

# Limitations of principal components in quantitative genetic association models for human studies

Yiqi Yao<sup>1</sup>, Alejandro Ochoa<sup>1,2,\*</sup>

<sup>1</sup> Department of Biostatistics and Bioinformatics, Duke University, Durham, NC 27705, USA

<sup>2</sup> Duke Center for Statistical Genetics and Genomics, Duke University, Durham, NC 27705, USA

\* Corresponding author: [alejandro.ochoa@duke.edu](mailto:alejandro.ochoa@duke.edu)

## Abstract

Principal Component Analysis (PCA) and the Linear Mixed-Effects Model (LMM), sometimes in combination, are the most common modern models for genetic association. Previous PCA-LMM comparisons give mixed results and unclear guidance, and have several limitations, including not varying the number of principal components (PCs), simulating overly simple population structures, and inconsistent use of real data and power evaluations. In this work, we thoroughly evaluate PCA and LMM both with varying number of PCs in new realistic genotype and complex trait simulations including admixed families, trees, and large real multiethnic human genotype datasets (1000 Genomes Project, the Human Genome Diversity Panel, and Human Origins) with simulated traits. We find that LMM without PCs performs best in all cases, with the largest effects in the family simulation and all real human datasets. We determined that the large gaps in PCA to LMM performance on the real human datasets is due to the high-dimensional family structure stemming from large numbers of distant relatives, and not from the smaller number of highly related pairs. While it was known that PCA fails on family data, here we report a strong effect on association of cryptic family relatedness in several genetically diverse human datasets, a problem that is not avoided with the common practice of pruning high-relatedness individual pairs. Overall, this work better characterizes the severe limitations of PCA compared to LMM in modeling the complex relatedness structures present in real multiethnic human data and its impact in association studies.

**Abbreviations:** PCA: principal component analysis; PCs: principal components; LMM: linear mixed-effects model; FES: fixed effect sizes; RC: random coefficients; MAF: minor allele frequency; WGS: whole genome sequencing.

## 1 Introduction

The goal of a genetic association study is to identify loci whose genotype variation is significantly correlated to given trait. Naive association tests assume that genotypes are drawn independently from a common allele frequency. This assumption does not hold for structured populations, which includes multiethnic cohorts and admixed individuals (ancient relatedness), and for family data (recent relatedness) [1]. When insufficient approaches are applied to data with relatedness, their association statistics are miscalibrated, resulting in excess false positives and loss of power [1–3]. Therefore, many specialized approaches have been developed for genetic association under relatedness, of which PCA and LMM are the most popular.

Genetic association with PCA consists of including the top eigenvectors of the population kinship matrix as covariates in a generalized linear model [4–6]. These top eigenvectors are commonly referred to as PCs in genetics [7], the convention adopted here, but in other fields PCs denote the projections of loci onto eigenvectors [8]. The direct ancestor of PCA association is structured association, in which inferred ancestry or admixture proportions are used as regression covariates [9]. These models are deeply connected because PCs map to ancestry empirically [10, 11] and theoretically [12–15], and they work as well as global ancestry in association studies but are estimated more easily [6, 7, 10, 16]. The strength of PCA is its simplicity, which as covariates can be readily included in more complex models, such as haplotype association [17] and polygenic models [18]. However, PCA assumes that relatedness is low-dimensional, which may limit its applicability. PCA is known to be inadequate for family data [7, 19, 20], which is called “cryptic relatedness” when it is unknown to the researchers, but no other troublesome cases have been confidently identified. Recent work has focused on developing more scalable versions of the PCA algorithm [21–25]. PCA remains a popular and powerful approach for association studies.

The other dominant association model under relatedness is the LMM, which includes a random

effect parametrized by the kinship matrix. Unlike PCA, LMM does not assume that relatedness is low-dimensional, and explicitly models families via the kinship matrix. Early LMMs required kinship matrices estimated from known pedigrees or which otherwise captured recent relatedness only [16, 26]. Modern LMMs estimate kinship from genotypes using a non-parametric estimator, often referred to as a genetic relationship matrix, that captures the combined covariance due to recent family relatedness and ancestral population structure [1, 27, 28]. The classic LMM assumes a quantitative (continuous) complex trait, the focus of our work. Although case-control (binary) traits and their underlying ascertainment are theoretically a challenge [29], LMMs have been applied successfully to balanced case-control studies [1, 30] and simulations [20, 31, 32], and have been adapted for unbalanced case-control studies [33]. However, LMMs tend to be considerably slower than PCA and other models, so much effort has focused on improving their runtime and scalability [27, 30, 33–41].

An LMM variant that incorporates PCs as fixed covariates is tested thoroughly in our work. Since PCs are the top eigenvectors of the same kinship matrix estimate used in modern LMMs [1, 42], then population structure is modeled twice in an LMM with PCs. However, some previous work has found the apparent redundancy of an LMM with PCs beneficial [20, 43], while others did not [44], and the approach continues to be used [45]. Recall that early LMMs used kinship to model family relatedness only, so population structure had to be modeled separately, in practice as admixture fractions instead of PCs [16, 26].

LMM and PCA are closely related models [1, 42], so similar performance is expected particularly under low-dimensional relatedness. Direct comparisons have yielded mixed results, with several studies finding superior performance for LMM (notably from papers promoting advances in LMMs) while many others report comparable performance (Table 1). No papers find that PCA outperforms LMM decisively, although PCA occasionally performs better in isolated and artificial cases or individual measures (often with unknown significance). Previous studies were generally divided those that employed simulated versus real genotypes (only one study used both). The simulated genotype studies, which tended to have low dimensionalities and differentiation ( $F_{ST}$ ), were more likely to report ties or mixed results (6/7), whereas real genotypes tended to clearly favor LMMs

(5/7). Similarly, 6/8 papers with quantitative traits favor LMMs, whereas 5/7 papers with case-control traits gave ties or mixed results (the only factor we do not explore). Additionally, although all previous evaluations measured type I error (or proxies such as inflation factors or QQ plots), a large fraction (5/13) did not measure power (including proxies such as ROC curves), and only two used more than one number of PCs for PCA. Lastly, no consensus has emerged as to why LMM might outperform PCA or vice versa [20, 32, 42, 49], or which features of the real datasets are critical for the LMM advantage other than cryptic relatedness, resulting in unclear guidance for using PCA. Hence, our work includes real and simulated genotypes with higher dimensionalities and differentiation matching that of multiethnic human cohorts, we vary the number of PCs, and measure robust proxies for type I error control and power.

In this work, we evaluate the PCA and LMM association models under various numbers of PCs (included in LMM too). We use genotype simulations (admixture, family, and tree models) and three real datasets: the 1000 Genomes Project [50, 51], the Human Genome Diversity Panel

**Table 1: Previous PCA-LMM evaluations in the literature.**

Publication	Sim. Genotypes			Real <sup>d</sup>	Trait <sup>e</sup>	Power	PCs ( <i>r</i> )	Best
	Type <sup>a</sup>	<i>K</i> <sup>b</sup>	<i>F</i> <sub>ST</sub> <sup>c</sup>					
Zhao et al. [16]				✓	Q	✓	8	LMM
Astle and Balding [1]	I	3	0.10		CC	✓	10	Tie
Kang et al. [30]				✓	Both		2-100	LMM
Price et al. [20]	I, F	2	0.01		CC		1	Mixed
Wu et al. [31]	I, A	2-4	0.01		CC	✓	10	Mixed
Liu et al. [44]	S, A	2-3	R		Q	✓	10	Tie
Sul and Eskin [32]	I	2	0.01		CC		1	Tie
Tucker, Price, and Berger [43]	I	2	0.05	✓	Both	✓	5	Tie
Yang et al. [29]				✓	CC	✓	5	Tie
Song, Hao, and Storey [46]	S, A	2-3	R		Q		3	LMM
Loh et al. [41]				✓	Q	✓	10	LMM
Liu et al. [47]				✓	Q	✓	3-6	LMM
Sul, Martin, and Eskin [48]				✓	Q		100	LMM
This work	A, T, F	10-243	$\leq 0.25$	✓	Q	✓	0-90	LMM

<sup>a</sup>Genotype simulation types. I: Independent subpopulations; S: subpopulations (with parameters drawn from real data); A: Admixture; T: Tree; F: Family.

<sup>b</sup>Model dimensionality (number of subpopulations or ancestries)

<sup>c</sup>R: simulated parameters based on real data, *F*<sub>ST</sub> not reported.

<sup>d</sup>Evaluations using unmodified real genotypes.

<sup>e</sup>Q: quantitative; CC: case-control.

(HGDP) [52–54], and Human Origins [55–58]. We simulate quantitative traits from two models: fixed effect sizes (FES; coefficients inverse to allele frequency) that matches real data [45, 59, 60] and corresponds to high pleiotropy and strong balancing selection [61] and strong negative selection [45, 60], which are appropriate assumptions for diseases; and random coefficients (RC; independent of allele frequency) that corresponds to neutral traits [45, 61]. LMM without PCs consistently performs best, and greatly outperforms PCA in the family simulation and in all real datasets. The tree simulations do not recapitulate the real data results, suggesting that family relatedness in real data is the reason for poor PCA performance. Lastly, removing up to 4th degree relatives in the real datasets recapitulates poor PCA performance, showing that the more numerous distant relatives explain the result, and suggesting that PCA is generally not an appropriate model for real data. All together, we find that LMMs without PCs are generally a preferable association model, and present novel simulation and evaluation approaches to measure the performance of these and other genetic association approaches.

## 2 Materials and Methods

### 2.1 The complex trait model and PCA and LMM approximations

Let  $x_{ij} \in \{0, 1, 2\}$  be the genotype at the biallelic locus  $i$  for individual  $j$ , which counts the number of reference alleles. Suppose there are  $n$  individuals and  $m$  loci,  $\mathbf{X} = (x_{ij})$  is their  $m \times n$  genotype matrix, and  $\mathbf{y}$  is the length- $n$  (column) vector of individual trait values. The additive linear model for a quantitative (continuous) trait is:

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{X}^\top \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (1)$$

where  $\mathbf{1}$  is a length- $n$  vector of ones,  $\alpha$  is the scalar intercept coefficient,  $\boldsymbol{\beta}$  is the length- $m$  vector of locus coefficients,  $\boldsymbol{\epsilon}$  is a length- $n$  vector of residuals, and  $\top$  denotes matrix transposition. The residuals follow  $\epsilon_j \sim \text{Normal}(0, \sigma^2)$  independently per individual  $j$ , for some  $\sigma^2$ . For simplicity, non-genetic covariates are omitted from this model (and the PCA and LMM counterparts) but are trivial to include without changing any of our theoretical results.

The full model of Eq. (1), which has a coefficient for each of the  $m$  loci, is overdetermined in current datasets where  $m \gg n$ . The PCA and LMM models, respectively, approximate the full model fit at a single locus  $i$ :

$$\text{PCA: } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{U}_r\boldsymbol{\gamma}_r + \boldsymbol{\epsilon}, \quad (2)$$

$$\text{LMM: } \mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{s} + \boldsymbol{\epsilon}, \quad \mathbf{s} \sim \text{Normal}(\mathbf{0}, 2\sigma_s^2 \boldsymbol{\Phi}^T), \quad (3)$$

where  $\mathbf{x}_i$  is the length- $n$  vector of genotypes at locus  $i$  only,  $\beta_i$  is the locus coefficient,  $\mathbf{U}_r$  is an  $n \times r$  matrix of PCs,  $\boldsymbol{\gamma}_r$  is the length- $r$  vector of PC coefficients,  $\mathbf{s}$  is a length- $n$  vector of random effects,  $\boldsymbol{\Phi}^T = (\varphi_{jk}^T)$  is the  $n \times n$  kinship matrix conditioned on the ancestral population  $T$ , and  $\sigma_s^2$  is a variance factor (do not confuse the ancestral population superscript  $T$  with the matrix transposition symbol  $\intercal$ ). Both models condition the regression of the focal locus  $i$  on an approximation of the total polygenic effect  $\mathbf{X}^\intercal \boldsymbol{\beta}$  with the same covariance structure, which is parametrized by the kinship matrix. Under the kinship model, genotypes are random variables obeying

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

where  $p_i^T$  is the ancestral allele frequency of locus  $i$  [1, 62–64]. Assuming independent loci, the covariance of the polygenic effect is

$$\text{Cov}(\mathbf{X}^\intercal \boldsymbol{\beta}) = 2\sigma_s^2 \boldsymbol{\Phi}^T, \quad \sigma_s^2 = \sum_{i=1}^m 2p_i^T(1 - p_i^T)\beta_i^2,$$

which is readily modeled by the LMM random effect  $\mathbf{s}$ . (The difference in mean is absorbed by the intercept.) Alternatively, consider the eigendecomposition of the kinship matrix  $\boldsymbol{\Phi}^T = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\intercal$  where  $\mathbf{U}$  is the  $n \times n$  eigenvector matrix and  $\boldsymbol{\Lambda}$  is the  $n \times n$  diagonal matrix of eigenvalues. The random effect can be written as

$$\mathbf{s} = \mathbf{U}\boldsymbol{\gamma}_{\text{LMM}}, \quad \boldsymbol{\gamma}_{\text{LMM}} \sim \text{Normal}(\mathbf{0}, 2\sigma_s^2 \boldsymbol{\Lambda}),$$

which follows from the affine transformation property of multivariate normal distributions. There-

fore, the PCA term  $\mathbf{U}_r \boldsymbol{\gamma}_r$  can be derived from the above equation under the additional assumption that the kinship matrix has dimensionality  $r$  and the coefficients  $\boldsymbol{\gamma}_r$  are fit without constraints. In contrast, the LMM uses all eigenvectors, while effectively shrinking their coefficients  $\boldsymbol{\gamma}_{\text{LMM}}$  as all random effects models do, although these parameters are marginalized [1, 42]. PCA has more parameters than LMM, so it may overfit more: ignoring the shared terms in Eqs. (2) and (3), PCA fits  $r$  parameters (length of  $\boldsymbol{\gamma}$ ), whereas LMMs fit only one ( $\sigma_s^2$ ).

