simulation, which were much larger, we did not see σ^2 estimation errors (Fig. S9), and the SRMSD_p and AUC_{PR} errors are instead well explained by the condition number of \mathbf{V} , which is defined for its influence of regression coefficient estimation accuracy (Fig. S12).

4 Discussion

Previous research showed that commonly used kinship estimators are biased, and that these biases can be large (Ochoa and Storey (2021); Fig. 1). We initiated the present work under the hypothesis that these kinship biases would affect association testing, but surprisingly find that association is unaffected by these kinship biases. We then prove theoretically that it is the intercept and relatedness (random effect or PCs) coefficients that compensate for the bias, and result in identical genetic effect coefficients and significance statistics.

One previously uncomfortable fact was that kinship estimates depend on the choice of ancestral population, which conditions the distributions of allele frequencies and genotypes, but the effect of this choice of association testing was not only unknown but completely disregarded. A corollary of our theoretical results is that changes of ancestral population, which behave algebraically like the reverse of the WG kinship bias, are also compensated for by the relatedness and intercept coefficients, so association testing is also invariant to the choice of ancestral population. Thus, although a choice of ancestral population is always being made when estimating kinship, this choice is fortunately inconsequential to association testing, as it ought to be since the relatedness structure overall is being conditioned upon in these tests.

Given that kinship bias type is not important for association studies, we are free to choose a kinship estimator based on other properties. The biased standard kinship matrix may be more desirable than the popkin estimator based on the numerical stability we observed in our simulations. In particular, while theory shows that the solutions should be the same for all estimators of the same

type, we find that popkin’s statistics disagree more often from the standard and WG estimators, namely LMM association with popkin MOR (admixed family simulation, Fig. 2, Fig. S2) and popkin ROM (1000 Genomes, Fig. S3). The standard kinship matrix is orthogonal to the intercept, because of the centering operation applied to obtain it in our theoretical results, whereas the popkin and true kinship matrices are not orthogonal to the intercept. Thus, PCA regression with the eigenvectors of the standard kinship matrix is more numerically stable (because more covariates are linearly independent) than the popkin counterpart. We believe that the observed popkin disagreements in LMMs are due to poor convergence of that algorithm in those cases.

We also found that all MOR estimators perform better in the LMM association (and overall) compared to the ROM versions in the 1000 Genomes evaluation. Perhaps this is expected because our trait simulation follows the “fixed effect sizes” model, in which rare variants have larger coefficients, and the MOR estimators also weigh rare variants more highly in estimating kinship coefficients. This effect was not observed in the admixed family simulation, where MOR and ROM versions gave similar kinship estimates and performed similarly, compared to the real data evaluation, where kinship estimates were also strikingly different. However, only the popkin ROM estimator is unbiased (Fig. 1B, Fig. S1), so it is unclear why the biased popkin MOR estimator performs better in this setting. One potential explanation is that our kinship model assumes that all variants were preexisting in the MRCA population, whereas rare variants in human data are known to be very recent mutations, and thus their effective kinship matrix is different than that of ancestral variants. Therefore, despite its biases, it is possible that the popkin MOR estimator is more accurately capturing the kinship matrix of these rare variants and thus modeling them better in association tests, particularly in LMMs where the effect is most pronounced.

Our conclusions extend to variants of the standard kinship estimator that weigh loci according to linkage disequilibrium (Speed et al., 2017; Wang et al., 2017), which have the same bias form since this bias is present in each individual locus (Ochoa and Storey, 2021). As shown in our proof, the more general form of the standard kinship estimator that weighs individuals to estimate ancestral allele frequencies \hat{p}_i^T is also subject to the same conclusions. Such weighted \hat{p}_i^T estimates include the best unbiased linear estimator (Astle and Balding, 2009; Thornton and McPeek, 2010).

In this study, we show empirically and theoretically that association tests are invariant to the use of common kinship estimators that are biased as well as a more recent unbiased estimator. The theoretical underpinnings of our proof show that the same is expected of any generalized linear model with the same setup, namely intercept and population structure with coefficients that are nuisance variables, which includes case/control models as well as the quantitative trait model we explicitly studied here. However, heritability estimation requires unbiased estimates of the random effect coefficient (σ^2), so our results prove that it will be biased when the standard kinship estimator is used, as it is using GCTA (Yang et al., 2011; Yang et al., 2014). Nevertheless, heritability estimation is a complex problem and its full study is beyond the scope of this work. Overall, we have described an unexpected robustness of association studies, and our theoretical understanding of this result may help guide future improvements for association and other related models.

Declaration of interests

The authors declare no competing interests.

Acknowledgments

This work was funded in part by the Duke University School of Medicine Whitehead Scholars Program, a gift from the Whitehead Charitable Foundation. The 1000 Genomes data were generated at the New York Genome Center with funds provided by NHGRI Grant 3UM1HG008901-03S1.

Web resources

plink2, <https://www.cog-genomics.org/plink/2.0/>

GCTA, <https://yanglab.westlake.edu.cn/software/gcta/>

bnpsd, <https://cran.r-project.org/package=bnpsd>

simfam, <https://cran.r-project.org/package=simfam>

simtrait, <https://cran.r-project.org/package=simtrait>

popkin, <https://cran.r-project.org/package=popkin>
popkinsuppl, <https://github.com/OchoaLab/popkinsuppl>

Data and code availability

The data and code generated during this study are available on GitHub at <https://github.com/OchoaLab/bias-assoc-paper>. The high-coverage version of the 1000 Genomes Project was downloaded from ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/working/20190425_NYGC_GATK/.

References

- 1000 Genomes Project Consortium et al. (2012). “An integrated map of genetic variation from 1,092 human genomes”. *Nature* 491(7422), pp. 56–65.
- Altschul, Stephen F., Raymond J. Carroll, and David J. Lipman (1989). “Weights for data related by a tree”. *Journal of Molecular Biology* 207(4), pp. 647–653.
- Astle, William and David J. Balding (2009). “Population Structure and Cryptic Relatedness in Genetic Association Studies”. *Statist. Sci.* 24(4), pp. 451–471.
- Aulchenko, Yurii S., Dirk-Jan de Koning, and Chris Haley (2007). “Genomewide rapid association using mixed model and regression: a fast and simple method for genomewide pedigree-based quantitative trait loci association analysis”. *Genetics* 177(1), pp. 577–585.
- Balding, D. J. and R. A. Nichols (1995). “A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity”. *Genetica* 96(1-2), pp. 3–12.
- Bhatia, Gaurav et al. (2013). “Estimating and interpreting FST: the impact of rare variants”. *Genome Res.* 23(9), pp. 1514–1521.
- Chang, Christopher C. et al. (2015). “Second-generation PLINK: rising to the challenge of larger and richer datasets”. *GigaScience* 4(1), p. 7.

- Consortium, The 1000 Genomes Project (2010). “A map of human genome variation from population-scale sequencing”. *Nature* 467(7319), pp. 1061–1073.
- Devlin, B. and Kathryn Roeder (1999). “Genomic Control for Association Studies”. *Biometrics* 55(4), pp. 997–1004.
- Emik, L. Otis and Clair E. Terrill (1949). “Systematic procedures for calculating inbreeding coefficients”. *J Hered* 40(2), pp. 51–55.
- Fairley, Susan et al. (2020). “The International Genome Sample Resource (IGSR) collection of open human genomic variation resources”. *Nucleic Acids Research* 48(D1), pp. D941–D947.
- García-Cortés, Luis Alberto (2015). “A novel recursive algorithm for the calculation of the detailed identity coefficients”. *Genetics Selection Evolution* 47(1), p. 33.
- Hoffman, Gabriel E. (2013). “Correcting for population structure and kinship using the linear mixed model: theory and extensions”. *PLoS ONE* 8(10), e75707.
- Jacquard, Albert (1970). *Structures génétiques des populations*. Paris: Masson et Cie.
- Kang, Hyun Min et al. (2008). “Efficient control of population structure in model organism association mapping”. *Genetics* 178(3), pp. 1709–1723.
- Kang, Hyun Min et al. (2010). “Variance component model to account for sample structure in genome-wide association studies”. *Nat. Genet.* 42(4), pp. 348–354.
- Lippert, Christoph et al. (2011). “FaST linear mixed models for genome-wide association studies”. *Nat. Methods* 8(10), pp. 833–835.
- Loh, Po-Ru et al. (2015). “Efficient Bayesian mixed-model analysis increases association power in large cohorts”. *Nat. Genet.* 47(3), pp. 284–290.
- Malécot, Gustave (1948). *Mathématiques de l'hérédité*. Masson et Cie.
- Ochoa, Alejandro and John D. Storey (2021). “Estimating FST and kinship for arbitrary population structures”. *PLoS Genet* 17(1), e1009241.
- Price, Alkes L. et al. (2006). “Principal components analysis corrects for stratification in genome-wide association studies”. *Nat. Genet.* 38(8), pp. 904–909.

- Rakovski, Cyril S. and Daniel O. Stram (2009). “A kinship-based modification of the armitage trend test to address hidden population structure and small differential genotyping errors”. *PLoS ONE* 4(6), e5825.
- Sherman, Jack and Winifred J. Morrison (1950). “Adjustment of an Inverse Matrix Corresponding to a Change in One Element of a Given Matrix”. *The Annals of Mathematical Statistics* 21(1), pp. 124–127.
- Speed, Doug and David J. Balding (2015). “Relatedness in the post-genomic era: is it still useful?” *Nat. Rev. Genet.* 16(1), pp. 33–44.
- Speed, Doug et al. (2012). “Improved heritability estimation from genome-wide SNPs”. *Am. J. Hum. Genet.* 91(6), pp. 1011–1021.
- Speed, Doug et al. (2017). “Reevaluation of SNP heritability in complex human traits”. *Nat Genet* 49(7), pp. 986–992.
- Sul, Jae Hoon, Lana S. Martin, and Eleazar Eskin (2018). “Population structure in genetic studies: Confounding factors and mixed models”. *PLoS Genet.* 14(12), e1007309.
- Svishcheva, Gulnara R. et al. (2012). “Rapid variance components–based method for whole-genome association analysis”. *Nat Genet* 44(10), pp. 1166–1170.
- Thornton, Timothy and Mary Sara McPeek (2010). “ROADTRIPS: case-control association testing with partially or completely unknown population and pedigree structure”. *Am. J. Hum. Genet.* 86(2), pp. 172–184.
- Voight, Benjamin F. and Jonathan K. Pritchard (2005). “Confounding from Cryptic Relatedness in Case-Control Association Studies”. *PLOS Genetics* 1(3), e32.
- Wang, Bowen, Serge Sverdlov, and Elizabeth Thompson (2017). “Efficient Estimation of Realized Kinship from SNP Genotypes”. *Genetics, genetics*.116.197004.
- Weir, Bruce S. and Jérôme Goudet (2017). “A Unified Characterization of Population Structure and Relatedness”. *Genetics* 206(4), pp. 2085–2103.
- Wright, Sewall (1922). “Coefficients of Inbreeding and Relationship”. *The American Naturalist* 56(645), pp. 330–338.

- Xie, C., D. D. Gessler, and S. Xu (1998). “Combining different line crosses for mapping quantitative trait loci using the identical by descent-based variance component method”. *Genetics* 149(2), pp. 1139–1146.
- Yang, Jian et al. (2010). “Common SNPs explain a large proportion of the heritability for human height”. *Nat. Genet.* 42(7), pp. 565–569.
- Yang, Jian et al. (2011). “GCTA: a tool for genome-wide complex trait analysis”. *Am. J. Hum. Genet.* 88(1), pp. 76–82.
- Yang, Jian et al. (2014). “Advantages and pitfalls in the application of mixed-model association methods”. *Nat Genet* 46(2), pp. 100–106.
- Yao, Yiqi and Alejandro Ochoa (2022). *Limitations of principal components in quantitative genetic association models for human studies*. Tech. rep. bioRxiv, p. 2022.03.25.485885.
- Yu, Jianming et al. (2006). “A unified mixed-model method for association mapping that accounts for multiple levels of relatedness”. *Nat. Genet.* 38(2), pp. 203–208.
- Zhou, Xiang and Matthew Stephens (2012). “Genome-wide efficient mixed-model analysis for association studies”. *Nat. Genet.* 44(7), pp. 821–824.

Appendices

A Justification for popkin generalizations

The popkin estimator in Eq. (1) has been generalized in this work to include locus weights w_i . The original formulation had $w_i = 1$ for all loci i (Ochoa and Storey, 2021). Recalling from that original work that

$$\mathbb{E} [(x_{ij} - 1)(x_{ik} - 1) - 1 | T] = 4p_i^T (1 - p_i^T) (\varphi_{jk}^T - 1),$$

then for fixed w_i we get

$$\begin{aligned}\mathbb{E}[A_{jk}|T] &= v_m^T (\varphi_{jk}^T - 1), \\ v_m^T &= \frac{4}{m} \sum_{i=1}^m w_i p_i^T (1 - p_i^T).\end{aligned}$$

Therefore, as before all the unknowns p_i^T and now also the (known) weights w_i collapse into a single parameter v_m^T , which is estimated under the assumption that the minimum kinship is zero, giving $\hat{A}_{\min} = -v_m^T$, so that

$$\hat{\varphi}_{jk}^{T,\text{popkin-ROM}} = 1 - \frac{A_{jk}}{\hat{A}_{\min}} \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \varphi_{jk}^T$$

as desired.

The MOR case of $w_i = (\hat{p}_i^T (1 - \hat{p}_i^T))^{-1}$ does not fit the previous case because this w_i is a random variable (it is a function of the genotypes). The term of interest $w_i((x_{ij} - 1)(x_{ik} - 1) - 1)$ is a ratio of random variables whose expectation does not have a closed form. In this case, we rely on the first-order approximation to this expectation, namely

$$\begin{aligned}\mathbb{E} \left[\frac{(x_{ij} - 1)(x_{ik} - 1) - 1}{\hat{p}_i^T (1 - \hat{p}_i^T)} \middle| T \right] &\approx \frac{\mathbb{E}[(x_{ij} - 1)(x_{ik} - 1) - 1|T]}{\mathbb{E}[\hat{p}_i^T (1 - \hat{p}_i^T)|T]} \\ &= \frac{4p_i^T (1 - p_i^T) (\varphi_{jk}^T - 1)}{p_i^T (1 - p_i^T) (1 - \bar{\varphi}^T)} \\ &= \frac{4(\varphi_{jk}^T - 1)}{1 - \bar{\varphi}^T},\end{aligned}$$

where the expectation of $\hat{p}_i^T (1 - \hat{p}_i^T)$ was calculated previously (Ochoa and Storey, 2021). In this case the expectation of A_{jk} , summing across loci, is also approximated by

$$\mathbb{E}[A_{jk}|T] \approx \frac{4(\varphi_{jk}^T - 1)}{1 - \bar{\varphi}^T}.$$

The same strategy as before applies to estimate the unknown factor $4/(1 - \bar{\varphi}^T)$, namely that if the

minimum kinship is zero then $\hat{A}_{\min} \approx -4/(1 - \bar{\varphi}^T)$, resulting in

$$\hat{\varphi}_{jk}^{T,\text{popkin-MOR}} = 1 - \frac{A_{jk}}{\hat{A}_{\min}} \approx \varphi_{jk}^T.$$

B Connection between popkin and standard kinship estimator

Since the connection we discovered holds when data is complete but not under missingness, to determine necessary conditions, here we introduce more complete forms of the estimators that handle missingness. The generalized popkin estimator (including both ROM and MOR special cases) is

$$\begin{aligned} A_{ijk} &= I_{ij}I_{ik}((x_{ij} - 1)(x_{ik} - 1) - 1), \\ A_{jk} &= \frac{1}{m_{jk}} \sum_{i=1}^m w_i A_{ijk}, \\ m_{jk} &= \sum_{i=1}^m I_{ij}I_{ik}, \end{aligned}$$

where $I_{ij} = 1$ if x_{ij} is not missing, 0 otherwise (this way missing x_{ij} can be treated as having any finite value and not contribute to the estimator). Note that only loci where both genotypes (x_{ij} and x_{ik}) are non-missing are included in the above average, and m_{jk} counts the total number of such loci. The ancestral allele frequency estimator with missingness is

$$\begin{aligned} \hat{p}_i^T &= \frac{1}{2n_i} \sum_{j=1}^n I_{ij}x_{ij}, \\ n_i &= \sum_{j=1}^n I_{ij}, \end{aligned}$$

which averages over individuals rather than loci, so its denominator is the number of non-missing individuals at this locus. Let us compute some averages of the generalized popkin estimator. Since the result we want holds at every locus separately, let us formulate the averages of interest at locus

i only:

$$\begin{aligned}\bar{A}_{ij} &= \frac{1}{n} \sum_{k=1}^n A_{ijk} = I_{ij} \frac{n_i}{n} ((x_{ij} - 1) (2\hat{p}_i^T - 1) - 1), \\ \bar{A}_i &= \frac{1}{n} \sum_{k=1}^n \bar{A}_{ik} = - \left(\frac{n_i}{n} \right)^2 4\hat{p}_i^T (1 - \hat{p}_i^T).\end{aligned}$$

Therefore, the combination of interest is:

$$\begin{aligned}A_{ijk} + \bar{A}_i - \bar{A}_{ij} - \bar{A}_{ik} &= I_{ij} I_{ik} (x_{ij} - 2\hat{p}_i^T) (x_{ik} - 2\hat{p}_i^T) \\ &\quad + \frac{n_i}{n} \left(I_{ij} - \frac{n_i}{n} \right) 4\hat{p}_i^T + \left(\left(\frac{n_i}{n} \right)^2 - I_{ij} I_{ik} \right) 4(\hat{p}_i^T)^2 \\ &\quad + I_{ij} \left(I_{ik} - \frac{n_i}{n} \right) x_{ij} (2\hat{p}_i^T - 1) + I_{ik} \left(I_{ij} - \frac{n_i}{n} \right) x_{ik} (2\hat{p}_i^T - 1).\end{aligned}$$

To arrive at the desired result of $I_{ij} I_{ik} (x_{ij} - 2\hat{p}_i^T) (x_{ik} - 2\hat{p}_i^T)$, which is the first term above, it is necessary for the rest of the terms to vanish for arbitrary values of \hat{p}_i^T , x_{ij} , and x_{ik} . Since $n_i > 0$ (there is at least one non-missing individual at every locus), the term $\frac{n_i}{n} (I_{ij} - \frac{n_i}{n}) 4\hat{p}_i^T$ vanishes if and only if $I_{ij} = \frac{n_i}{n}$, and since $I_{jk} = 0$ does not solve this equation (because $n_i > 0$) the only other case is $I_{jk} = 1$, which requires $n_i = n$, so no individuals can have missing data at this locus. Thus,

$$A_{ijk} + \bar{A}_i - \bar{A}_{ij} - \bar{A}_{ik} = I_{ij} I_{ik} (x_{ij} - 2\hat{p}_i^T) (x_{ik} - 2\hat{p}_i^T)$$

if and only if there is no missing data at locus i . The other desired result of

$$\bar{A}_i = -4\hat{p}_i^T (1 - \hat{p}_i^T)$$

also requires $n_i = n$.