In practice, the kinship matrix used for PCA and LMM is estimated with variations of a method-of-moments formula applied to standardized genotypes  $\mathbf{X}_S$ , which is derived from Eq. (4):

$$\mathbf{X}_S = \left( \frac{x_{ij} - 2\hat{p}_i^T}{\sqrt{4\hat{p}_i^T (1 - \hat{p}_i^T)}} \right), \quad \hat{\Phi}^T = \frac{1}{m} \mathbf{X}_S^\top \mathbf{X}_S, \quad (5)$$

where the unknown  $p_i^T$  is estimated by  $\hat{p}_i^T = \frac{1}{2n} \sum_{j=1}^n x_{ij}$  [5, 27, 30, 33, 37, 39, 41, 48]. However, this kinship estimator has a complex bias that differs for every individual pair, which arises due to the use of this estimated  $\hat{p}_i^T$  [28, 65]. Nevertheless, in PCA and LMM these biased estimates perform as well as unbiased ones, an observation that will be explored in future work (data not shown).

We selected fast and robust software implementing the basic PCA and LMM models. PCA association was performed with `plink2` [66]. The quantitative trait association model is a linear regression with covariates, evaluated using the t-test. PCs were calculated with `plink2`, which equal the top eigenvectors of Eq. (5) after removing loci with MAF < 0.1.

LMM association was performed using GCTA [29, 37]. Kinship equals Eq. (5) except self-kinship uses a different formula. PCs were calculated using GCTA from its kinship estimate. Association significance is evaluated with a score test. GCTA with large numbers of PCs (small simulation only) had convergence and singularity errors in some replicates, which were treated as missing data.

## 2.2 Simulations

Every simulation was replicated 50 times, drawing anew all genotypes (except for real datasets) and traits. Below we use the notation  $f_A^B$  for the inbreeding coefficient of a subpopulation  $A$  from

another subpopulation  $B$  ancestral to  $A$ . In the special case of the *total* inbreeding of  $A$ ,  $f_A^T$ ,  $T$  is an overall ancestral population (ancestral to every individual under consideration, such as the most recent common ancestor (MRCA) population).

### 2.2.1 Genotype simulation from the admixture model

The basic admixture model is as described previously [28] and is implemented in the R package `bnpstd`. Large and Family have  $n = 1,000$  individuals, while Small has  $n = 100$ . The number of loci is  $m = 100,000$ . Individuals are admixed from  $K = 10$  intermediate subpopulations, or ancestries. Each subpopulation  $S_u$  ( $u \in \{1, \dots, K\}$ ) is at coordinate  $u$  and has an inbreeding coefficient  $f_{S_u}^T = u\tau$  for some  $\tau$ . Ancestry proportions  $q_{ju}$  for individual  $j$  and  $S_u$  arise from a random walk with spread  $\sigma$  on the 1D geography, and  $\tau$  and  $\sigma$  are fit to give  $F_{ST} = 0.1$  and mean kinship  $\bar{\theta}^T = 0.5F_{ST}$  for the admixed individuals [28]. Random ancestral allele frequencies  $p_i^T$ , subpopulation allele frequencies  $p_i^{S_u}$ , individual-specific allele frequencies  $\pi_{ij}$ , 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 [67] with mean  $p_i^T$  and variance  $p_i^T (1 - p_i^T) f_{S_u}^T$ . Fixed loci ( $i$  where  $x_{ij} = 0$  for all  $j$ , or  $x_{ij} = 2$  for all  $j$ ) are drawn again from the model, starting from  $p_i^T$ , iterating until no loci are fixed. Each replicate draws a genotypes starting from  $p_i^T$ .

As a brief aside, we prove that global ancestry proportions as covariates is equivalent in expectation to using PCs under the admixture model. Note that the latent space of  $\mathbf{X}$ , given by  $(\pi_{ij})$ , has  $K$  dimensions (number of columns of  $\mathbf{Q} = (q_{ju})$ ), so the top  $K$  PCs span this space. Since associations include an intercept term ( $\mathbf{1}\alpha$  in Eq. (2)), estimated PCs are orthogonal to  $\mathbf{1}$  (note

$\hat{\Phi}^T \mathbf{1} = \mathbf{0}$  because  $\mathbf{X}_S \mathbf{1} = \mathbf{0}$ ), and the sum of rows of  $\mathbf{Q}$  sums to one, then only  $K - 1$  PCs (plus intercept) are needed to span the latent space of this admixture model.

### 2.2.2 Genotype simulation from random admixed families

We simulated a pedigree with admixed founders, no close relative pairings, assortative mating based on a 1D geography (to preserve admixture structure), random family sizes, and arbitrary numbers of generations (20 here). This simulation is implemented in the R package `simfam`. Generations are drawn iteratively. Generation 1 has  $n = 1000$  individuals from the above admixture simulation ordered by their 1D geography. Local kinship measures pedigree relatedness; in the first generation, everybody is locally unrelated and outbred. Individuals are randomly assigned sex. In the next generation, individuals are paired iteratively, removing random males from the pool of available males and pairing them with the nearest available female with local kinship  $< 1/4^3$  (stay unpaired if there are no matches), until there are no more available males or females. Let  $n = 1000$  be the desired population size,  $n_m = 1$  the minimum number of children and  $n_f$  the number of families (paired parents) in the current generation, then the number of additional children (beyond the minimum) is drawn from Poisson( $n/n_f - n_m$ ). Let  $\delta$  be the difference between desired and current population sizes. If  $\delta > 0$ , then  $\delta$  random families are incremented by 1. If  $\delta < 0$ , then  $|\delta|$  random families with at least  $n_m + 1$  children are decremented by 1. If  $|\delta|$  exceeds the number of families, all families are incremented or decremented as needed and the process is iterated. Children are assigned sex randomly, and are reordered by the average coordinate of their parents. Children draw alleles from their parents independently per locus. A new random pedigree is drawn for each replicate, as well as new founder genotypes from the admixture model.

### 2.2.3 Genotype simulation from a tree model

This model draws subpopulations allele frequencies from a hierarchical model parametrized by a tree, which is also implemented in `bnpst` and relies on `ape` for general tree data structures and methods [68]. The ancestral population  $T$  is the root, and each node is a subpopulation  $S_w$  indexed arbitrarily. Each edge between  $S_w$  and its parent population  $P_w$  has an inbreeding coefficient  $f_{S_w}^{P_w}$ .

$p_i^T$  are drawn from a given distribution (constructed to mimic each real dataset in Appendix A). Given the allele frequencies  $p_i^{P_w}$  of the parent population,  $S_w$ 's allele frequencies are drawn from:

$$p_i^{S_w} | p_i^{P_w} \sim \text{Beta} \left( p_i^{P_w} \left( \frac{1}{f_{S_w}^{P_w}} - 1 \right), (1 - p_i^{P_w}) \left( \frac{1}{f_{S_w}^{P_w}} - 1 \right) \right).$$

Individuals  $j$  in  $S_w$  draw genotypes from its allele frequency:  $x_{ij} | p_i^{S_w} \sim \text{Binomial}(2, p_i^{S_w})$ . Loci with MAF < 0.01 are drawn again starting from the  $p_i^T$  distribution, iterating until no such loci remain.

#### 2.2.4 Fitting tree to real data

We developed new methods to fit trees to real data based on unbiased kinship estimates from `popkin`, implemented in `bnpstd`. A tree with given inbreeding edges  $f_{S_w}^{P_w}$  gives rise to a coancestry matrix  $\vartheta_{uv}^T$  for a subpopulation pair  $(S_u, S_v)$ , and the goal is to recover the inbreeding edges from coancestry estimates. Coancestry values are total inbreeding coefficients of the MRCA population of each subpopulation pair. Therefore, we calculate  $f_{S_w}^T$  for every  $S_w$  recursively from the root as follows. Nodes with parent  $P_w = T$  are already as desired. Given  $f_{P_w}^T$ , the desired  $f_{S_w}^T$  is calculated via the additive edge  $\delta_w$  [28]:

$$f_{S_w}^T = f_{P_w}^T + \delta_w, \quad \delta_w = f_{S_w}^{P_w} (1 - f_{P_w}^T). \quad (6)$$

These  $\delta_w \geq 0$  because  $0 \leq f_{S_w}^{P_w}, f_{P_w}^T \leq 1$  for every  $w$ . Inbreeding edges can be recovered from additive edges:  $f_{S_w}^{P_w} = \delta_w / (1 - f_{P_w}^T)$ . Overall, coancestry values are sums of  $\delta_w$  over common ancestor nodes,

$$\vartheta_{uv}^T = \sum_w \delta_w I_w(u, v), \quad (7)$$

where the sum includes all  $w$ , and  $I_w(u, v)$  equals 1 if  $S_w$  is a common ancestor of  $S_u, S_v$ , 0 otherwise. Note that  $I_w(u, v)$  reflects tree topology and  $\delta_w$  edge values.

To estimate population-level coancestry, first kinship ( $\phi_{jk}^T$ ) is estimated using `popkin` [28]. In-

dividual coancestry ( $\hat{\theta}_{jk}^T$ ) is estimated from kinship using

$$\hat{\theta}_{jk}^T = \begin{cases} \hat{\varphi}_{jk}^T & \text{if } k \neq j, \\ \hat{f}_j^T = 2\hat{\varphi}_{jj}^T - 1 & \text{if } k = j. \end{cases} \quad (8)$$

Lastly, coancestry  $\hat{\vartheta}_{uv}^T$  between subpopulations are averages of individual coancestry values:

$$\hat{\vartheta}_{uv}^T = \frac{1}{|S_u||S_v|} \sum_{j \in S_u} \sum_{k \in S_v} \hat{\theta}_{jk}^T.$$

Topology is estimated with hierarchical clustering using the weighted pair group method with arithmetic mean [69], with distance function  $d(S_u, S_v) = \max \left\{ \hat{\vartheta}_{uv}^T \right\} - \hat{\vartheta}_{uv}^T$ , which succeeds due to the monotonic relationship between node depth and coancestry (Eq. (7)). This algorithm recovers the true topology from the true coancestry values, and performs well for estimates from genotypes.

To estimate tree edge lengths, first  $\delta_w$  are estimated from  $\hat{\vartheta}_{uv}^T$  and the topology using Eq. (7) and non-negative least squares linear regression [70] (implemented in `nnls` [71]) to yield non-negative  $\delta_w$ , and  $f_{S_w}^{P_w}$  are calculated from  $\delta_w$  by reversing Eq. (6). To account for small biases in coancestry estimation, an intercept term  $\delta_0$  is included ( $I_0(u, v) = 1$  for all  $u, v$ ), and when converting  $\delta_w$  to  $f_{S_w}^{P_w}$ ,  $\delta_0$  is treated as an additional edge to the root, but is ignored when drawing allele frequencies from the tree.

### 2.2.5 Trait Simulation

Traits are simulated from the quantitative trait model of Eq. (1), with novel bias corrections for simulating the desired heritability from real data relying on the unbiased kinship estimator `popkin` [28]. This simulation is implemented in the R package `simtrait`. All simulations have a narrow-sense heritability of  $h^2 = 0.8$  and  $\epsilon_j \sim \text{Normal}(0, 1 - h^2)$ . To balance power while varying  $n$ , the number of causal loci is  $m_1 = n/10$ . The set of causal loci  $C$  is drawn anew for each replicate, from loci with MAF  $\geq 0.01$  to avoid rare causal variants (inappropriate for PCA and LMM). Letting  $v_i^T = p_i^T (1 - p_i^T)$ , the effect size of locus  $i$  equals  $2v_i^T \beta_i^2$ , its contribution of the trait variance [72].

Under the *fixed effect sizes* (FES) model, initial causal coefficients are

$$\beta_i = \frac{1}{\sqrt{2v_i^T}}$$

for known  $p_i^T$ ; otherwise  $v_i^T$  is replaced by the unbiased estimator [28]  $\hat{v}_i^T = \hat{p}_i^T (1 - \hat{p}_i^T) / (1 - \bar{\varphi}^T)$ , where  $\bar{\varphi}^T$  is the mean kinship estimated with `popkin`. Each causal locus is multiplied by -1 with probability 0.5. Alternatively, under the *random coefficients* (RC) model, initial causal coefficients are drawn independently from  $\beta_i \sim \text{Normal}(0, 1)$ . For both models, the initial genetic variance is  $\sigma_0^2 = \sum_{i \in C} 2v_i^T \beta_i^2$ , replacing  $v_i^T$  with  $\hat{v}_i^T$  for unknown  $p_i^T$  (so  $\sigma_0^2$  is an unbiased estimate), so we multiply every initial  $\beta_i$  by  $\frac{h}{\sigma_0}$  to have the desired heritability. Lastly, for known  $p_i^T$ , the intercept coefficient is  $\alpha = -\sum_{i \in C} 2p_i^T \beta_i$ . When  $p_i^T$  are unknown,  $\hat{p}_i^T$  should not replace  $p_i^T$  since that distorts the trait covariance (for the same reason the standard kinship estimator in Eq. (5) is biased), which is avoided with

$$\alpha = -\frac{2}{m_1} \left( \sum_{i \in C} \hat{p}_i^T \right) \left( \sum_{i \in C} \beta_i \right).$$

### 2.3 Real human genotype datasets

The three datasets were processed as before [65] (summarized below), except with an additional filter so loci are in approximate linkage equilibrium and rare variants are removed. All processing was performed with `plink2` [66], and analysis was uniquely enabled by the R packages `BEDMatrix` [73] and `genio`. Each dataset groups individuals in a two-level hierarchy, which we call continental and fine-grained subpopulations, respectively. Final dataset sizes are in Table 2.

We obtained the full (including non-public) Human Origins by contacting the authors and agreeing to their usage restrictions. The Pacific data [58] was obtained separately from the rest [56, 57], and datasets were merged using the intersection of loci. We removed ancient individuals, and individuals from singleton and non-native subpopulations. Non-autosomal loci were removed. Our analysis of the WGS version of HGDP [54] was restricted to autosomal biallelic SNP loci with filter “PASS”. Our analysis of the high-coverage NYGC version of 1000 Genomes [74] was restricted to autosomal biallelic SNP loci with filter “PASS”.

Since our evaluations assume uncorrelated loci, we filtered each dataset with `plink2` using parameters “`--indep-pairwise 1000kb 0.3`”, which iteratively removes loci that have a greater than 0.3 correlation coefficient with another locus that is within 1000kb, stopping until no such loci remain. Since all real datasets have numerous rare variants, while PCA and LMM are not able to detect associations involving rare variants, we removed all loci with  $\text{MAF} < 0.01$ . Kinship dimensionality and eigenvalues were calculated from `popkin` kinship estimates. Eigenvalues were assigned p-values with `twstats` of the Eigensoft package [7], and dimensionality was estimated as the largest number of consecutive eigenvalue from the start that all satisfy  $p < 0.01$  (p-values did not increase monotonically). For the evaluation with close relatives removed, each dataset was filtered with `plink2` with option “`--king-cutoff`” with cutoff  $0.02209709 (= 2^{-11/2})$  for removing up to 4th degree relatives using KING-robust [75], and  $\text{MAF} < 0.01$  is reapplied (Table S1).