Assuming now no missingness, transforming the popkin estimates as desired gives

$$\begin{aligned}
\frac{\hat{\varphi}_{jk}^{T,\text{popkin}} + \bar{\hat{\varphi}}^{T,\text{popkin}} - \bar{\hat{\varphi}}_j^{T,\text{popkin}} - \bar{\hat{\varphi}}_k^{T,\text{popkin}}}{1 - \bar{\hat{\varphi}}^{T,\text{popkin}}} &= \frac{A_{jk} + \bar{A} - \bar{A}_j - \bar{A}_k}{-\bar{A}} \\
&= \frac{\sum_{i=1}^m w_i (A_{ijk} + \bar{A}_i - \bar{A}_{ij} - \bar{A}_{ik})}{-\sum_{i=1}^m w_i \bar{A}_i} \\
&= \frac{\sum_{i=1}^m w_i (x_{ij} - 2\hat{p}_i^T) (x_{ik} - 2\hat{p}_i^T)}{\sum_{i=1}^m w_i 4\hat{p}_i^T (1 - \hat{p}_i^T)}.
\end{aligned}$$

Therefore, if popkin-ROM is input ($w_i = 1$), this transformation yields std-ROM. On the other hand, if popkin-MOR is used ($w_i^{-1} = \hat{p}_i^T (1 - \hat{p}_i^T)$), the transformation yields std-MOR.

C Mean kinship inequalities

Denote the mean of the diagonal kinship terms as $\bar{\delta}^T = \frac{1}{n} \sum_{j=1}^n \varphi_{jj}^T$. Here we prove that

$$0 \leq \tilde{\varphi}^T \leq \bar{\varphi}^T \leq \bar{\delta}^T \leq 1,$$

with each of $\tilde{\varphi}^T = \bar{\varphi}^T$ and $\bar{\varphi}^T = \bar{\delta}^T$ if and only if all kinship values are equal.

The inequalities $0 \leq \tilde{\varphi}^T \leq \bar{\varphi}^T \leq \bar{\delta}^T \leq 1$ follow directly from previous work, applied to a kinship matrix rather than a coancestry matrix as done originally, as the proof required solely a covariance matrix with values between 0 and 1 (Ochoa and Storey, 2021). Recall that $\tilde{\varphi}^T$ is defined in Eq. (7). The lower bound $0 \leq \tilde{\varphi}^T$ follows since every kinship value is non-negative. Note that $\bar{\varphi}^T$ and $\tilde{\varphi}^T$ are related by

$$\bar{\varphi}^T = \frac{\tilde{\varphi}^T(n-1) + \bar{\delta}^T}{n}. \quad (17)$$

Applying $\bar{\varphi}^T \leq \bar{\delta}^T$ to Eq. (17) and simplifying yields $\tilde{\varphi}^T \leq \bar{\delta}^T$. Lastly, since $\bar{\varphi}^T - \tilde{\varphi}^T = (\bar{\delta}^T - \tilde{\varphi}^T)/n$ (from rearranging Eq. (17)), it also follows that $\tilde{\varphi}^T \leq \bar{\varphi}^T$, as desired. Furthermore, $\tilde{\varphi}^T = \bar{\varphi}^T$ holds if and only if all $\varphi_{jk}^T = \bar{\delta}^T$, since that is necessary and sufficient for $\bar{\varphi}^T = \bar{\delta}^T$.

D Derivation of WG bias factorization

Here we rewrite the WG bias function of Eq. (6) as a factorization of the form of Eq. (12). It is easy to see that $c = 1 - \tilde{\varphi}^T$. Expanding Eq. (12) gives

$$\begin{aligned}\mathbf{B}\Phi^T\mathbf{B}^\top &= (\mathbf{I} - \mathbf{1}\mathbf{b}^\top)\Phi^T(\mathbf{I} - \mathbf{b}\mathbf{1}^\top) \\ &= \Phi^T - \mathbf{1}(\Phi^T\mathbf{b})^\top - (\Phi^T\mathbf{b})\mathbf{1}^\top + \mathbf{J}(\mathbf{b}^\top\Phi^T\mathbf{b}),\end{aligned}$$

where $\mathbf{b}^\top\Phi^T\mathbf{b}$ is a scalar and $\Phi^T\mathbf{b}$ a vector. Equating the above to Eq. (6) and rearranging, we obtain

$$\mathbf{J}(\tilde{\varphi}^T + (\mathbf{b}^\top\Phi^T\mathbf{b})) = \mathbf{1}(\Phi^T\mathbf{b})^\top + (\Phi^T\mathbf{b})\mathbf{1}^\top.$$

Since $\tilde{\varphi}^T + (\mathbf{b}^\top\Phi^T\mathbf{b})$ is a scalar and $\mathbf{J} = \mathbf{1}\mathbf{1}^\top$, we can see that the solution requires the right side to also be a constant matrix, which is only achieved if $\Phi^T\mathbf{b} \propto \mathbf{1}$. We choose the scaling factor for the last $\mathbf{1}$ to be $q(\mathbf{1}^\top(\Phi^T)^{-1}\mathbf{1})^{-1}$ as this simplifies notation later, and solving for \mathbf{b} results in Eq. (14). To solve for q , we replace \mathbf{b} from Eq. (14) into the above equation, which after rearranging results in

$$q^2 - 2q + \tilde{\varphi}^T \left(\mathbf{1}^\top (\Phi^T)^{-1} \mathbf{1} \right) = 0.$$

The solution to the above quadratic equation is given by Eq. (15), as desired.

E Minimum weighted mean kinship

Consider the weighted mean kinship value $\mathbf{w}^\top\Phi^T\mathbf{w}$, where \mathbf{w} are weights that sum to one ($\mathbf{w}^\top\mathbf{1} = 1$). The ordinary (or unweighted) mean kinship $\bar{\varphi}^T$ as given in the rest of this work is the special case with $\mathbf{w} = \frac{1}{n}\mathbf{1}$. The value of the weights, constrained to sum to one, that minimize the weighted

mean kinship, is found by solving the Lagrangian multiplier problem

$$G = \mathbf{w}^\top \Phi^T \mathbf{w} + \lambda(\mathbf{w}^\top \mathbf{1} - 1).$$

The derivatives are the constraint and $\frac{dG}{d\mathbf{w}} = 2\Phi^T \mathbf{w} + \lambda \mathbf{1} = \mathbf{0}$. The optimal weights thus satisfy $\mathbf{w} = \frac{-\lambda}{2} (\Phi^T)^{-1} \mathbf{1}$. Multiplying by $\mathbf{1}^\top$, since $\mathbf{1}^\top \mathbf{w} = 1$, allows us to solve for $\lambda^{-1} = -\frac{1}{2} \mathbf{1}^\top (\Phi^T)^{-1} \mathbf{1}$.

Thus, the optimal weights are

$$\mathbf{w} = \frac{(\Phi^T)^{-1} \mathbf{1}}{\mathbf{1}^\top (\Phi^T)^{-1} \mathbf{1}},$$

a solution that recurs as optimal weights in related settings (Altschul et al., 1989; Astle and Balding, 2009). Therefore, the minimum weighted mean kinship is, and satisfies,

$$\mathbf{w}^\top \Phi^T \mathbf{w} = \frac{1}{\mathbf{1}^\top (\Phi^T)^{-1} \mathbf{1}} \leq \bar{\varphi}^T.$$

F Proof that WG bias results in zero intercept shift under LMM generalized least squares estimation

For this section suppose that variance components have been estimated, so $\mathbf{V} = 2\sigma^2 \Phi^T + \sigma_\epsilon^2 \mathbf{I}$ is given, assume it is non-singular, and rewrite the LMM association model as

$$\mathbf{y} = \mathbf{Z}\beta + \boldsymbol{\epsilon}_V, \quad \boldsymbol{\epsilon}_V \sim \text{Normal}(\mathbf{0}, \mathbf{V}),$$

where the design matrix $\mathbf{Z} = (\mathbf{1}, \mathbf{x}_i, \dots)$ contains the intercept, genotype and now additional covariates, and $\beta = (\alpha, \beta_i, \dots)$ are their coefficients. The generalized least squares coefficients estimate, used by GCTA and other LMMs, is

$$\hat{\beta} = (\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{y}.$$

Now suppose \mathbf{V} corresponds to some kinship matrix Φ^T while \mathbf{V}' corresponds to $\Phi^{T'} = F^{\text{WG}}(\Phi^T)$, and \mathbf{V}' is also non-singular. Our strategy involves repeated application of the Sherman-Morrison

formula for calculating inverses of matrices after a rank-1 update, which for a symmetric update of a matrix \mathbf{A} with a vector \mathbf{z} and a scalar b takes the form (Sherman and Morrison, 1950)

$$(\mathbf{A} + b\mathbf{z}\mathbf{z}^\top)^{-1} = \mathbf{A}^{-1} - \frac{b}{1 + b(\mathbf{z}^\top \mathbf{A}^{-1} \mathbf{z})} (\mathbf{A}^{-1} \mathbf{z}) (\mathbf{A}^{-1} \mathbf{z})^\top.$$

Since $F^{\text{WG}}(\boldsymbol{\Phi}^T)$ is a rank-1 update of $\boldsymbol{\Phi}^T$ by Eq. (6), then \mathbf{V}' is also a rank-1 update of \mathbf{V} :

$$\begin{aligned} \mathbf{V}' &= 2\sigma^{2\prime} \boldsymbol{\Phi}^{T'} + \sigma_\epsilon^2 \mathbf{I} \\ &= 2\sigma^2 (\boldsymbol{\Phi}^T - \tilde{\varphi}^T \mathbf{1} \mathbf{1}^\top) + \sigma_\epsilon^2 \mathbf{I} \\ &= \mathbf{V} - d \mathbf{1} \mathbf{1}^\top, \end{aligned}$$

where $d = 2\sigma^2 \tilde{\varphi}^T$ and we used $\sigma^{2\prime} = (1 - \tilde{\varphi}^T) \sigma^2$. Therefore,

$$(\mathbf{V}')^{-1} = \mathbf{V}^{-1} + e \mathbf{V}^{-1} \mathbf{1} (\mathbf{V}^{-1} \mathbf{1})^\top,$$

where $e = d / (1 - d(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1}))$. Therefore the following remains a rank-1 update,

$$\mathbf{Z}^\top (\mathbf{V}')^{-1} \mathbf{Z} = \mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z} + e \mathbf{u} \mathbf{u}^\top,$$

where $\mathbf{u} = \mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{1}$ is a column vector the length of the number of covariates (including intercept and genotype). Therefore,

$$(\mathbf{Z}^\top (\mathbf{V}')^{-1} \mathbf{Z})^{-1} = (\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1} - g \mathbf{v} \mathbf{v}^\top,$$

where $\mathbf{v} = (\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1} \mathbf{u}$ and $g = e / (1 + e(\mathbf{u}^\top \mathbf{v}))$. Noting that the first column of \mathbf{Z} is the intercept, and denoting \mathbf{Z}_{-1} the matrix without the first column (so $\mathbf{Z} = (\mathbf{1}, \mathbf{Z}_{-1})$), a major simplification

occurs:

$$\begin{aligned}
\mathbf{v} &= (\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{1} \\
&= \begin{pmatrix} \mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1} & \mathbf{1}^\top \mathbf{V}^{-1} \mathbf{Z}_{-1} \\ \mathbf{Z}_{-1}^\top \mathbf{V}^{-1} \mathbf{1} & \mathbf{Z}_{-1}^\top \mathbf{V}^{-1} \mathbf{Z}_{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1} \\ \mathbf{Z}_{-1}^\top \mathbf{V}^{-1} \mathbf{1} \end{pmatrix} \\
&= \begin{pmatrix} 1 \\ \mathbf{0} \end{pmatrix},
\end{aligned}$$

where $\mathbf{0}$ is a vector the length of the number of covariates minus one (exclude the intercept). In other words, \mathbf{v} is the first column of the identity matrix because $\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{1}$ is the first column of $\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z}$. As a consequence, $\mathbf{Z}\mathbf{v} = \mathbf{1}$, so $\mathbf{u}^\top \mathbf{v} = \mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1}$ and

$$\begin{aligned}
g &= \frac{e}{1 + e(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})} \\
&= \frac{\frac{d}{1 - d(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})}}{1 + \frac{d}{1 - d(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})}(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})} \\
&= d.
\end{aligned}$$

The final step yields the coefficient estimates as a rank-1 update:

$$\begin{aligned}
\hat{\beta}' &= (\mathbf{Z}^\top (\mathbf{V}')^{-1} \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{V}')^{-1} \mathbf{y} \\
&= ((\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1} - d\mathbf{v}\mathbf{v}^\top) \mathbf{Z}^\top (\mathbf{V}^{-1} + e\mathbf{V}^{-1} \mathbf{1} (\mathbf{V}^{-1} \mathbf{1})^\top) \mathbf{y} \\
&= \hat{\beta} + e\mathbf{v} (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{y}) - d\mathbf{v} (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{y}) - d\mathbf{v} (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1}) (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{y}) \\
&= \hat{\beta} + \mathbf{v} (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{y}) (e - d - de (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})).
\end{aligned}$$

The last factor above vanishes:

$$\begin{aligned}
e - d - de (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1}) &= \frac{d}{1 - d(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})} - d - d \frac{d}{1 - d(\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1})} (\mathbf{1}^\top \mathbf{V}^{-1} \mathbf{1}) \\
&= 0.
\end{aligned}$$

Therefore, $\hat{\beta}' = \hat{\beta}$, which shows that all fixed effect coefficients, including the intercept, are invariant to using a WG estimate instead of the unbiased kinship matrix when estimated with generalized least squares.

Furthermore, since the diagonal values of $(\mathbf{Z}^\top (\mathbf{V}')^{-1} \mathbf{Z})^{-1}$, which correspond to $\text{Var}(\hat{\beta}'_k)$ for each k , are the same as those of $(\mathbf{Z}^\top \mathbf{V}^{-1} \mathbf{Z})^{-1}$ except for the first one corresponding to the intercept, then the Wald test statistic of the k th covariate coefficients, given by $\hat{\beta}_k^2 / \text{Var}(\hat{\beta}_k)$, and their p-values, are also the same for $k \neq 1$ for WG bias as for the unbiased kinship matrix.

G Characterization of WG bias positive definite subspaces, and eigenvector with negative eigenvalue

We found that WG estimates are non-positive semidefinite, and that this is due to a single eigenvalue that is negative. While we did not find a closed form for this eigenvector, we prove that it must have non-zero projections onto two orthogonal subspaces: S_1 spanned by the intercept $\mathbf{1}$, and S_2 its complement (orthogonal to $\mathbf{1}$). This is because, starting from a positive definite kinship matrix Φ^T , the WG-bias transformed matrix $\Phi^{T'} = F^{\text{WG}}(\Phi^T)$ is positive definite in each of the subspaces S_1, S_2 , so the eigenvector with the negative eigenvalue cannot be wholly in just one of the two subspaces, so it must have non-zero projections to both. Recall from Eq. (6) that $\Phi^{T'} = \frac{1}{1-\tilde{\varphi}^T} (\Phi^T - \tilde{\varphi}^T \mathbf{J})$. We shall not consider $\Phi^T = \mathbf{J}$ as a valid kinship matrix, which therefore ensures that $\tilde{\varphi}^T < 1$ as there is at least one kinship value with $\varphi_{jk}^T < 1$. In both subspaces we will prove that $\mathbf{v} \in S_i$ and $\mathbf{v} \neq \mathbf{0}$ implies $\mathbf{v}^\top \Phi^{T'} \mathbf{v} > 0$ which proves that $\Phi^{T'}$ is positive definite in that subspace.

We begin by considering $\mathbf{v} \in S_2$, which satisfy $\mathbf{1}^\top \mathbf{v} = \mathbf{0}$, and by hypothesis $\mathbf{v} \neq \mathbf{0}$. Therefore $\mathbf{v}^\top \mathbf{J} \mathbf{v} = 0$ in this subspace, which results in

$$\mathbf{v}^\top \Phi^{T'} \mathbf{v} = \frac{1}{1-\tilde{\varphi}^T} \mathbf{v}^\top \Phi^T \mathbf{v} > 0,$$

where the final inequality follows since the original kinship matrix is positive definite and $1-\tilde{\varphi}^T > 0$.

Now consider $\mathbf{v} \in S_1$, which are necessarily of the form $\mathbf{v} = v\mathbf{1}$, and by hypothesis $v \neq 0$.

Therefore

$$\mathbf{v}^\top \boldsymbol{\Phi}^{T'} \mathbf{v} = \frac{v^2}{1 - \tilde{\varphi}^T} (\mathbf{1}^\top \boldsymbol{\Phi}^T \mathbf{1} - \tilde{\varphi}^T n^2) = \frac{v^2 n^2}{1 - \tilde{\varphi}^T} (\bar{\varphi}^T - \tilde{\varphi}^T),$$

where $\bar{\varphi}^T$ is the overall mean kinship value, while $\tilde{\varphi}^T$ is the mean of the off-diagonal kinship values only (Eq. (7)). Note that $v^2, n^2, 1 - \tilde{\varphi}^T > 0$, so the desired result follows if $\tilde{\varphi}^T < \bar{\varphi}^T$, which is proven in Appendix C. In general it is true that $\tilde{\varphi}^T \leq \bar{\varphi}^T$, and $\tilde{\varphi}^T = \bar{\varphi}^T$ occurs if and only if the kinship matrix has the degenerate form $\boldsymbol{\Phi}^T = \bar{\varphi}^T \mathbf{J}$, which is a singular matrix not expected in practice (in this case $\boldsymbol{\Phi}^{T'}$ is a matrix full of zeroes).

Lastly, consider the weights that minimize the weighted mean kinship found in Appendix E, calculated for the unbiased kinship matrix $\boldsymbol{\Phi}^T$:

$$\mathbf{w} = \frac{(\boldsymbol{\Phi}^T)^{-1} \mathbf{1}}{\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1}}.$$

When applied to the WG-biased kinship matrix $\boldsymbol{\Phi}^{T'}$, given the previous calculation and also that $\mathbf{w}^\top \mathbf{1} = 1$, we see that

$$\begin{aligned} \mathbf{w}^\top \boldsymbol{\Phi}^{T'} \mathbf{w} &= \frac{1}{1 - \tilde{\varphi}^T} (\mathbf{w}^\top \boldsymbol{\Phi}^T \mathbf{w} - \tilde{\varphi}^T \mathbf{w}^\top \mathbf{1} \mathbf{1}^\top \mathbf{w}) \\ &= \frac{1}{1 - \tilde{\varphi}^T} \left(\frac{1}{\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1}} - \tilde{\varphi}^T \right) \end{aligned}$$

Thus, we again find that a WG-biased matrix is non-positive semidefinite if $1 / (\mathbf{1}^\top (\boldsymbol{\Phi}^T)^{-1} \mathbf{1}) \leq \tilde{\varphi}^T$, which empirically was always true in our evaluations. This \mathbf{w} is not an eigenvector (even if normalized), but empirically it had a large projection to the desired eigenvector with the negative eigenvalue.