## 2.4 Evaluation of performance

All approaches are evaluated in two orthogonal dimensions:  $\text{SRMSD}_p$  quantifies p-value uniformity, and  $\text{AUC}_{\text{PR}}$  measures causal locus classification performance and reflects power while ranking mis-calibrated models fairly. These measures are more robust alternatives to previous measures from the literature (see Appendix B), and are implemented in `simtrait`.

P-values for continuous test statistics have a uniform distribution when the null hypothesis holds, a crucial assumption for type I error and FDR control [76, 77]. We use the Signed Root Mean Square Deviation ( $\text{SRMSD}_p$ ) to measure the difference between the observed null p-value quantiles and the expected uniform quantiles:

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

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

conservative p-values, and  $\text{SRMSD}_p < 0$  are conservative p-values. The maximum  $\text{SRMSD}_p$  is achieved when all p-values are zero (the limit of anti-conservative p-values), which for infinite loci approaches

$$\text{SRMSD}_p \rightarrow \sqrt{\int_0^1 u^2 du} = \frac{1}{\sqrt{3}} \approx 0.577.$$

The same worst-case value (with negative sign) occurs for all p-values of 1.

Precision and recall are standard performance measures for binary classifiers that do not require calibrated p-values [78]. Given the total numbers of true positives (TP), false positives (FP) and false negatives (FN) at some threshold or parameter  $t$ , precision and recall are

$$\begin{aligned}\text{Precision}(t) &= \frac{\text{TP}(t)}{\text{TP}(t) + \text{FP}(t)}, \\ \text{Recall}(t) &= \frac{\text{TP}(t)}{\text{TP}(t) + \text{FN}(t)}.\end{aligned}$$

Precision and Recall trace a curve as  $t$  is varied, and the area under this curve is  $\text{AUC}_{\text{PR}}$ . We use the R package `PRROC` to integrate the correct non-linear piecewise function when interpolating between points. A model obtains the maximum  $\text{AUC}_{\text{PR}} = 1$  if there is a  $t$  that classifies all loci perfectly. In contrast, the worst models, which classify at random, have an expected precision ( $= \text{AUC}_{\text{PR}}$ ) equal to the overall proportion of causal loci:  $\frac{m_1}{m}$ .

## 3 Results

### 3.1 Overview of evaluations

We use three real genotype datasets and simulated genotypes from six population structure scenarios to cover various features of interest (Table 2). We introduce them in sets of three, as they appear in the rest of our results. Population kinship matrices, which combine population and family relatedness, are estimated without bias using `popkin` [28] (Fig. 1). The first set of three simulated genotypes are based on an admixture model with 10 ancestries (Fig. 1A) [14, 28, 79]. The “large” version (1000 individuals) illustrates asymptotic performance, while the “small” simulation (100 individuals) illustrates model overfitting. The “family” simulation has admixed founders and draws

a 20-generation random pedigree with assortative mating, resulting in a complex joint family and ancestry structure in the last generation (Fig. 1B). The second set of three are the real human datasets representing global human diversity: Human Origins (Fig. 1D), HGDP (Fig. 1G), and 1000 Genomes (Fig. 1J), which are enriched for small minor allele frequencies even after  $\text{MAF} < 1\%$  filter (Fig. 1C). Last are tree simulations (Fig. 1F,I,L) fit to the kinship (Fig. 1E,H,K) and MAF (Fig. 1C) of each real human dataset, which by design do not have family structure.

All traits in this work are simulated. We repeated all evaluations on two additive quantitative trait models, *fixed effect sizes* (FES) and *random coefficients* (RC), which differ in how causal coefficients are constructed. The FES model captures the rough inverse relationship between coefficient and minor allele frequency that arises under strong negative and balancing selection and has been observed in numerous diseases and other traits [45, 59–61], so it is the focus of our results. The

Table 2: **Features of simulated and real human genotype datasets.**

Dataset	Type	Loci ( $m$ )	Ind. ( $n$ )	Subpops. <sup>a</sup> ( $K$ )	Causal loci <sup>b</sup> ( $m_1$ )	$F_{ST}$ <sup>c</sup>
Admix. Large sim.	Admix.	100,000	1000	10	100	0.1
Admix. Small sim.	Admix.	100,000	100	10	10	0.1
Admix. Family sim.	Admix.+Pedig.	100,000	1000	10	100	0.1
Human Origins	Real	190,394	2922	11-243	292	0.28
HGDP	Real	924,892	929	7-54	93	0.28
1000 Genomes	Real	1,111,266	2504	5-26	250	0.22
Human Origins sim.	Tree	190,394	2922	243	292	0.23
HGDP sim.	Tree	924,892	929	54	93	0.25
1000 Genomes sim.	Tree	1,111,266	2504	26	250	0.21

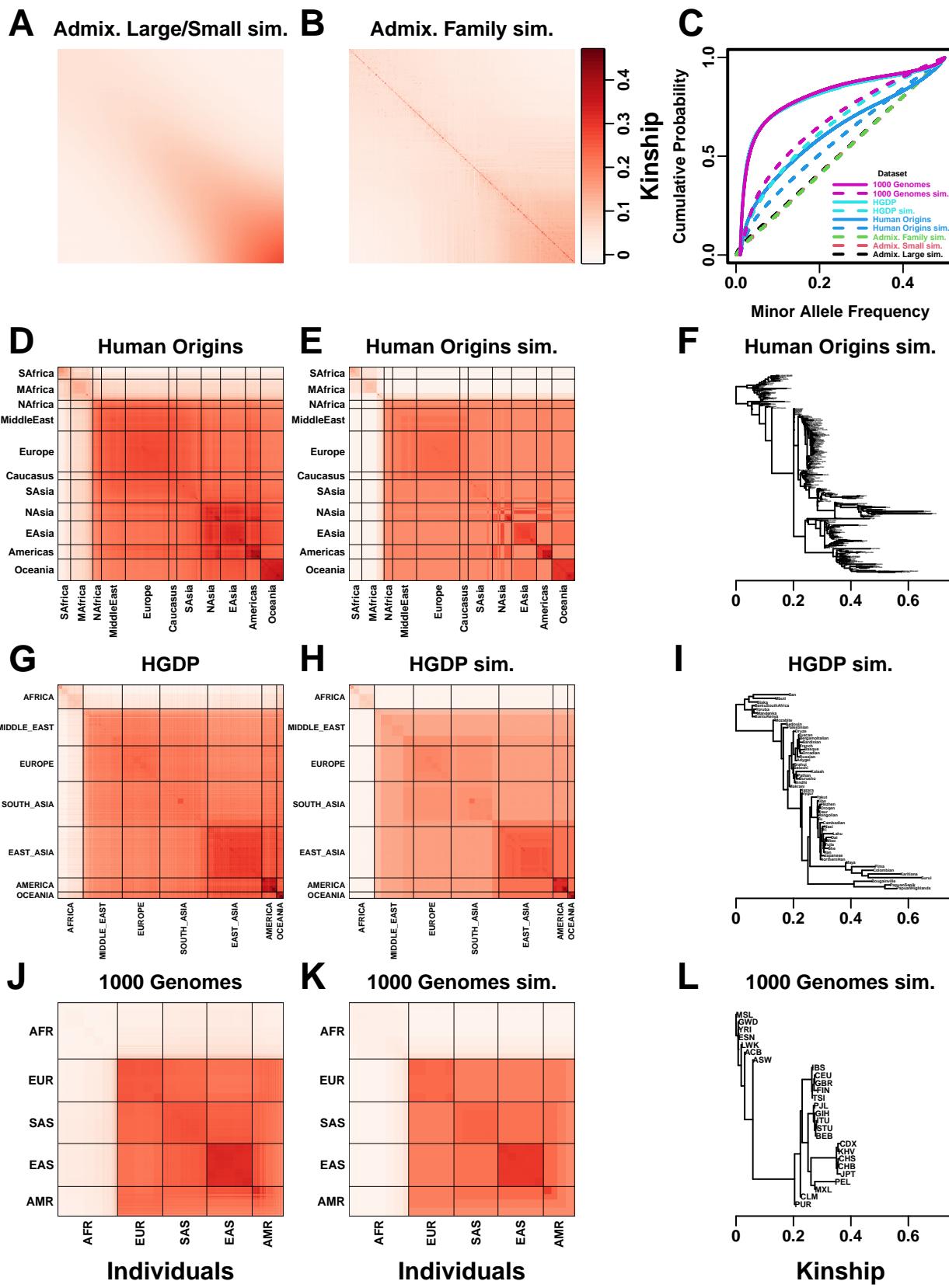
<sup>a</sup>For admixed family, ignores dimensionality of 20 generation pedigree structure. For real datasets, lower range is continental subpopulations, upper range is number of fine-grained subpopulations.

<sup>b</sup> $m_1 = n/10$  to balance power across datasets.

<sup>c</sup>Model parameter for simulations, estimated value on real datasets.

---

Figure 1 (following page): **Population structures of simulated and real human genotype datasets.** First two columns are population kinship matrices as heatmaps: individuals along x- and y-axes, kinship as color. Diagonal shows inbreeding values. **A.** Admixture scenario for both Large and Small simulations. **B.** Last generation of 20-generation admixed family, shows larger kinship values near diagonal corresponding to siblings, first cousins, etc. **C.** Minor allele frequency (MAF) distributions. Real datasets and tree simulations had  $\text{MAF} \geq 0.01$  filter. **D.** Human Origins is an array dataset of a large diversity of global populations. **G.** Human Genome Diversity Panel (HGDP) is a WGS dataset from global native populations. **J.** 1000 Genomes Project is a WGS dataset of global cosmopolitan populations. **F,I,L.** Trees between subpopulations fit to real data. **E,H,K.** Simulations from trees fit to the real data recapitulate subpopulation structure.



RC model draws coefficients independent of allele frequency, corresponding to neutral traits [45, 61], which results in a wider effect size distribution that reduces association power and effective polygenicity compared to FES.

We evaluate using two complementary measures: (1) SRMSD<sub>p</sub> (p-value signed root mean square deviation) measures p-value calibration (closer to zero is better), and (2) AUC<sub>PR</sub> (precision-recall area under the curve) measures causal locus classification performance (higher is better; Fig. 2). SRMSD<sub>p</sub> is a more robust alternative to the common inflation factor  $\lambda$  and type I error control measures; there is a correspondence between  $\lambda$  and SRMSD<sub>p</sub>, with SRMSD<sub>p</sub> > 0.01 giving  $\lambda > 1.06$  (Fig. S1) and thus evidence of miscalibration close to the rule of thumb of  $\lambda > 1.05$  [20]. AUC<sub>PR</sub> has been used to evaluate association models [80], and reflects statistical power while being robust to miscalibrated models (Appendix B).

Both PCA and LMM were evaluated in each replicate dataset including a number of PCs  $r$  between 0 and 90 as fixed covariates. In terms of p-value calibration, for PCA the best number of PCs  $r$  (minimizing mean |SRMSD<sub>p</sub>| over replicates) is typically large across all datasets, but much smaller “min”  $r$  values often performed as well (numbers in parentheses in Table 3). Most cases had a mean |SRMSD<sub>p</sub>| < 0.01 (marked with asterisks in Table 3), whose p-values are effectively calibrated. However, PCA best and min  $r$  values tended to be large on the family simulation and real datasets, and those cases were often miscalibrated. In contrast, for LMM,  $r = 0$  (no PCs) was always best, and was always calibrated. Comparing LMM with  $r = 0$  to PCA with its best  $r$ , LMM always had significantly smaller |SRMSD<sub>p</sub>| than PCA or was statistically tied. For AUC<sub>PR</sub> and PCA, the best  $r$  was always smaller than the best  $r$  for |SRMSD<sub>p</sub>|, so there is often a tradeoff between calibrated p-values versus classification performance. For LMM there is no tradeoff, as  $r = 0$  had AUC<sub>PR</sub> not significantly different from the best  $r$  in all cases except two (the min  $r$  was 2 for both 1000 Genomes simulation with FES trait and 1000 Genomes real dataset with RC trait). Lastly, LMM with its best  $r$  always had significantly greater AUC<sub>PR</sub> than PCA with its best  $r$  except for one tie (HGDP with RC trait).

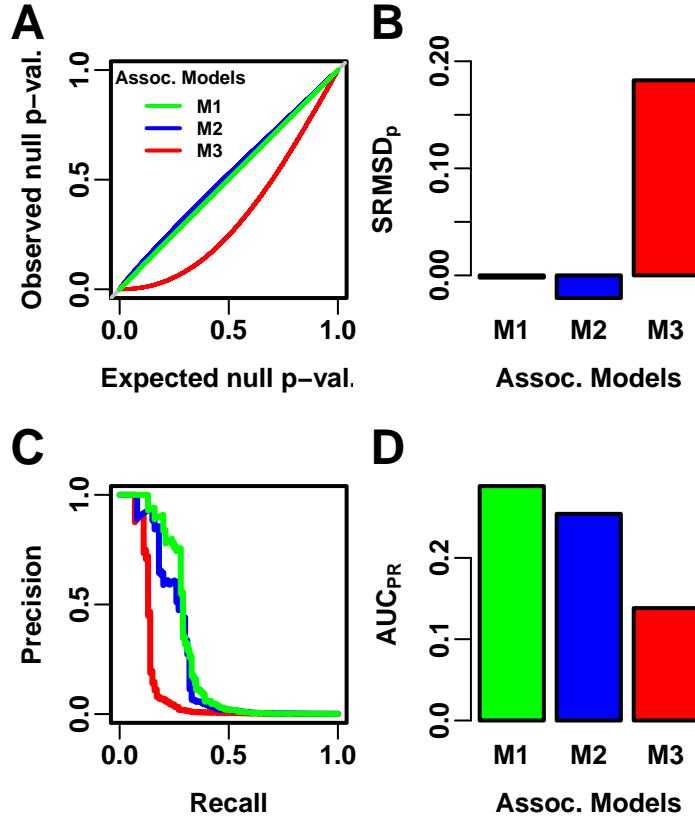


Figure 2: **Illustration of evaluation measures.** Three archetypal models illustrate our complementary measures: M1 is ideal, M2 overfits slightly, M3 is naive. **A.** QQ plot of p-values of “null” (non-causal) loci. M1 has desired uniform p-values, M2/M3 are miscalibrated. **B.** SRMSD<sub>p</sub> (p-value Signed Root Mean Square Deviation) measures signed distance between observed and expected null p-values (closer to zero is better). **C.** Precision and Recall (PR) measure causal locus classification performance (higher is better). **D.** AUC<sub>PR</sub> (Area Under the PR Curve) reflects power (higher is better).