Supplemental figures

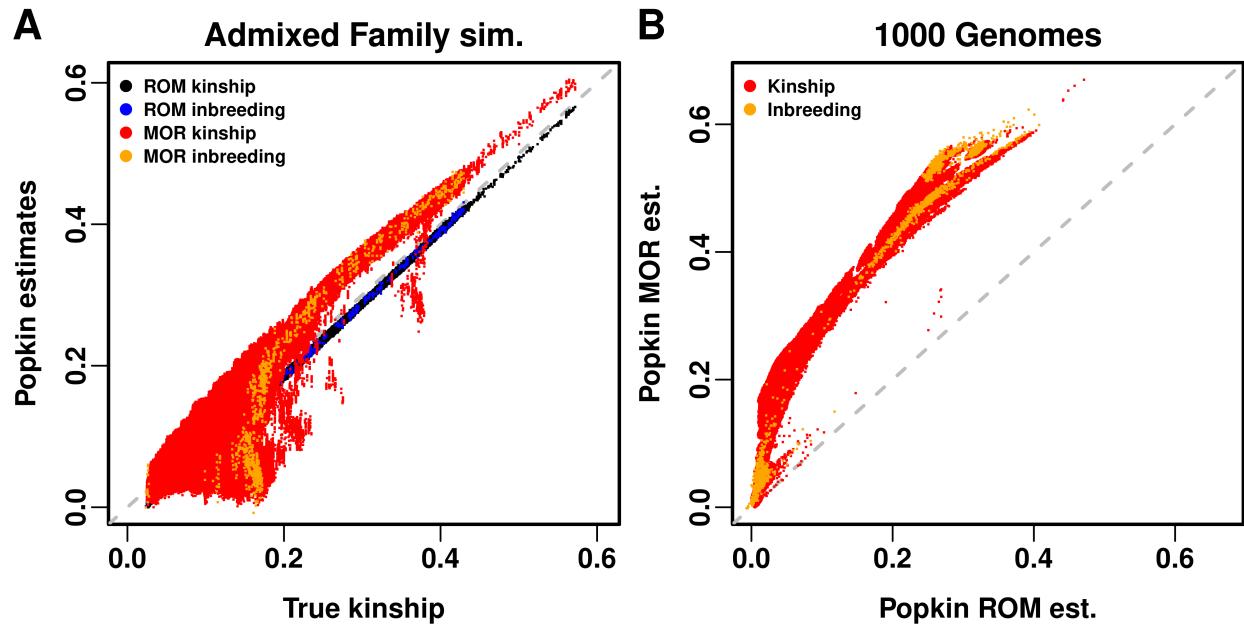


Figure S1: **Comparison of popkin ROM and MOR estimates.** Kinship (off-diagonal of matrix) and inbreeding (transformed diagonal) are plotted in different colors, which shows that their biases (if any) overlap. **A.** In admixed family simulation, both estimates are compared against true kinship. Popkin ROM has a negligible bias, due to the minimum true kinship of the simulation being slightly larger than zero. Popkin MOR has considerable biases, tending to be upward though not always. **B.** In 1000 Genomes, since true kinship is unknown, popkin ROM takes its place. Popkin MOR biases take on a similar shape as panel A.

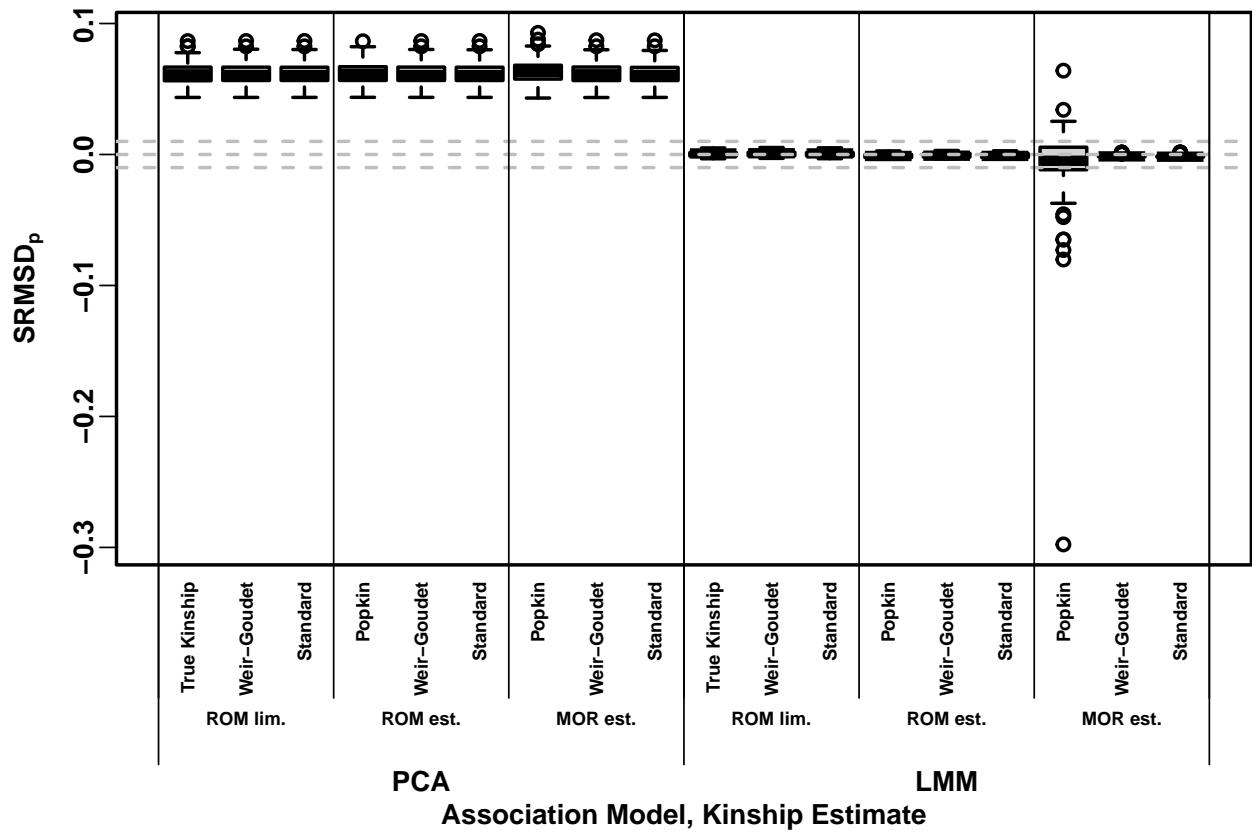


Figure S2: **Signed Root Mean Square Deviation of null p-values (SRMSD_p) on the admixed family simulation.** Same methods and simulation as Fig. 2, see that for more information. $|\text{SRMSD}_p| < 0.01$ (area between gray dashed lines) is considered calibrated. All PCA runs are miscalibrated by similar amounts, whereas most LMM runs are calibrated with few exceptions.

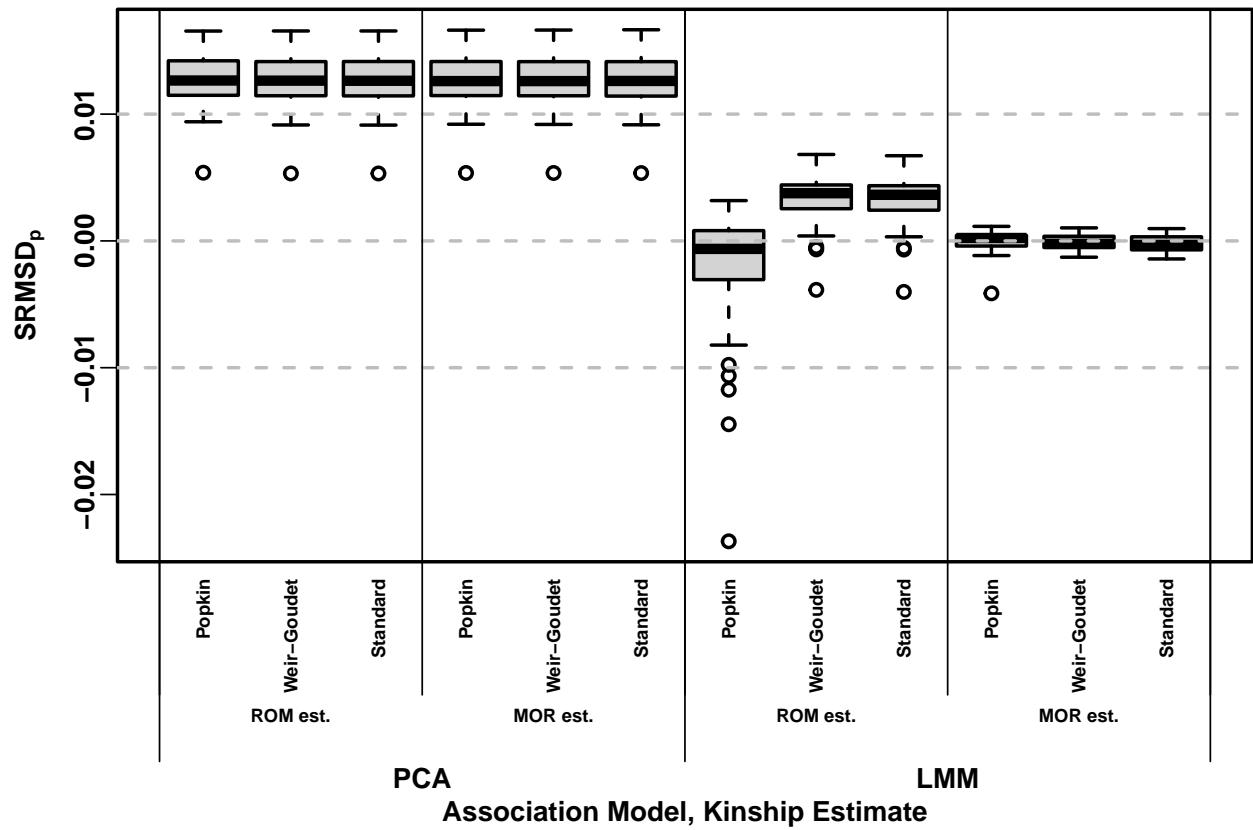


Figure S3: **Signed Root Mean Square Deviation of null p-values (SRMSD_p) on 1000 Genomes.** Same methods and simulation as Fig. 5, and y-axis statistic and conclusions of Fig. S2, see those for more information.

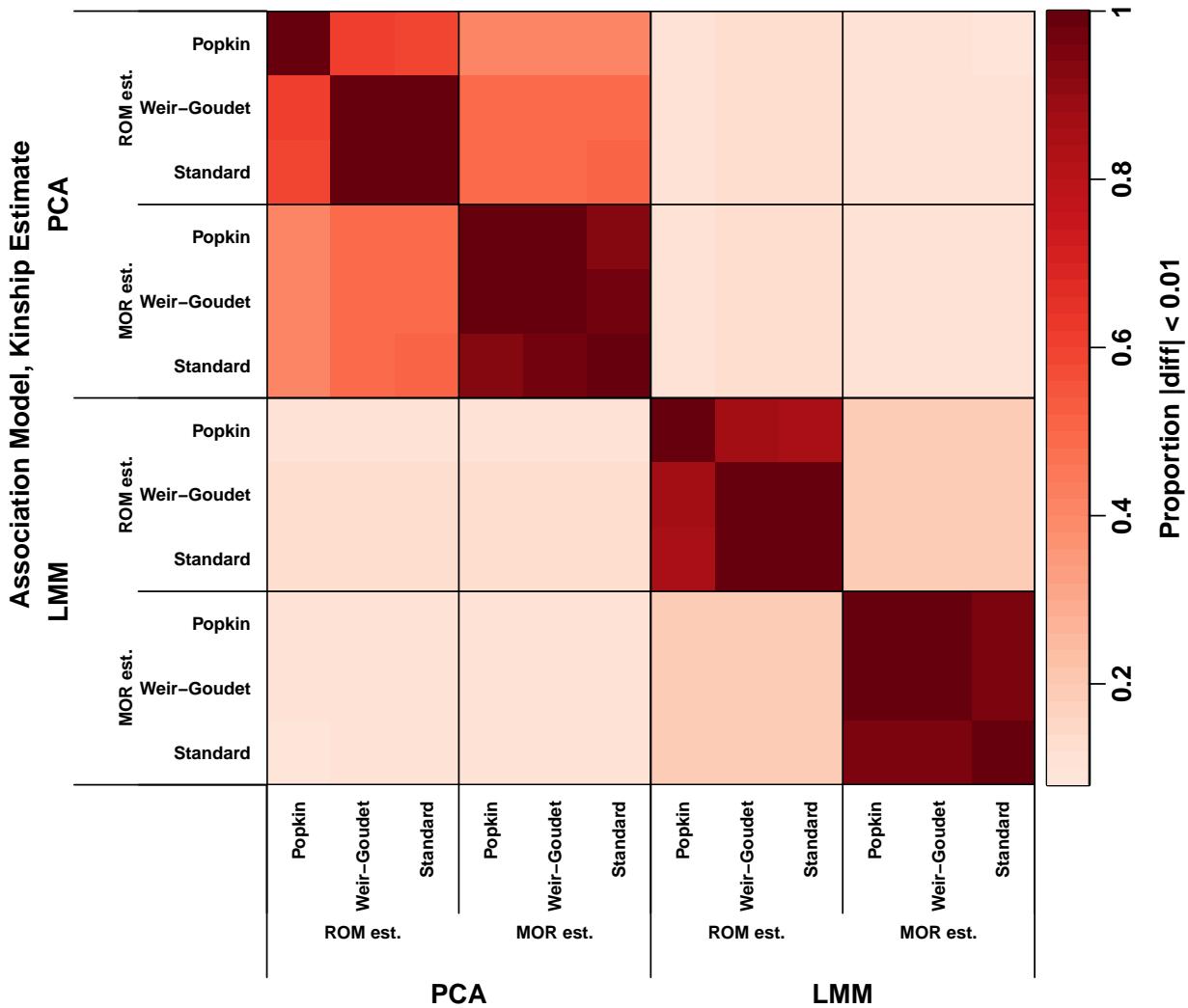


Figure S4: Approximate agreement between p-values on 1000 Genomes. See Fig. 3 for more details.

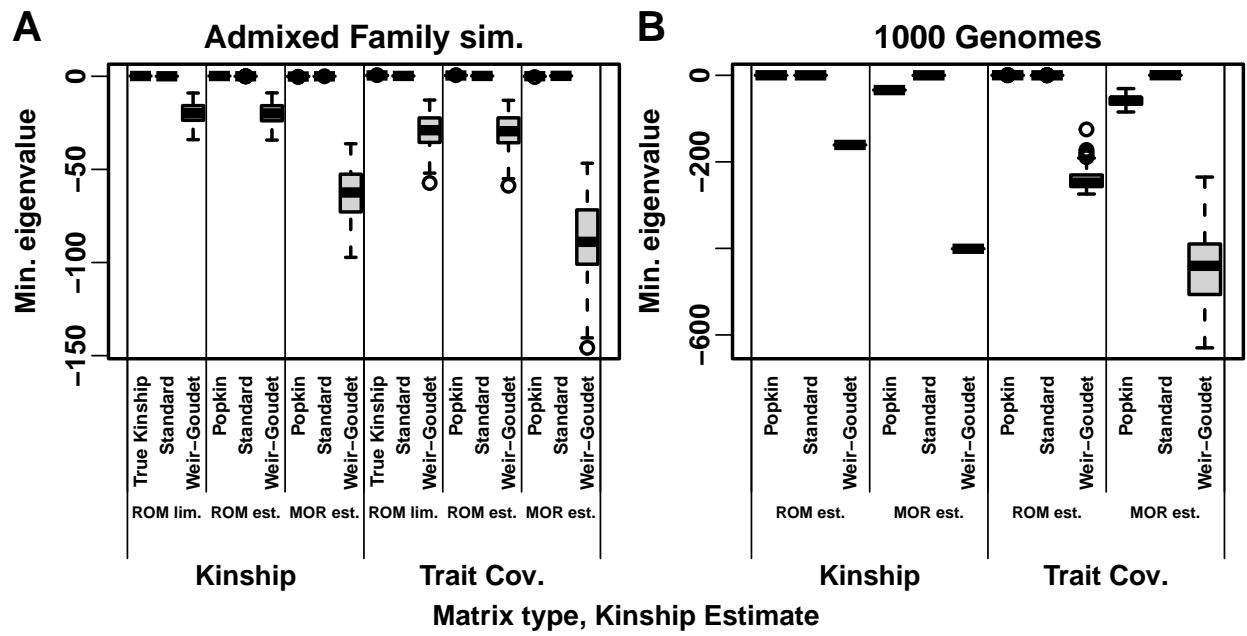


Figure S5: **Minimum eigenvalue of kinship and trait covariance matrices.** “Trait Cov.” corresponds to \mathbf{V} (see text). Each distribution is over the 100 replicates of each simulation (only 1000 Genomes “Kinship” has a single value because the genotypes did not change across replicates, but “Trait Cov.” always varies per replicate). All WG matrices has very large negative eigenvalues, and Popkin MOR has negative eigenvalues as well; in these cases a non-positive semidefinite kinship matrix always resulted in a non-positive semidefinite \mathbf{V} .

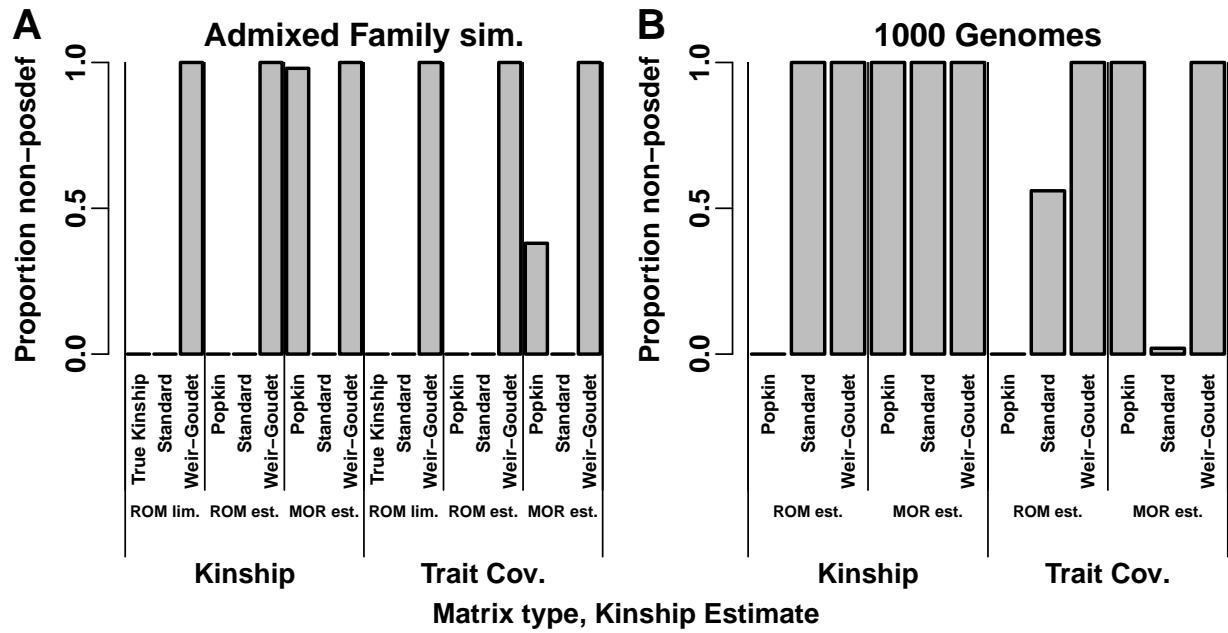


Figure S6: **Proportion of kinship and trait covariance matrices with negative eigenvalues.** An eigenvalue was considered negative if it was below -10^{-7} to allow for limited machine precision. Proportion is calculated over the 100 replicates of each simulation (only “Kinship” has a single value because the genotypes did not change across replicates, but “Trait Cov.” always varies per replicate).