### 3.2 Evaluations in admixture simulations

Now we look more closely at results per dataset. The complete SRMSD<sub>p</sub> and AUC<sub>PR</sub> distributions for the admixture simulations and FES traits are in Fig. 3. RC traits gave qualitatively similar results (Fig. S2).

In the large admixture simulation, the SRMSD<sub>p</sub> of PCA is largest when  $r = 0$  (no PCs) and decreases rapidly to near zero at  $r = 3$ , where it stays for up to  $r = 90$  (Fig. 3A). Thus, PCA has calibrated p-values for  $r \geq 3$ , smaller than the theoretical optimum for this simulation of  $r = K - 1 = 9$ . In contrast, the SRMSD<sub>p</sub> for LMM starts near zero for  $r = 0$ , but becomes negative

Table 3: Overview of PCA and LMM evaluation results

Dataset	Trait model <sup>a</sup>	SRMSD <sub>p</sub>			AUC <sub>PR</sub>		
		Best (min <sup>b</sup> ) PCs PCA	LMM	Best <sup>c</sup>	Best (min <sup>b</sup> ) PCs PCA	LMM	Best <sup>c</sup>
Admix. Large sim.	FES	84* (3*)	0*	tie	3	3 (0)	LMM
Admix. Small sim.	FES	4* (2*)	0*	LMM	4 (1)	0	LMM
Admix. Family sim.	FES	90 (87)	0*	LMM	83 (34)	0	LMM
Human Origins	FES	90 (87)	0*	LMM	34 (9)	1 (0)	LMM
HGDP	FES	87* (34*)	0*	LMM	19 (16)	1 (0)	LMM
1000 Genomes	FES	39 (32)	0*	LMM	8	1 (0)	LMM
Human Origins sim.	FES	90* (80*)	0*	tie	47 (36)	0	LMM
HGDP sim.	FES	43* (20*)	0*	LMM	17 (15)	0	LMM
1000 Genomes sim.	FES	77* (15*)	0*	LMM	16 (6)	2	LMM
Admix. Large sim.	RC	89* (3*)	0*	tie	3	2 (0)	LMM
Admix. Small sim.	RC	8* (2*)	0*	tie (LMM)	1 (0)	0	LMM
Admix. Family sim.	RC	90 (88)	0*	LMM	74 (28)	0	LMM
Human Origins	RC	89* (79*)	0*	LMM	34 (18)	5 (0)	LMM
HGDP	RC	77* (30*)	0*	LMM	19 (13)	3 (0)	tie (LMM)
1000 Genomes	RC	37* (27*)	0*	LMM	19 (4)	9 (2)	LMM
Human Origins sim.	RC	89* (85*)	0*	tie	45 (25)	0	LMM
HGDP sim.	RC	30* (23*)	0*	LMM	18 (15)	5 (0)	LMM
1000 Genomes sim.	RC	90* (16)	0*	LMM	10 (6)	2 (0)	LMM

<sup>a</sup>FES: Fixed Effect Sizes, RC: Random Coefficients.

<sup>b</sup>Smallest  $r$  (number of PCs) whose distribution ( $|\text{SRMSD}_p|$  or  $\text{AUC}_{\text{PR}}$ ) was not significantly different (Wilcoxon paired 1-tailed  $p > 0.01$ ) from best  $r$  (if any).

<sup>c</sup>Tie if distributions of best PCA and LMM version did not differ significantly (Wilcoxon paired 1-tailed  $p > 0.01$ ). Same result for “min” except cases in parentheses.

\* $r$  for which mean  $|\text{SRMSD}_p| < 0.01$ .

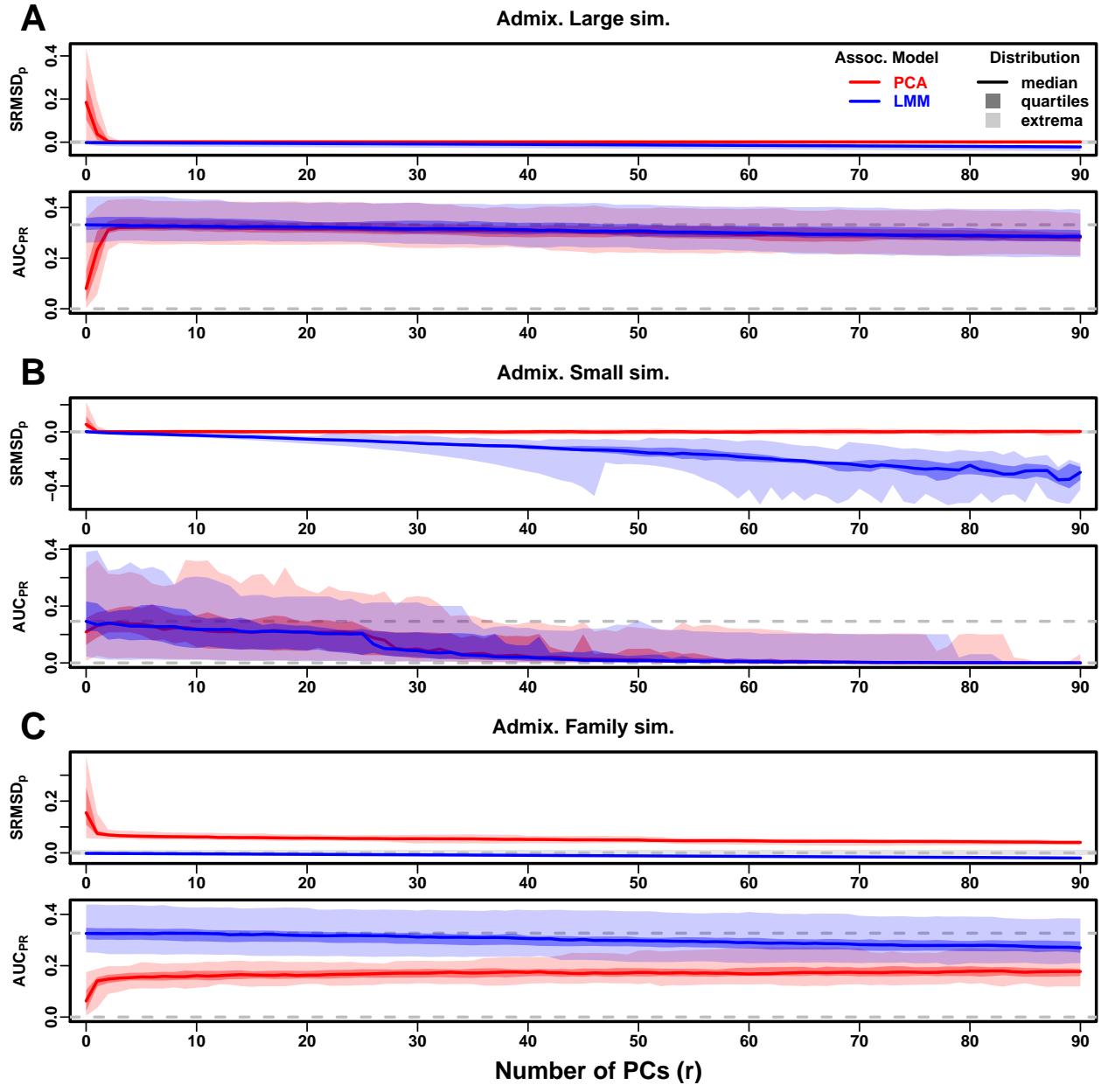
as  $r$  increases (p-values are conservative). The  $\text{AUC}_{\text{PR}}$  distribution of PCA is similarly worst at  $r = 0$ , increases rapidly and peaks at  $r = 3$ , then decreases slowly for  $r > 3$ , while the  $\text{AUC}_{\text{PR}}$  distribution for LMM starts near its maximum at  $r = 0$  and decreases with  $r$ . Although the  $\text{AUC}_{\text{PR}}$  distributions for LMM and PCA overlap considerably at each  $r$ , LMM with  $r = 0$  has significantly greater  $\text{AUC}_{\text{PR}}$  values than PCA with  $r = 3$  (Table 3). However, qualitatively PCA performs nearly as well as LMM in this simulation.

The observed robustness to large  $r$  led us to consider smaller sample sizes. A model with large numbers of parameters  $r$  should overfit more as  $r$  approaches the sample size  $n$ . Rather than increase  $r$  beyond 90, we reduce individuals to  $n = 100$ , which is small for typical association studies but may occur in studies of rare diseases, pilot studies, or other constraints. To compensate for the loss of power due to reducing  $n$ , we also reduce the number of causal loci (fixed ratio  $n/m_1 = 10$ ), which increases per-locus effect sizes. We found a large decrease in performance for both models as  $r$  increases, and best performance for  $r = 1$  for PCA and  $r = 0$  for LMM (Fig. 3B). Remarkably, LMM attains much larger negative  $\text{SRMSD}_p$  values than in our other evaluations. LMM with  $r = 0$  is significantly better than PCA ( $r = 1$  to 4) in both metrics (Table 3), but qualitatively the difference is negligible.

The family simulation adds a 20-generation random family to our large admixture simulation. Only the last generation is studied for association, which contains numerous siblings, first cousins, etc., with the initial admixture structure preserved by geographically-biased mating. Our evaluation reveals a sizable gap in both metrics between LMM and PCA across all  $r$  (Fig. 3C). LMM again performs best with  $r = 0$  and achieves mean  $|\text{SRMSD}_p| < 0.01$ . However, PCA does not achieve mean  $|\text{SRMSD}_p| < 0.01$  at any  $r$ , and its best mean  $\text{AUC}_{\text{PR}}$  is considerably worse than that of LMM. Thus, LMM is conclusively superior to PCA, and the only calibrated model, when there is family structure.

### 3.3 Evaluations in real human genotype datasets

Next we repeat our evaluations with real human genotype data, which differs from our simulations in allele frequency distributions and more complex population structures with greater differentiation,



**Figure 3: Evaluations in admixture simulations.** Traits simulated from FES model. PCA and LMM models have varying number of PCs ( $r \in \{0, \dots, 90\}$  on x-axis), with the distributions (y-axis) of SRMSD<sub>p</sub> (top subpanel) and AUC<sub>PR</sub> (bottom subpanel) for 50 replicates. Best performance is zero SRMSD<sub>p</sub> and large AUC<sub>PR</sub>. Zero and maximum median AUC<sub>PR</sub> values are marked with horizontal gray dashed lines, and  $|\text{SRMSD}_p| < 0.01$  is marked with a light gray area. LMM performs best with  $r = 0$ , PCA with various  $r$ .

**A.** Large simulation ( $n = 1,000$  individuals). **B.** Small simulation ( $n = 100$ ) shows overfitting for large  $r$ . **C.** Family simulation ( $n = 1,000$ ) has admixed founders and large numbers of close relatives from a realistic random 20-generation pedigree. PCA performs poorly compared to LMM: SRMSD<sub>p</sub> > 0 for all  $r$  and large AUC<sub>PR</sub> gap.

numerous correlated subpopulations, and potential cryptic family relatedness.

Human Origins has the greatest number and diversity of subpopulations. The SRMSD<sub>p</sub> and AUC<sub>PR</sub> distributions in this dataset and FES traits (Fig. 4A) most resemble those from the family simulation (Fig. 3C). In particular, while LMM with  $r = 0$  performed optimally (both metrics) and satisfies mean  $|\text{SRMSD}_p| < 0.01$ , PCA maintained  $\text{SRMSD}_p > 0.01$  for all  $r$  and its AUC<sub>PR</sub> were all considerably smaller than the best AUC<sub>PR</sub> of LMM.

HGDP has the fewest individuals among real datasets, but compared to Human Origins contains more loci and low-frequency variants. Performance (Fig. 4B) was intermediate between the admixture and family simulations. In particular, here both LMM ( $r = 0$ ) and PCA ( $r \geq 31$ ) achieve mean  $|\text{SRMSD}_p| < 0.01$  (p-values are calibrated). However, there is a sizable AUC<sub>PR</sub> gap between LMM and PCA. Maximum AUC<sub>PR</sub> values were lowest in HGDP compared to the two other real datasets.

1000 Genomes has the fewest subpopulations but largest number of individuals per subpopulation. Thus, although this dataset has the simplest subpopulation structure among the real datasets, we find SRMSD<sub>p</sub> and AUC<sub>PR</sub> distributions (Fig. 4C) that again most resemble our earlier family simulation, with mean  $|\text{SRMSD}_p| < 0.01$  for LMM only and large AUC<sub>PR</sub> gaps between LMM and PCA.

Our results are qualitatively different for RC traits, which had smaller AUC<sub>PR</sub> gaps between LMM and PCA (Fig. S3). Maximum AUC<sub>PR</sub> were smaller in RC compared to FES in Human Origins and 1000 Genomes, suggesting lower power for RC traits across association models. Nevertheless, LMM with  $r = 0$  was significantly better than PCA for all metrics in the real datasets and RC traits (Table 3).

### 3.4 Evaluations in tree simulations fit to human data

To better understand which features of the real datasets lead to the large differences in performance between LMM and PCA, we carried out tree simulations. Human subpopulations are related roughly by trees, which induce the strongest correlations and have numerous tips, so we fit trees to each real dataset and tested if data simulated from these complex tree structures could recapitulate our

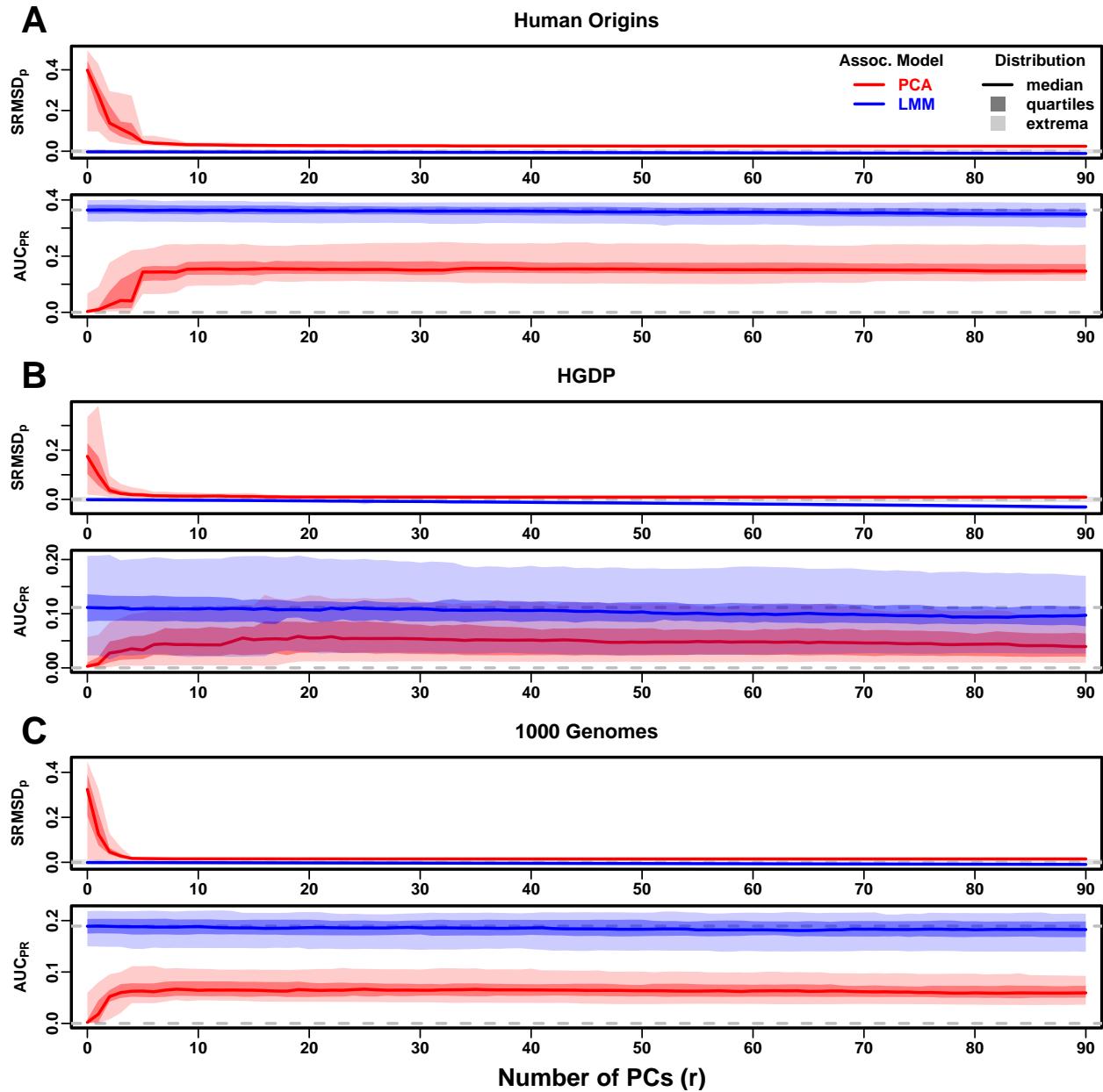


Figure 4: **Evaluations in real human genotype datasets.** Traits simulated from FES model. Same setup as Fig. 3, see that for details. These datasets strongly favor LMM with no PCs over PCA, with distributions that most resemble the family simulation. **A.** Human Origins. **B.** Human Genome Diversity Panel (HGDP). **C.** 1000 Genomes Project.

previous results (Fig. 1). These tree simulations also feature non-uniform ancestral allele frequency distributions, which recapitulated some of the skew for smaller minor allele frequencies of the real datasets (Fig. 1C). The SRMSD<sub>p</sub> and AUC<sub>PR</sub> distributions for these tree simulations (Fig. 5) resembled our admixture simulation more than either the family simulation (Fig. 3) or real data results (Fig. 4). Both LMM with  $r = 0$  and PCA (various  $r$ ) achieve mean  $|\text{SRMSD}_p| < 0.01$  (Table 3). The AUC<sub>PR</sub> distributions of both LMM and PCA track closely as  $r$  is varied, although there is a small gap resulting in LMM ( $r = 0$ ) besting PCA in all three simulations. The results are qualitatively similar for RC traits (Fig. S4 and Table 3). Overall, these tree simulations do not recapitulate the large LMM advantage over PCA observed on the real data.

### 3.5 Numerous distant relatives explain poor PCA performance in real data

In principle, PCA performance should be determined by the dimensionality of relatedness, since PCA is a low-dimensional model whereas LMM can model high-dimensional relatedness without overfitting. We used the Tracy-Widom test [7] with  $p < 0.01$  to estimate dimensionality as the number of significant PCs (Fig. S5A). The true dimensionality of our simulations is slightly underestimated (Table 2), but we confirm that the family simulation has the greatest dimensionality, and real datasets have greater estimates than their respective tree simulations, which confirms our hypothesis to some extent. However, estimated dimensionalities do not separate real datasets from tree simulations, as required to predict the observed PCA performance. Moreover, the HGDP and 1000 Genomes dimensionality estimates are 46 and 61, respectively, yet PCA performed poorly for all  $r \leq 90$  numbers of PCs (Fig. 4). The top eigenvalue explained a proportion of variance proportional to  $F_{\text{ST}}$  (Table 2), but the rest of the top 10 eigenvalues show no clear differences between datasets, except the small simulation had larger variances explained per eigenvalue (expected since it has fewer eigenvalues; Fig. S5C). Comparing cumulative variance explained versus rank fraction across all eigenvalues, all datasets increase from their starting point almost linearly until they reach 1, except the family simulation has much greater variance explained by mid-rank eigenvalues (Fig. S5B). Overall, there is no separation between real datasets (where PCA performed poorly) and tree simulations (where PCA performed relatively well) in terms of their eigenvalues or

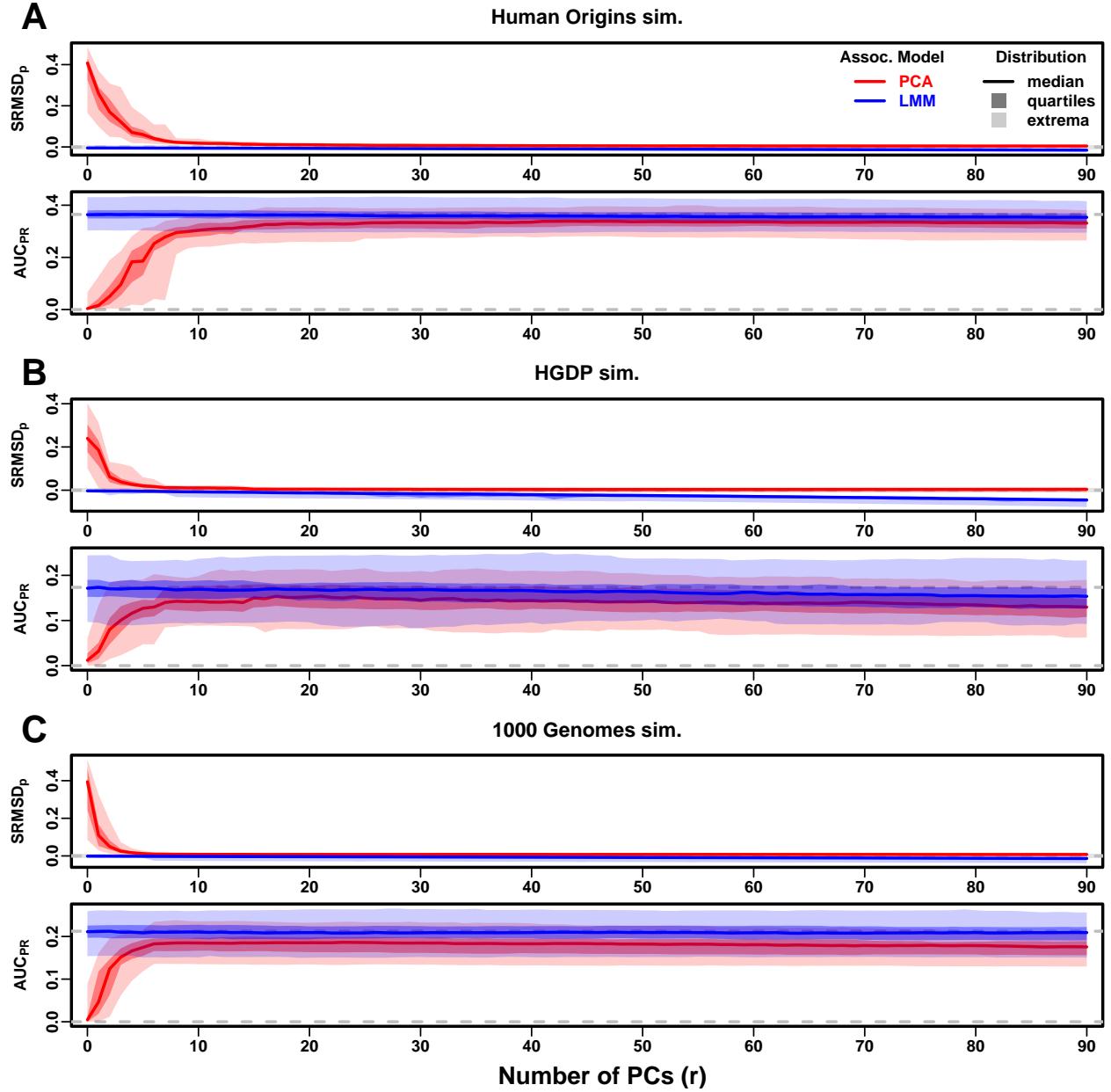


Figure 5: **Evaluations in tree simulations fit to human data.** Traits simulated from FES model. Same setup as Fig. 3, see that for details. These tree simulations, which exclude family structure by design, do not explain the large gaps in LMM-PCA performance observed in the real data. **A.** Human Origins tree simulation. **B.** Human Genome Diversity Panel (HGDP) tree simulation. **C.** 1000 Genomes Project tree simulation.

dimensionality estimates.

Local kinship, which is recent relatedness due to family structure excluding population structure, is the presumed cause of the LMM to PCA performance gap observed in real datasets but not their tree simulation counterparts. Instead of inferring local kinship through increased dimensionality, as attempted in the last section, here we measure it directly using the KING-robust estimator [75]. We observe more large local kinship in the real datasets and the family simulation compared to the other simulations (Fig. 6). However, for real data this distribution depends on the subpopulation structure, since locally related pairs are most likely in the same subpopulation. Therefore, the only comparable curve to each real dataset is their corresponding tree simulation, which matches subpopulation structure. In all real datasets we identified highly related individual pairs with kinship above the 4th degree relative threshold of 0.022 [75, 81]. However, these highly related pairs are vastly outnumbered by more distant pairs with evident non-zero local kinship as compared to the extreme tree simulation values.

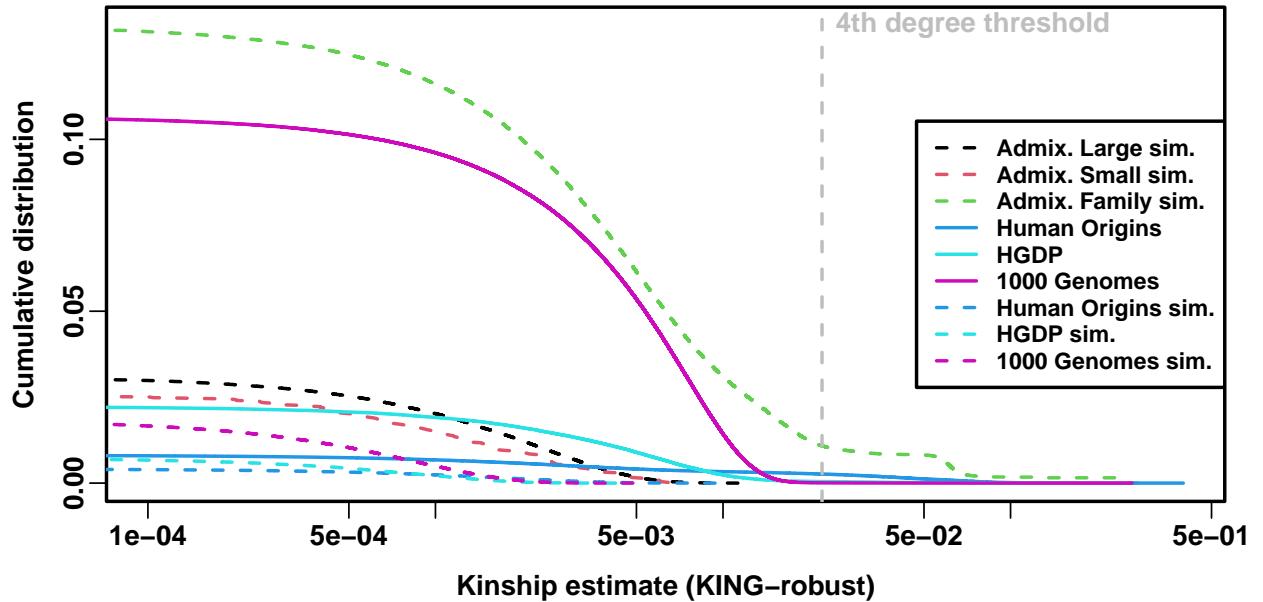


Figure 6: **Local kinship distribution.** Curves are complementary cumulative distribution of lower triangular kinship matrix (self kinship excluded) from KING-robust estimator. Note log x-axis; negative estimates are counted but not shown. Most values are below 4th degree relative threshold. Each real dataset has a greater cumulative than its tree simulations.

To try to improve PCA performance, we followed the standard practice of removing 4th degree relatives, which reduced sample sizes between 5% and 10% (Table S1). Only  $r = 0$  for LMM and  $r = 20$  for PCA were tested, as these performed well in our earlier evaluation, and only FES traits were tested because they previously displayed the large PCA-LMM performance gap. LMM significantly outperforms PCA in all these cases (Wilcoxon paired 1-tailed  $p < 0.01$ ; Fig. 7). Notably, PCA still had miscalibrated p-values in Human Origins and 1000 Genomes ( $|\text{SRMSD}_p| > 0.01$ ). Otherwise,  $\text{AUC}_{\text{PR}}$  and  $\text{SRMSD}_p$  ranges were similar here as in our earlier evaluation. Therefore, the removal of the small number of highly related individual pairs had a negligible effect in PCA performance, so the larger number of more distantly related pairs explain the poor PCA performance in the real datasets.