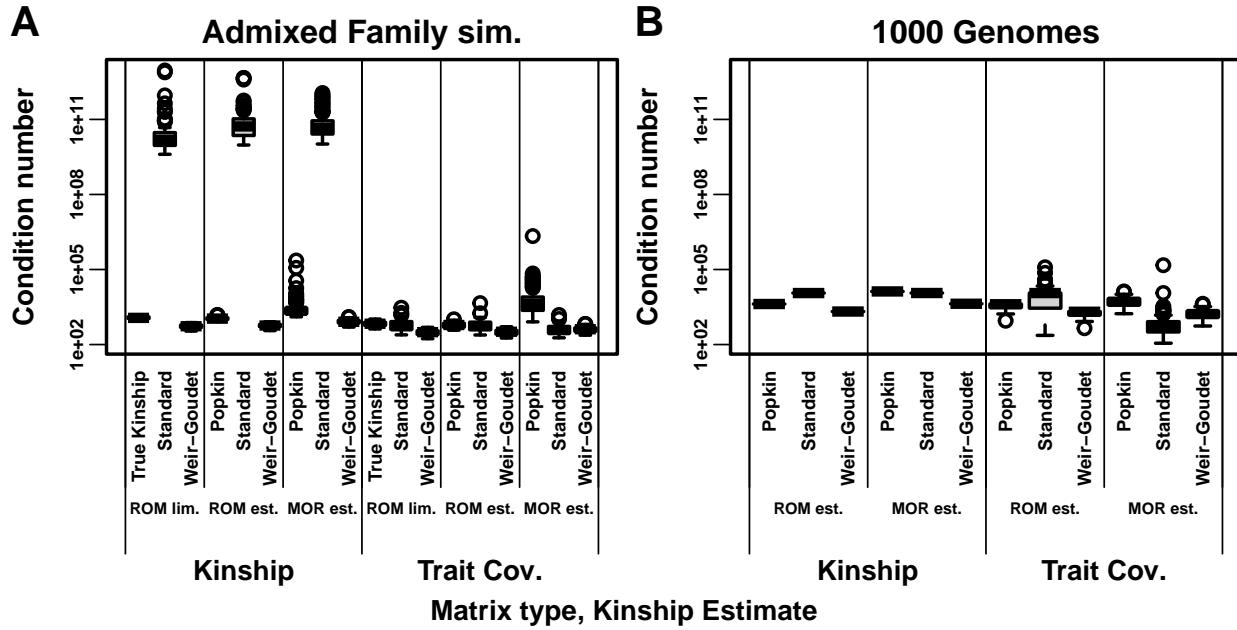


Figure S7: Condition numbers of kinship and trait covariance matrices. Larger condition numbers reflect ill-conditioned problems such as near singularity. Each distribution is over the 100 replicates of each simulation (only 1000 Genomes “Kinship” has a single value because the genotypes did not change across replicates, but “Trait Cov.” always varies per replicate).

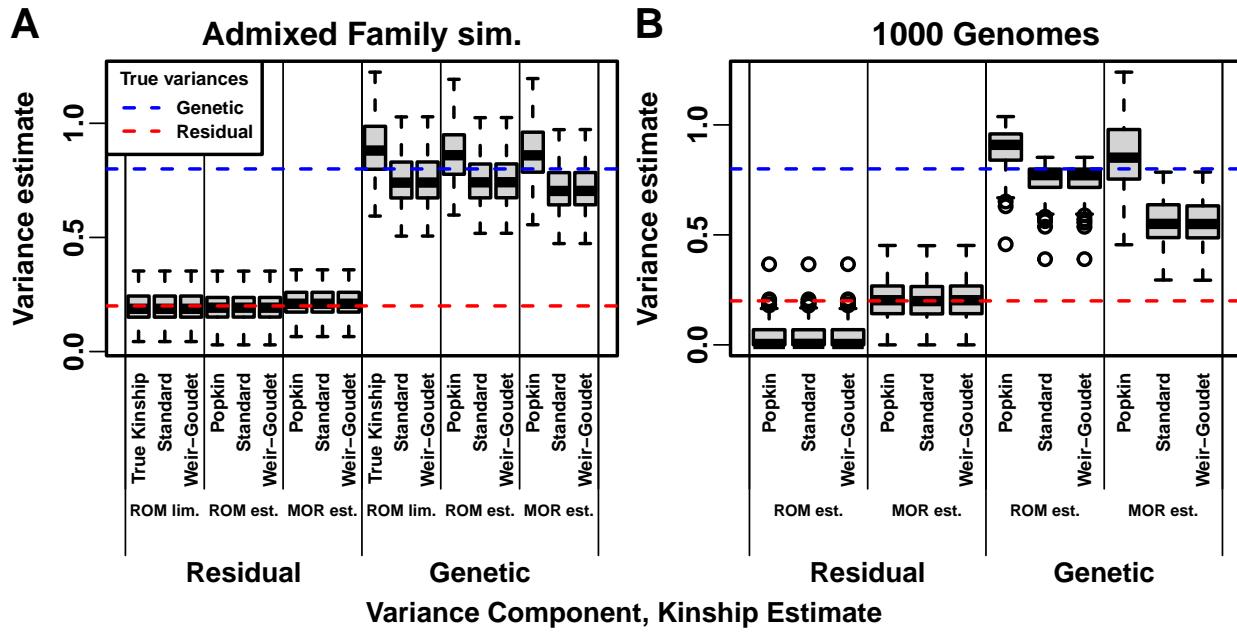


Figure S8: Variance component estimates from GCTA across evaluations. True variances (red and blue dashed lines) are the parameter values of our trait simulations. Each distribution is over the 100 replicates of each simulation.

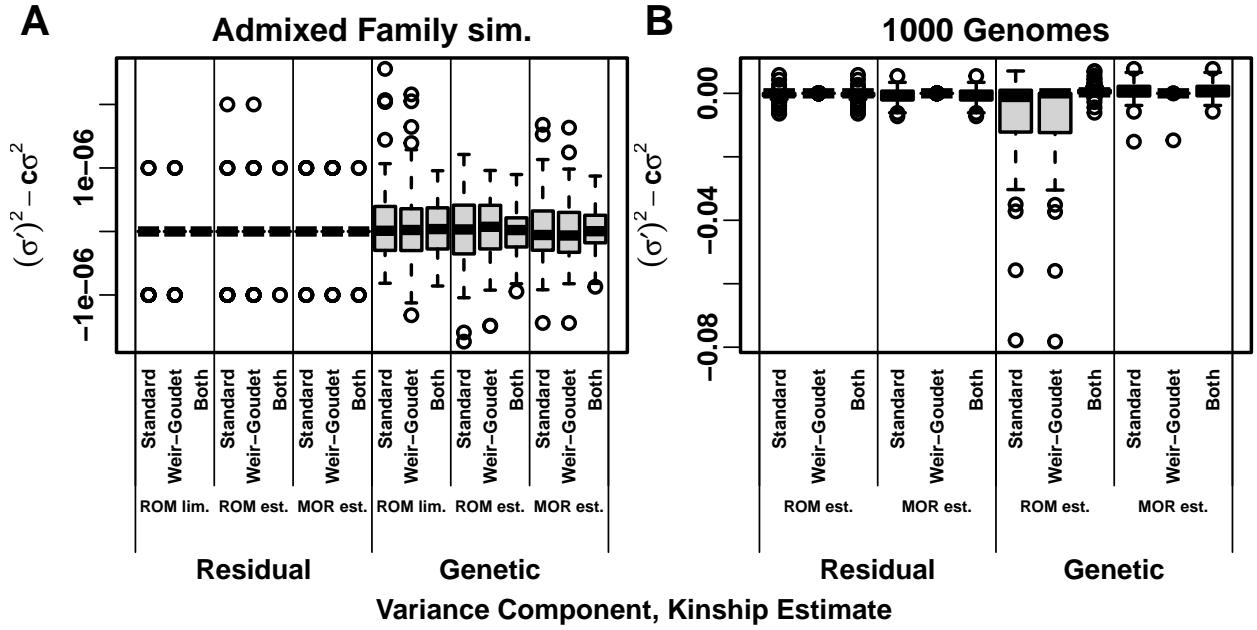


Figure S9: **Variance component prediction errors across evaluations.** In all “Residual” cases the error is $\sigma_{\epsilon}^{2\prime} - \sigma_{\epsilon}^2$ and $c = 1$, as the prediction is that they are identical (excess perfect zeroes are due to limited precision of GCTA outputs). For genetic variance, Standard uses $c = 1 - \bar{\varphi}^T$, WG uses $c = 1 - \tilde{\varphi}^T$, in both cases their estimate is $\sigma^{2\prime}$ while σ^2 is True (for ROM lim.) or Popkin (for corresponding ROM or MOR estimate). However, “Both” compares Standard ($\sigma^{2\prime}$) to WG (σ^2), and genetic variance uses $c = (1 - \bar{\varphi}^T) / (1 - \tilde{\varphi}^T)$. Each distribution is over the 100 replicates of each simulation.