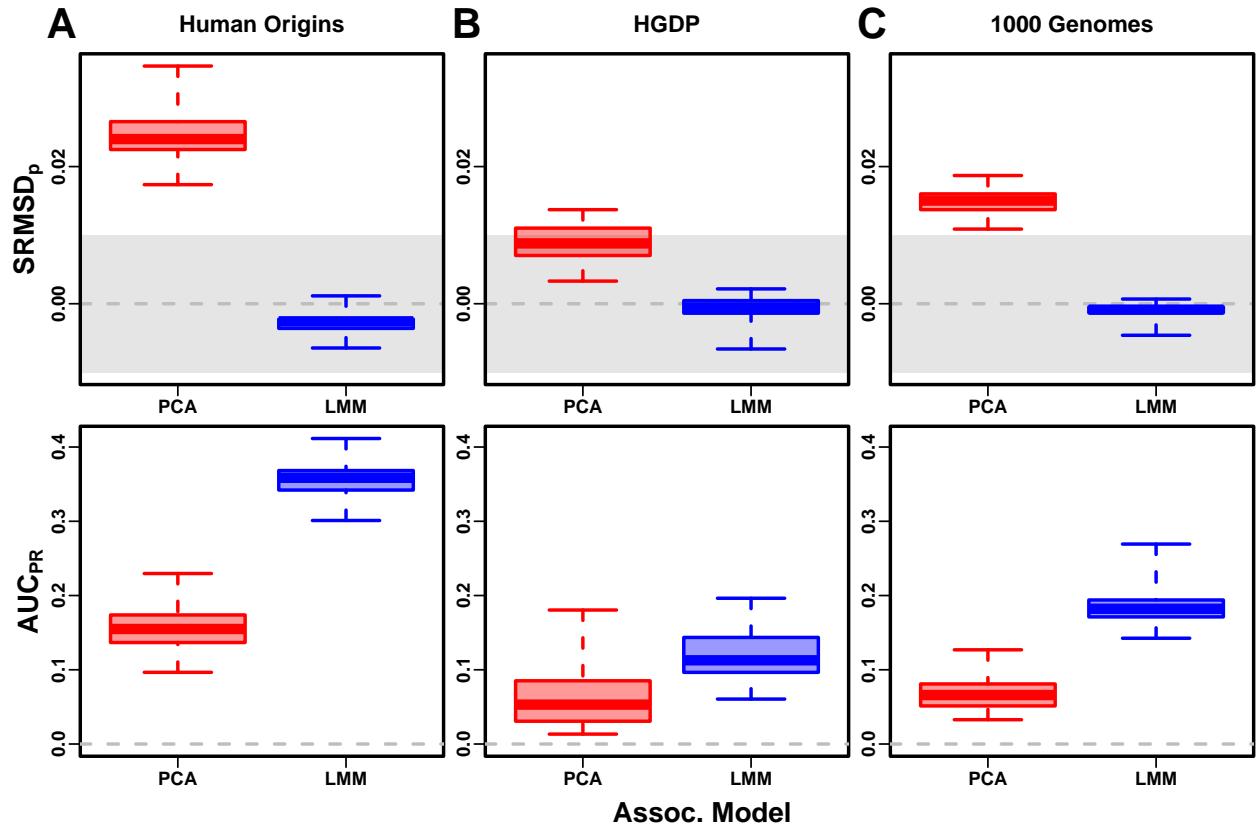


Figure 7: **LMM and PCA performance in real datasets excluding 4th degree relatives.** LMM had  $r = 0$  PCs and PCA had  $r = 20$ . FES traits only. Each dataset is a column, rows are metrics. First row has  $|\text{SRMSD}_p| < 0.01$  band marked as gray area.

## 4 Discussion

Our evaluations conclusively determined that LMM without PCs performs better than PCA (for any number of PCs) across all scenarios, including all real and simulated genotypes and two trait simulation models. Although the addition of a few PCs to LMM does not greatly hurt its performance (except for small sample sizes), they generally did not improve it either (Table 3), which agrees with previous observations [44] but contradicts others [16, 20]. Our findings make sense since PCs are the eigenvectors of the same kinship matrix that parametrizes random effects, so including both is redundant.

Previous studies found that PCA was better calibrated than LMM for unusually differentiated markers [20, 29, 31], which as simulated were an artificial scenario not based on a population genetics model, and are otherwise believed to be unusual [32, 49]. Our evaluations on real human data, which contain such loci in relevant proportions if they exist, do not replicate that result. Cryptic relatedness strongly favors LMM, an advantage that probably outweighs this potential PCA benefit in real data.

Relative to LMM, the behavior of PCA fell between two extremes. When PCA performed well, there was a small number of PCs with both calibrated p-values and  $AUC_{PR}$  near that of LMM without PCs. Conversely, PCA performed poorly when no number of PCs had either calibrated p-values or acceptably large  $AUC_{PR}$ . There were no PCA cases where high numbers of PCs optimized an acceptable  $AUC_{PR}$ , or miscalibrated p-values but high  $AUC_{PR}$ . PCA performed well in the admixture simulations (without families, both trait models), real human genotypes with RC traits, and, to a lesser extent, the tree simulations (both trait models). Conversely, PCA performed poorly in the admixed family simulation (both trait models) and the real human genotypes with FES traits.

PCA assumes that genetic relatedness is low-dimensional, whereas LMM can handle high-dimensional relatedness. Thus, PCA performs well in the admixture simulation, which is explicitly low-dimensional (see Models and Methods), and our tree simulations, which had few nodes with long branches so a low-dimensional approximation suffices. Conversely, PCA performs poorly under family structure because its kinship matrix is high-dimensional (Fig. S5). However, estimating the dimensionality of real datasets is challenging because estimated eigenvalues have biased distribu-

tions. Dimensionality estimated using the Tracy-Widom test [7] did not fully predict the datasets that PCA performs well on. In contrast, estimated local kinship finds considerable cryptic relatedness in all real human datasets and better explains why PCA performs poorly there. The trait model also influences the relative performance of PCA, so genotype-only parameters (eigenvalues or local kinship) alone do not tell the full story.

PCA is at best underpowered relative to LMMs, and at worst miscalibrated regardless of the numbers of PCs included, in real human genotype tests. Among our simulations, such poor performance occurred only in the admixed family. Local kinship estimates reveal considerable family relatedness in the real datasets absent in the corresponding tree simulations. Admixture is also absent in our tree simulations, but our simulations and theory show that admixture is handled well by PCA. Hundreds of close relative pairs have been identified in 1000 Genomes [82–85], but their removal does not improve PCA performance sufficiently in our tests, so the larger number of more distantly related pairs are PCA’s most serious obstacle in practice. Distant relatives are expected to be numerous in any large human dataset [86, 87]. Our FES trait tests show that cryptic relatedness is more challenging when rarer variants have larger coefficients. Overall, the high dimensionality induced by cryptic relatedness is the key challenge for PCA association in modern datasets that is readily overcome by LMM.

Our tests also found PCA robust to large numbers of PCs, far beyond the optimal choice, agreeing with previous anecdotal observations [5, 30], in contrast to using too few PCs for which there is a large performance penalty. The exception was the small sample size simulation, where only small numbers of PCs performed well. In contrast, LMM is simpler since there is no need to choose the number of PCs. However, an LMM with a large number of covariates may have conservative p-values (as observed for LMM with large numbers of PCs), a weakness of the score test used by the LMM we evaluated that may be overcome with other statistical tests. Simulations or post hoc evaluations remain crucial for ensuring that statistics are calibrated.

The largest limitation of our work is that we only considered quantitative traits. We noted that previous evaluations involving case-control traits tended to report PCA-LMM ties or mixed results, an observation potentially confounded by the use of low-dimensional simulations without family

relatedness (Table 1). An additional concern is case-control ascertainment bias, which appears to affect LMMs more severely, although recent work appears to solve this problem [29, 33]. Future evaluations should aim to include our simulations and real datasets, to ensure that previous results were not biased in favor of PCA by employing unrealistic low-dimensional genotype simulations, or by not simulating large coefficients for rare variants expected for diseases by various selection models.

Overall, our results lead us to recommend LMM over PCA for association studies in general. Although PCA offer flexibility and speed compared to LMM, additional work is required to ensure that PCA is adequate, including removal of close relatives (lowering sample size and wasting resources) followed by simulations or other evaluations of statistics, and even then PCA may perform poorly in terms of both type I error control and power. The large numbers of distant relatives expected of any real dataset all but ensures that PCA will perform poorly compared to LMM. Our findings also suggest that related applications such as polygenic models may enjoy gains in power and accuracy by employing an LMM instead of PCA to model relatedness [18, 80]. PCA remains indispensable across population genetics, from visualizing population structure and performing quality control to its deep connection to admixture models, but the time has come to limit its use in association testing in favor of LMM or other, richer models capable of modeling all forms of relatedness.

## 5 Appendices

### 5.1 Appendix A: Fitting ancestral allele frequency distribution to real data

We calculated  $\hat{p}_i^T$  distributions of each real dataset. However, differentiation increases the variance of these sample  $\hat{p}_i^T$  relative to the true  $p_i^T$  [28]. We present a new algorithm for constructing an “undifferentiated” distribution based on the input data but with the lower variance of the true ancestral distribution. Suppose the  $p_i^T$  distribution over loci  $i$  satisfies  $E[p_i^T] = \frac{1}{2}$  and  $\text{Var}(p_i^T) = V^T$ . The sample allele frequency  $\hat{p}_i^T$ , conditioned on  $p_i^T$ , satisfies

$$E[\hat{p}_i^T | p_i^T] = p_i^T, \quad \text{Var}(\hat{p}_i^T | p_i^T) = p_i^T (1 - p_i^T) \bar{\varphi}^T,$$

where  $\bar{\varphi}^T = \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n \varphi_{jk}^T$  is the mean kinship over all individual [28]. The unconditional moments of  $\hat{p}_i^T$  follow from the laws of total expectation and variance:  $E[\hat{p}_i^T] = \frac{1}{2}$  and

$$W^T = \text{Var}(\hat{p}_i^T) = \bar{\varphi}^T \frac{1}{4} + (1 - \bar{\varphi}^T) V^T.$$

Since  $V^T \leq \frac{1}{4}$  and  $\bar{\varphi}^T \geq 0$ , then  $W^T \geq V^T$ . Thus, the goal is to construct a new distribution with the original, lower variance of

$$V^T = \frac{W^T - \frac{1}{4}\bar{\varphi}^T}{1 - \bar{\varphi}^T}. \quad (9)$$

We use the unbiased estimator  $\hat{W}^T = \frac{1}{m} \sum_{i=1}^m (\hat{p}_i^T - \frac{1}{2})^2$ , while  $\bar{\varphi}^T$  is calculated from the tree parameters: the subpopulation coancestry matrix (Eq. (7)), expanded from subpopulations to individuals, the diagonal converted to kinship (reversing Eq. (8)), and the matrix averaged. However, since our model ignores the MAF filters imposed in our simulations,  $\bar{\varphi}^T$  was adjusted. For Human Origins the true model  $\bar{\varphi}^T$  of 0.143 was used. For 1000 Genomes and HGDP the true  $\bar{\varphi}^T$  are 0.126 and 0.124, respectively, but 0.4 for both produced a better fit.

Lastly, we construct new allele frequencies,

$$p' = w\hat{p}_i^T + (1 - w)q,$$

by a weighted average of  $\hat{p}_i^T$  and  $q \in (0, 1)$  drawn independently from a different distribution.  $E[q] = \frac{1}{2}$  is required to have  $E[p'] = \frac{1}{2}$ . The resulting variance is

$$\text{Var}(p') = w^2 W^T + (1 - w)^2 \text{Var}(q),$$

which we equate to the desired  $V^T$  (Eq. (9)) and solve for  $w$ . For simplicity, we also set  $\text{Var}(q) = V^T$ , which is achieved with:

$$q \sim \text{Beta}\left(\frac{1}{2} \left(\frac{1}{4V^T} - 1\right), \frac{1}{2} \left(\frac{1}{4V^T} - 1\right)\right).$$

Although  $w = 0$  yields  $\text{Var}(p') = V^T$ , we use the second root of the quadratic equation to use  $\hat{p}_i^T$ :

$$w = \frac{2V^T}{W^T + V^T}.$$

## 5.2 Appendix B: comparisons between SRMSD<sub>p</sub>, AUC<sub>PR</sub>, and evaluation measures from the literature

### 5.2.1 The inflation factor $\lambda$

Test statistic inflation has been used to measure model calibration [1, 20]. The inflation factor  $\lambda$  is defined as the median  $\chi^2$  association statistic divided by theoretical median under the null hypothesis [2]. To compare p-values from non- $\chi^2$  tests (such as t-statistics),  $\lambda$  can be calculated from p-values using

$$\lambda = \frac{F^{-1}(1 - p_{\text{median}})}{F^{-1}(1 - u_{\text{median}})},$$

where  $p_{\text{median}}$  is the median observed p-value (including causal loci),  $u_{\text{median}} = \frac{1}{2}$  is its null expectation, and  $F$  is the  $\chi^2$  cumulative density function ( $F^{-1}$  is the quantile function).

To compare  $\lambda$  and SRMSD<sub>p</sub> directly, for simplicity assume that all p-values are null. In this case, calibrated p-values give  $\lambda = 1$  and SRMSD<sub>p</sub> = 0. However, non-uniform p-values with the expected median, such as from genomic control [2], result in  $\lambda = 1$ , but SRMSD<sub>p</sub> ≠ 0 except for uniform p-values, a key flaw of  $\lambda$  that SRMSD<sub>p</sub> overcomes. Inflated statistics (anti-conservative p-values) give  $\lambda > 1$  and SRMSD<sub>p</sub> > 0. Deflated statistics (conservative p-values) give  $\lambda < 1$  and SRMSD<sub>p</sub> < 0. Thus,  $\lambda \neq 1$  always implies SRMSD<sub>p</sub> ≠ 0 (where  $\lambda - 1$  and SRMSD<sub>p</sub> have the same sign), but not the other way around. Overall,  $\lambda$  depends only on the median p-value, while SRMSD<sub>p</sub> uses the complete distribution. However, SRMSD<sub>p</sub> requires knowing which loci are null, so unlike  $\lambda$  it is only applicable to simulated traits.

### 5.2.2 Empirical comparison of SRMSD<sub>p</sub> and $\lambda$

There is a near one-to-one correspondence between  $\lambda$  and SRMSD<sub>p</sub> in our data (Fig. S1). PCA tended to be inflated ( $\lambda > 1$  and SRMSD<sub>p</sub> > 0) whereas LMM tended to be deflated ( $\lambda < 1$  and

$\text{SRMSD}_p < 0$ ), otherwise the data for both models fall on the same contiguous curve. We fit a sigmoidal function to this data,

$$\text{SRMSD}_p(\lambda) = a \frac{\lambda^b - 1}{\lambda^b + 1}, \quad (10)$$

which for  $a, b > 0$  satisfies  $\text{SRMSD}_p(\lambda = 1) = 0$  and reflects  $\log(\lambda)$  about zero ( $\lambda = 1$ ):

$$\text{SRMSD}_p(\log(\lambda) = -x) = -\text{SRMSD}_p(\log(\lambda) = x).$$

We fit this model to  $\lambda > 1$  only since it was less noisy and of greater interest, and obtained the curve shown in Fig. S1 with  $a = 0.566$  and  $b = 0.616$ . The value  $\lambda = 1.05$ , a common threshold for benign inflation [20], corresponds to  $\text{SRMSD}_p = 0.0085$  according to Eq. (10). Conversely,  $\text{SRMSD}_p = 0.01$ , serving as a simpler rule of thumb, corresponds to  $\lambda = 1.06$ .

### 5.2.3 Type I error rate

The type I error rate is the proportion of null p-values with  $p \leq t$ . Calibrated p-values have type I error rate near  $t$ , which may be evaluated with a binomial test. This measure may give different results for different  $t$ , for example be significantly miscalibrated only for large  $t$  (due to lack of power for smaller  $t$ ). In contrast,  $\text{SRMSD}_p = 0$  guarantees calibrated type I error rates at all  $t$ , while large  $|\text{SRMSD}_p|$  indicates incorrect type I errors for a range of  $t$ .

### 5.2.4 Statistical power and comparison to AUC<sub>PR</sub>

Power is the probability that a test is declared significant when the alternative hypothesis  $H_1$  holds. At a p-value threshold  $t$ , power equals

$$F(t) = \Pr(p < t | H_1).$$

$F(t)$  is a cumulative function, so it is monotonically increasing and has an inverse. Like type I error control, power may rank models differently depending on  $t$ .

Power is not meaningful when p-values are not calibrated. To establish a clear connection to

$\text{AUC}_{\text{PR}}$ , assume calibrated (uniform) null p-values:  $\Pr(p < t | H_0) = t$ . TPs, FPs, and FNs at  $t$  are

$$\text{TP}(t) = m\pi_1 F(t),$$

$$\text{FP}(t) = m\pi_0 t,$$

$$\text{FN}(t) = m\pi_1(1 - F(t)),$$

where  $\pi_0 = \Pr(H_0)$  is the proportion of null cases and  $\pi_1 = 1 - \pi_0$  of alternative cases. Therefore,

$$\text{Precision}(t) = \frac{\pi_1 F(t)}{\pi_1 F(t) + \pi_0 t},$$

$$\text{Recall}(t) = F(t).$$

Noting that  $t = F^{-1}(\text{Recall})$ , precision can be written as a function of recall, the power function, and constants:

$$\text{Precision}(\text{Recall}) = \frac{\pi_1 \text{Recall}}{\pi_1 \text{Recall} + \pi_0 F^{-1}(\text{Recall})}.$$

This last form leads most clearly to  $\text{AUC}_{\text{PR}} = \int_0^1 \text{Precision}(\text{Recall}) d\text{Recall}$ .

Lastly, consider a simple yet common case in which model  $A$  is uniformly more powerful than model  $B$ :  $F_A(t) > F_B(t)$  for every  $t$ . Therefore  $F_A^{-1}(\text{Recall}) < F_B^{-1}(\text{Recall})$  for every recall value. This ensures that the precision of  $A$  is greater than that of  $B$  at every recall value, so  $\text{AUC}_{\text{PR}}$  is greater for  $A$  than  $B$ . Thus,  $\text{AUC}_{\text{PR}}$  ranks calibrated models according to power.

## Declaration of interests

The authors declare no competing interests.

## Acknowledgments

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/>  
Eigensoft, <https://github.com/DReichLab/EIG>  
g bnpsd, <https://cran.r-project.org/package=bnpsd>  
simfam, <https://cran.r-project.org/package=simfam>  
simtrait, <https://cran.r-project.org/package=simtrait>  
genio, <https://cran.r-project.org/package=genio>  
popkin, <https://cran.r-project.org/package=popkin>  
ape, <https://cran.r-project.org/package=ape>  
nnls, <https://cran.r-project.org/package=nnls>  
PRROC, <https://cran.r-project.org/package=PRROC>  
BEDMatrix, <https://cran.r-project.org/package=BEDMatrix>

## Data and code availability

The data and code generated during this study are available on GitHub at <https://github.com/OchoaLab/pca-assoc-paper>. The public subset of Human Origins is available on the Reich Lab website at <https://reich.hms.harvard.edu/datasets>; non-public samples have to be requested from David Reich. The WGS version of HGDP was downloaded from the Wellcome Sanger Institute FTP site at [ftp://ngs.sanger.ac.uk/production/hgdp/hgdp\\_wgs.20190516/](ftp://ngs.sanger.ac.uk/production/hgdp/hgdp_wgs.20190516/). 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/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/working/20190425_NYGC_GATK/).

## References

- [1] W. Astle and D. J. Balding. “Population Structure and Cryptic Relatedness in Genetic Association Studies”. *Statist. Sci.* 24(4) (2009), pp. 451–471.

- [2] B. Devlin and K. Roeder. “Genomic Control for Association Studies”. *Biometrics* 55(4) (1999), pp. 997–1004.
- [3] B. F. Voight and J. K. Pritchard. “Confounding from Cryptic Relatedness in Case-Control Association Studies”. *PLOS Genetics* 1(3) (2005), e32.
- [4] S. Zhang, X. Zhu, and H. Zhao. “On a semiparametric test to detect associations between quantitative traits and candidate genes using unrelated individuals”. *Genetic Epidemiology* 24(1) (2003), pp. 44–56.
- [5] A. L. Price et al. “Principal components analysis corrects for stratification in genome-wide association studies”. *Nat. Genet.* 38(8) (2006), pp. 904–909.
- [6] M. Bouaziz, C. Ambroise, and M. Guedj. “Accounting for Population Stratification in Practice: A Comparison of the Main Strategies Dedicated to Genome-Wide Association Studies”. *PLOS ONE* 6(12) (2011), e28845.
- [7] N. Patterson, A. L. Price, and D. Reich. “Population Structure and Eigenanalysis”. *PLoS Genet* 2(12) (2006), e190.
- [8] I. T. Jolliffe. *Principal Component Analysis*. 2nd ed. New York: Springer-Verlag, 2002.
- [9] J. K. Pritchard et al. “Association Mapping in Structured Populations”. *The American Journal of Human Genetics* 67(1) (2000), pp. 170–181.
- [10] D. H. Alexander, J. Novembre, and K. Lange. “Fast model-based estimation of ancestry in unrelated individuals”. *Genome Res.* 19(9) (2009), pp. 1655–1664.
- [11] Q. Zhou, L. Zhao, and Y. Guan. “Strong Selection at MHC in Mexicans since Admixture”. *PLoS Genet.* 12(2) (2016), e1005847.
- [12] G. McVean. “A genealogical interpretation of principal components analysis”. *PLoS Genet* 5(10) (2009), e1000686.
- [13] X. Zheng and B. S. Weir. “Eigenanalysis of SNP data with an identity by descent interpretation”. *Theor Popul Biol* 107 (2016), pp. 65–76.

- [14] I. Cabreros and J. D. Storey. “A Likelihood-Free Estimator of Population Structure Bridging Admixture Models and Principal Components Analysis”. *Genetics* 212(4) (2019), pp. 1009–1029.
- [15] A. M. Chiu et al. “Inferring population structure in biobank-scale genomic data”. *The American Journal of Human Genetics* 0(0) (2022).
- [16] K. Zhao et al. “An Arabidopsis Example of Association Mapping in Structured Samples”. *PLOS Genetics* 3(1) (2007), e4.
- [17] H. Xu and Y. Guan. “Detecting Local Haplotype Sharing and Haplotype Association”. *Genetics* 197(3) (2014), pp. 823–838.
- [18] J. Qian et al. “A fast and scalable framework for large-scale and ultrahigh-dimensional sparse regression with application to the UK Biobank”. *PLOS Genetics* 16(10) (2020), e1009141.
- [19] T. Thornton and M. S. McPeek. “ROADTRIPS: case-control association testing with partially or completely unknown population and pedigree structure”. *Am. J. Hum. Genet.* 86(2) (2010), pp. 172–184.
- [20] A. L. Price et al. “New approaches to population stratification in genome-wide association studies”. *Nature Reviews Genetics* 11(7) (2010), pp. 459–463.
- [21] S. Lee et al. “Sparse Principal Component Analysis for Identifying Ancestry-Informative Markers in Genome-Wide Association Studies”. *Genetic Epidemiology* 36(4) (2012), pp. 293–302.
- [22] G. Abraham and M. Inouye. “Fast Principal Component Analysis of Large-Scale Genome-Wide Data”. *PLOS ONE* 9(4) (2014), e93766.
- [23] K. Galinsky et al. “Fast Principal-Component Analysis Reveals Convergent Evolution of ADH1B in Europe and East Asia”. *The American Journal of Human Genetics* 98(3) (2016), pp. 456–472.
- [24] G. Abraham, Y. Qiu, and M. Inouye. “FlashPCA2: principal component analysis of Biobank-scale genotype datasets”. *Bioinformatics* 33(17) (2017), pp. 2776–2778.

- [25] A. Agrawal et al. “Scalable probabilistic PCA for large-scale genetic variation data”. *PLOS Genetics* 16(5) (2020), e1008773.
- [26] J. Yu et al. “A unified mixed-model method for association mapping that accounts for multiple levels of relatedness”. *Nat. Genet.* 38(2) (2006), pp. 203–208.
- [27] H. M. Kang et al. “Efficient control of population structure in model organism association mapping”. *Genetics* 178(3) (2008), pp. 1709–1723.
- [28] A. Ochoa and J. D. Storey. “Estimating FST and kinship for arbitrary population structures”. *PLoS Genet* 17(1) (2021), e1009241.
- [29] J. Yang et al. “Advantages and pitfalls in the application of mixed-model association methods”. *Nat Genet* 46(2) (2014), pp. 100–106.
- [30] H. M. Kang et al. “Variance component model to account for sample structure in genome-wide association studies”. *Nat. Genet.* 42(4) (2010), pp. 348–354.
- [31] C. Wu et al. “A Comparison of Association Methods Correcting for Population Stratification in Case–Control Studies”. *Annals of Human Genetics* 75(3) (2011), pp. 418–427.
- [32] J. H. Sul and E. Eskin. “Mixed models can correct for population structure for genomic regions under selection”. *Nature Reviews Genetics* 14(4) (2013), p. 300.
- [33] W. Zhou et al. “Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies”. *Nat Genet* 50(9) (2018), pp. 1335–1341.
- [34] Y. S. Aulchenko, D.-J. de Koning, and C. Haley. “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) (2007), pp. 577–585.
- [35] Z. Zhang et al. “Mixed linear model approach adapted for genome-wide association studies”. *Nat Genet* 42(4) (2010), pp. 355–360.
- [36] C. Lippert et al. “FaST linear mixed models for genome-wide association studies”. *Nat. Methods* 8(10) (2011), pp. 833–835.

- [37] J. Yang et al. “GCTA: a tool for genome-wide complex trait analysis”. *Am. J. Hum. Genet.* 88(1) (2011), pp. 76–82.
- [38] J. Listgarten et al. “Improved linear mixed models for genome-wide association studies”. *Nat Methods* 9(6) (2012), pp. 525–526.
- [39] X. Zhou and M. Stephens. “Genome-wide efficient mixed-model analysis for association studies”. *Nat. Genet.* 44(7) (2012), pp. 821–824.
- [40] G. R. Svishcheva et al. “Rapid variance components-based method for whole-genome association analysis”. *Nat Genet* 44(10) (2012), pp. 1166–1170.
- [41] P.-R. Loh et al. “Efficient Bayesian mixed-model analysis increases association power in large cohorts”. *Nat. Genet.* 47(3) (2015), pp. 284–290.
- [42] G. E. Hoffman. “Correcting for population structure and kinship using the linear mixed model: theory and extensions”. *PLoS ONE* 8(10) (2013), e75707.
- [43] G. Tucker, A. L. Price, and B. Berger. “Improving the Power of GWAS and Avoiding Confounding from Population Stratification with PC-Select”. *Genetics* 197(3) (2014), pp. 1045–1049.
- [44] N. Liu et al. “Controlling Population Structure in Human Genetic Association Studies with Samples of Unrelated Individuals”. *Stat Interface* 4(3) (2011), pp. 317–326.
- [45] J. Zeng et al. “Signatures of negative selection in the genetic architecture of human complex traits”. *Nature Genetics* 50(5) (2018), pp. 746–753.
- [46] M. Song, W. Hao, and J. D. Storey. “Testing for genetic associations in arbitrarily structured populations”. *Nat. Genet.* 47(5) (2015), pp. 550–554.
- [47] X. Liu et al. “Iterative Usage of Fixed and Random Effect Models for Powerful and Efficient Genome-Wide Association Studies”. *PLOS Genet* 12(2) (2016), e1005767.
- [48] J. H. Sul, L. S. Martin, and E. Eskin. “Population structure in genetic studies: Confounding factors and mixed models”. *PLoS Genet.* 14(12) (2018), e1007309.
- [49] A. L. Price et al. “Response to Sul and Eskin”. *Nature Reviews Genetics* 14(4) (2013), p. 300.

- [50] T. G. P. Consortium. “A map of human genome variation from population-scale sequencing”. *Nature* 467(7319) (2010), pp. 1061–1073.
- [51] 1000 Genomes Project Consortium et al. “An integrated map of genetic variation from 1,092 human genomes”. *Nature* 491(7422) (2012), pp. 56–65.
- [52] H. M. Cann et al. “A human genome diversity cell line panel”. *Science* 296(5566) (2002), pp. 261–262.
- [53] N. A. Rosenberg et al. “Genetic Structure of Human Populations”. *Science* 298(5602) (2002), pp. 2381–2385.
- [54] A. Bergström et al. “Insights into human genetic variation and population history from 929 diverse genomes”. *Science* 367(6484) (2020).
- [55] N. Patterson et al. “Ancient admixture in human history”. *Genetics* 192(3) (2012), pp. 1065–1093.
- [56] I. Lazaridis et al. “Ancient human genomes suggest three ancestral populations for present-day Europeans”. *Nature* 513(7518) (2014), pp. 409–413.
- [57] I. Lazaridis et al. “Genomic insights into the origin of farming in the ancient Near East”. *Nature* 536(7617) (2016), pp. 419–424.
- [58] P. Skoglund et al. “Genomic insights into the peopling of the Southwest Pacific”. *Nature* 538(7626) (2016), pp. 510–513.
- [59] J.-H. Park et al. “Distribution of allele frequencies and effect sizes and their interrelationships for common genetic susceptibility variants”. *PNAS* 108(44) (2011), pp. 18026–18031.
- [60] L. J. O’Connor et al. “Extreme Polygenicity of Complex Traits Is Explained by Negative Selection”. *The American Journal of Human Genetics* 0(0) (2019).
- [61] Y. B. Simons et al. “A population genetic interpretation of GWAS findings for human quantitative traits”. *PLOS Biology* 16(3) (2018), e2002985.
- [62] G. Malécot. *Mathématiques de l'hérédité*. Masson et Cie, 1948.
- [63] S. Wright. “The genetical structure of populations”. *Ann Eugen* 15(4) (1951), pp. 323–354.