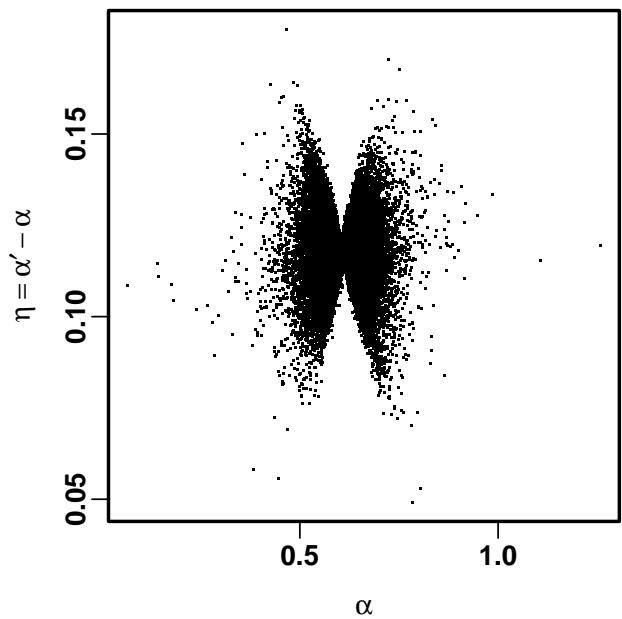


Figure S10: **Comparison of intercept coefficients from True vs Standard ROM limit kinship matrices.** Note coefficient differences are small relative to their absolute value.

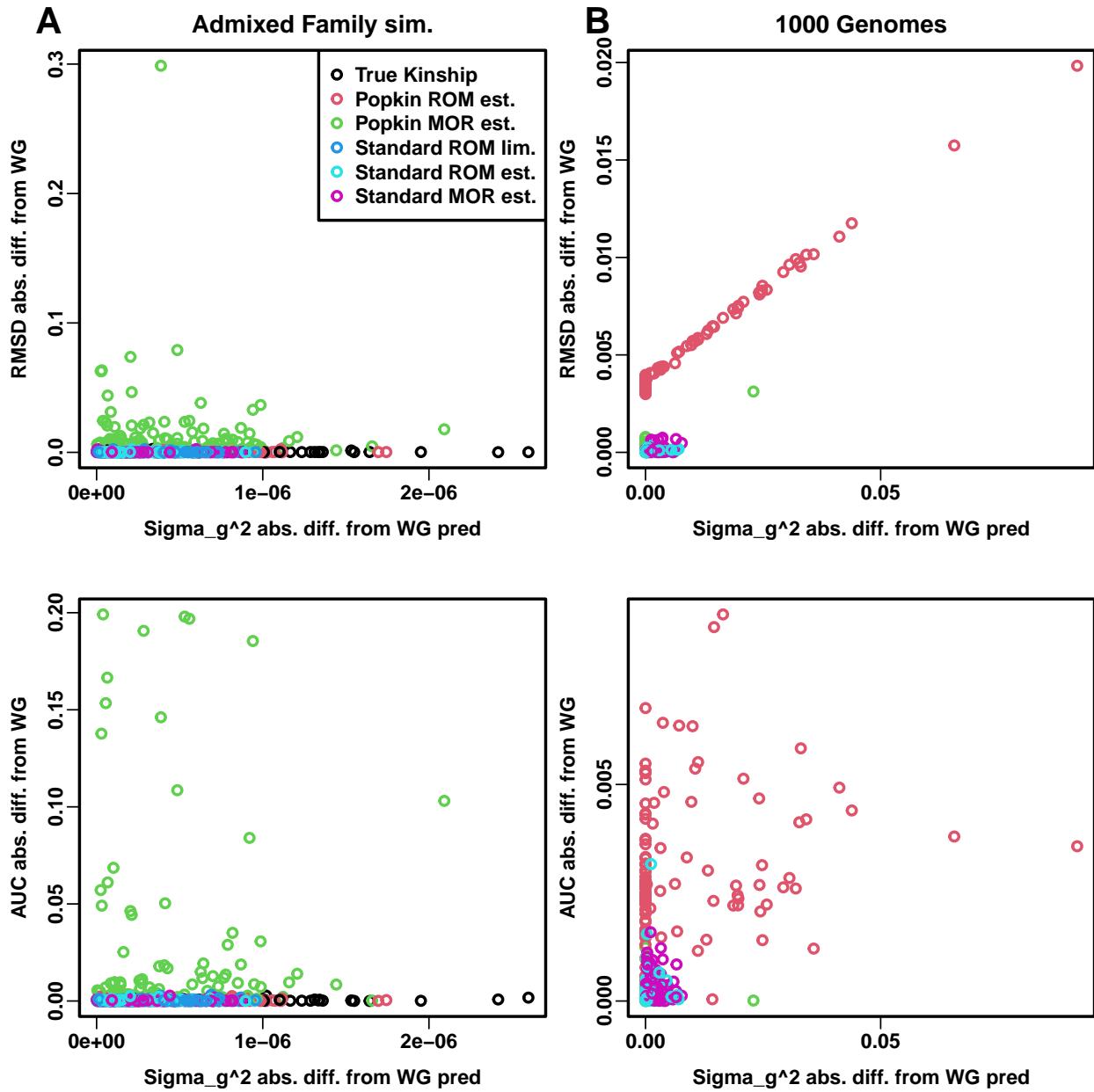


Figure S11: **AUC_{PR} and SRMSD_p prediction errors explained by variance component errors.** Genetic variance component (σ^2) absolute error was calculated based on the formulas in Fig. S9 but using WG as reference. AUC_{PR} and SRMSD_p are expected to be the same between WG, Standard, and True or Popkin (within each locus weight type). Popkin ROM prediction errors in 1000 Genomes are explained by σ^2 error, but does not explain errors in the admixed family simulation.

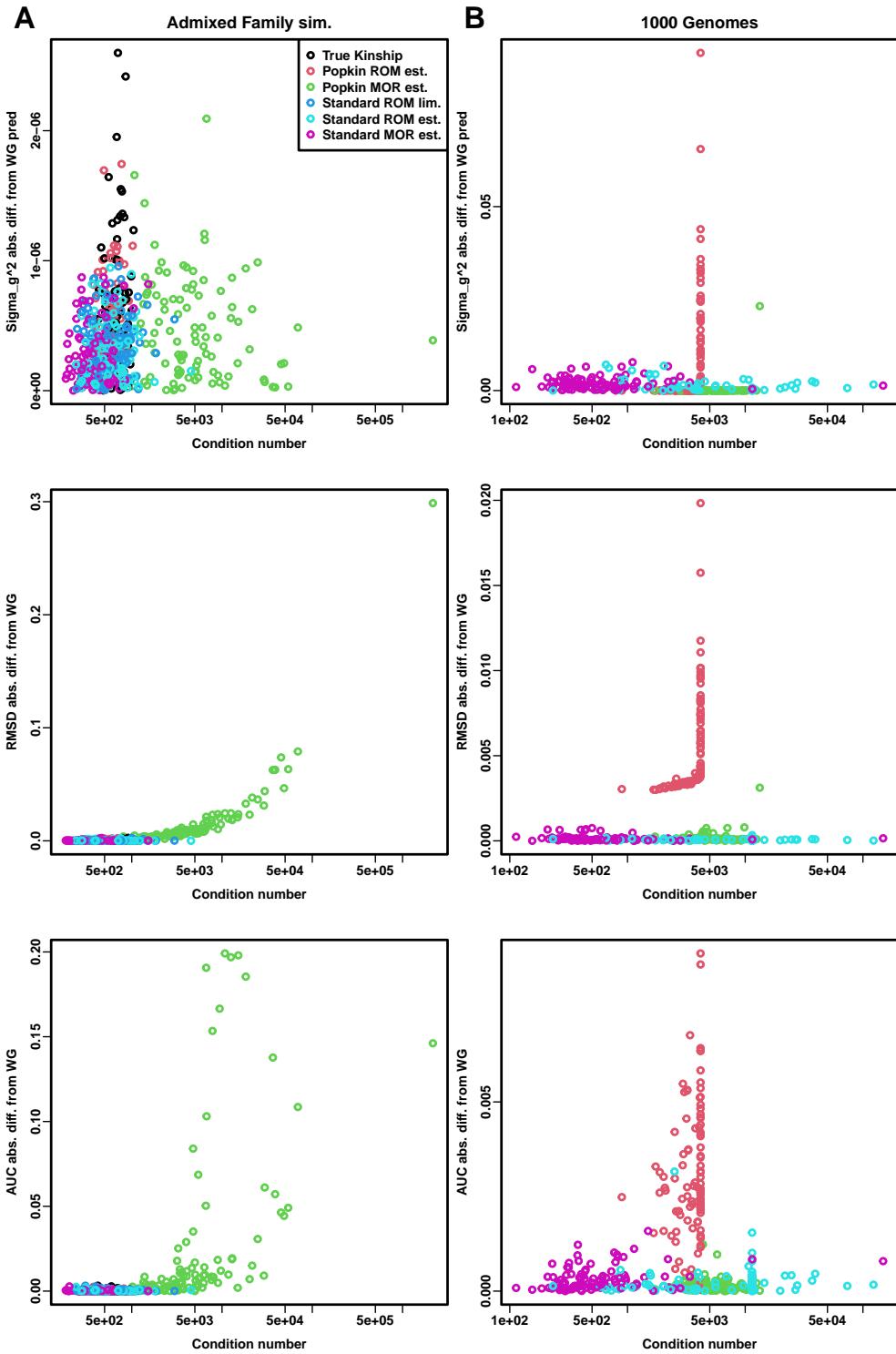


Figure S12: **AUC_{PR}** and **SRMSD_p** prediction errors explained by the condition number of \mathbf{V} . AUC_{PR} and SRMSD_p are expected to be the same between WG, Standard, and True or Popkin (within each locus weight type). Popkin MOR prediction errors in the admixed family simulation are explained by the condition number of \mathbf{V} , but does not explain errors in 1000 Genomes.