- [64] A. Jacquard. *Structures génétiques des populations*. Paris: Masson et Cie, 1970.
- [65] A. Ochoa and J. D. Storey. “New kinship and  $F_{ST}$  estimates reveal higher levels of differentiation in the global human population”. *bioRxiv* (10.1101/653279) (2019).
- [66] C. C. Chang et al. “Second-generation PLINK: rising to the challenge of larger and richer datasets”. *GigaScience* 4(1) (2015), p. 7.
- [67] D. J. Balding and R. A. Nichols. “A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity”. *Genetica* 96(1-2) (1995), pp. 3–12.
- [68] E. Paradis and K. Schliep. “ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R”. *Bioinformatics* 35 (2019), pp. 526–528.
- [69] R. R. Sokal and C. D. Michener. “A statistical method for evaluating systematic relationships.” *Univ. Kansas, Sci. Bull.* 38 (1958), pp. 1409–1438.
- [70] C. L. Lawson and R. J. Hanson. “Solving least squares problems prentice-hall”. *Englewood Cliffs* (1974).
- [71] K. M. Mullen and I. H. M. v. Stokkum. *nnls: The Lawson-Hanson algorithm for non-negative least squares (NNLS)*. 2012.
- [72] J.-H. Park et al. “Estimation of effect size distribution from genome-wide association studies and implications for future discoveries”. *Nature Genetics* 42(7) (2010), pp. 570–575.
- [73] A. Grueneberg and G. d. l. Campos. “BGData - A Suite of R Packages for Genomic Analysis with Big Data”. *G3: Genes, Genomes, Genetics* 9(5) (2019), pp. 1377–1383.
- [74] S. Fairley et al. “The International Genome Sample Resource (IGSR) collection of open human genomic variation resources”. *Nucleic Acids Research* 48(D1) (2020), pp. D941–D947.
- [75] A. Manichaikul et al. “Robust relationship inference in genome-wide association studies”. *Bioinformatics* 26(22) (2010), pp. 2867–2873.
- [76] J. D. Storey. “The positive false discovery rate: a Bayesian interpretation and the q-value”. *Ann. Statist.* 31(6) (2003), pp. 2013–2035.

- [77] J. D. Storey and R. Tibshirani. “Statistical significance for genomewide studies”. *Proceedings of the National Academy of Sciences of the United States of America* 100(16) (2003), pp. 9440–9445.
- [78] J. Grau, I. Grosse, and J. Keilwagen. “PRROC: computing and visualizing precision-recall and receiver operating characteristic curves in R”. *Bioinformatics* 31(15) (2015), pp. 2595–2597.
- [79] P. Gopalan et al. “Scaling probabilistic models of genetic variation to millions of humans”. *Nat. Genet.* 48(12) (2016), pp. 1587–1590.
- [80] B. Rakitsch et al. “A Lasso multi-marker mixed model for association mapping with population structure correction”. *Bioinformatics* 29(2) (2013), pp. 206–214.
- [81] M. Conomos et al. “Model-free Estimation of Recent Genetic Relatedness”. *The American Journal of Human Genetics* 98(1) (2016), pp. 127–148.
- [82] S. Gazal et al. “High level of inbreeding in final phase of 1000 Genomes Project”. *Sci Rep* 5(1) (2015), p. 17453.
- [83] A. Al-Khudhair et al. “Inference of Distant Genetic Relations in Humans Using “1000 Genomes””. *Genome Biology and Evolution* 7(2) (2015), pp. 481–492.
- [84] L. Fedorova et al. “Atlas of Cryptic Genetic Relatedness Among 1000 Human Genomes”. *Genome Biology and Evolution* 8(3) (2016), pp. 777–790.
- [85] D. Schlauch, H. Fier, and C. Lange. “Identification of genetic outliers due to sub-structure and cryptic relationships”. *Bioinformatics* 33(13) (2017), pp. 1972–1979.
- [86] B. M. Henn et al. “Cryptic Distant Relatives Are Common in Both Isolated and Cosmopolitan Genetic Samples”. *PLOS ONE* 7(4) (2012), e34267.
- [87] V. Shchur and R. Nielsen. “On the number of siblings and p-th cousins in a large population sample”. *J Math Biol* 77(5) (2018), pp. 1279–1298.

## S1 Supplementary figures

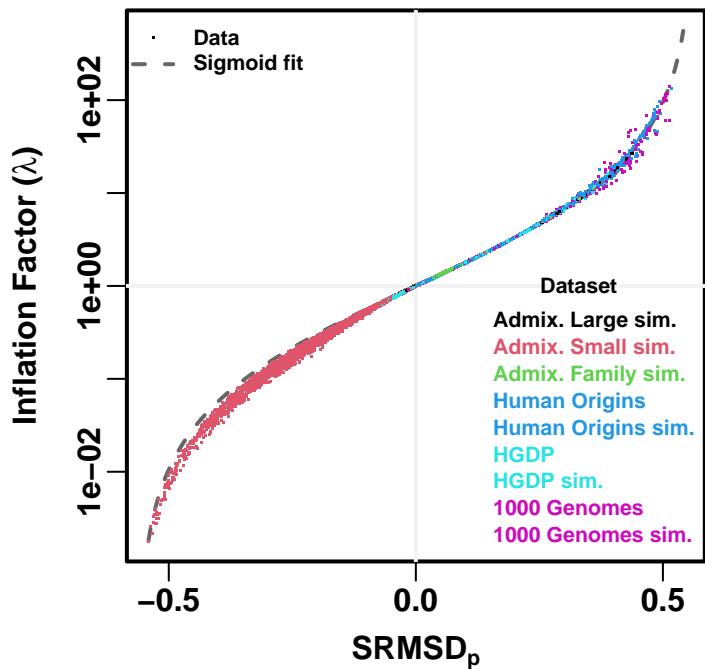


Figure S1: **Comparison between SRMSD<sub>p</sub> and inflation factor.** Each point is a pair of statistics for one replicate, one association model (PCA or LMM with some number of PCs  $r$ ), one trait model (FES vs RC), and one dataset (color coded by dataset). Note log y-axis ( $\lambda$ ). The sigmoidal curve in Eq. (10) is fit to the data.

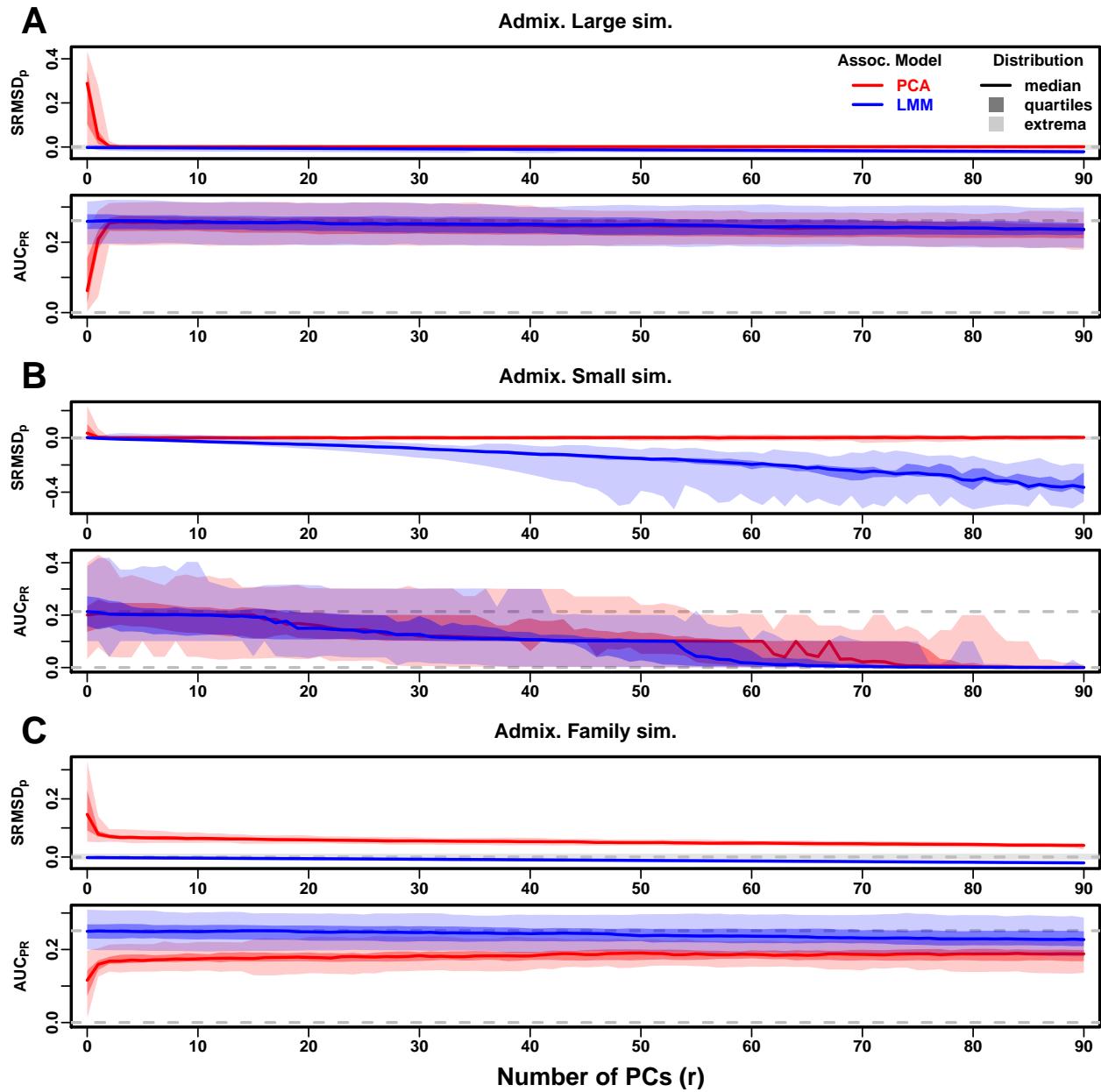


Figure S2: **Evaluations in admixture simulations.** Traits simulated from RC model, otherwise the same as Fig. 3.

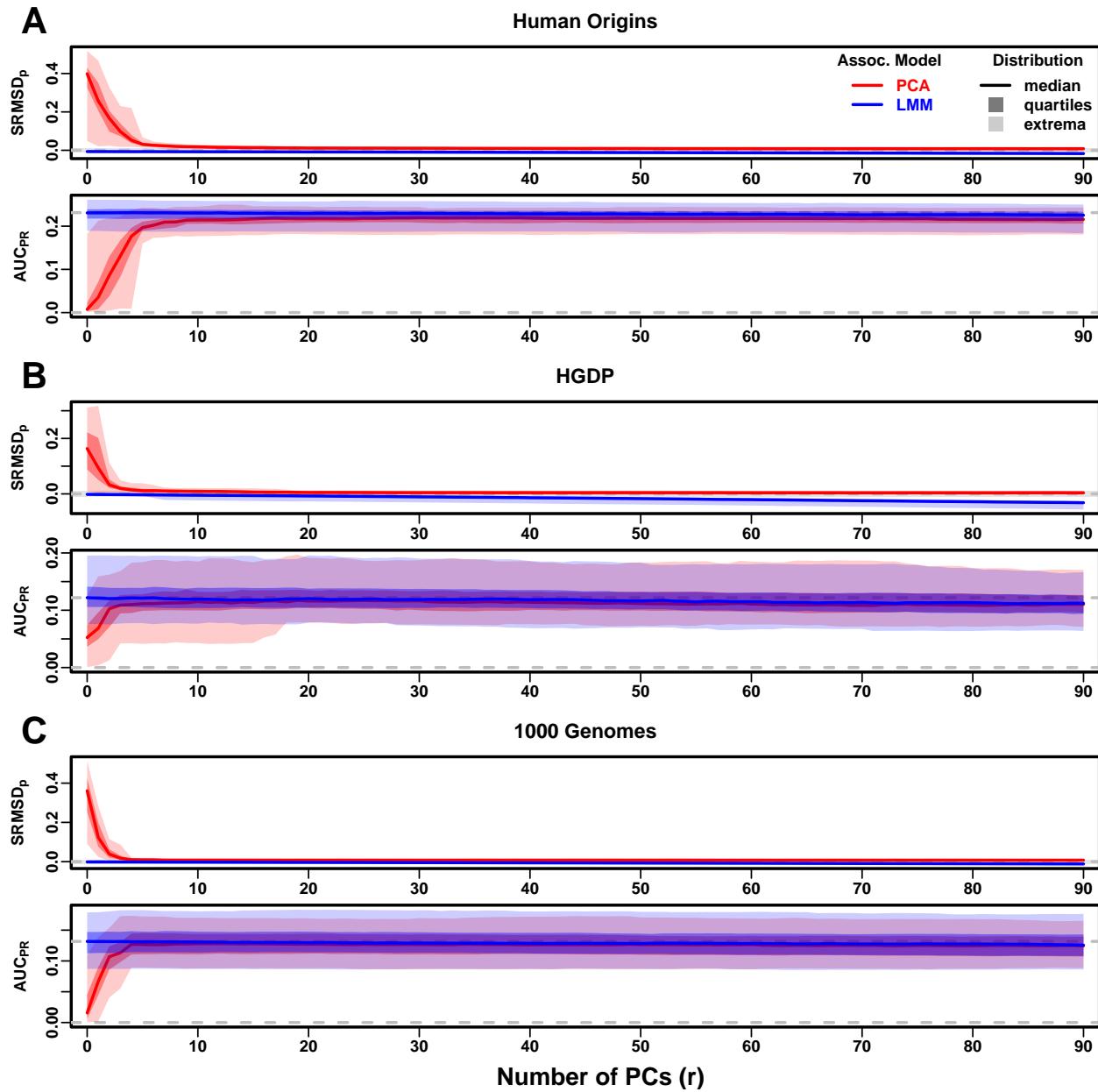


Figure S3: Evaluations in real human genotype datasets. Traits simulated from RC model, otherwise the same as Fig. 4.

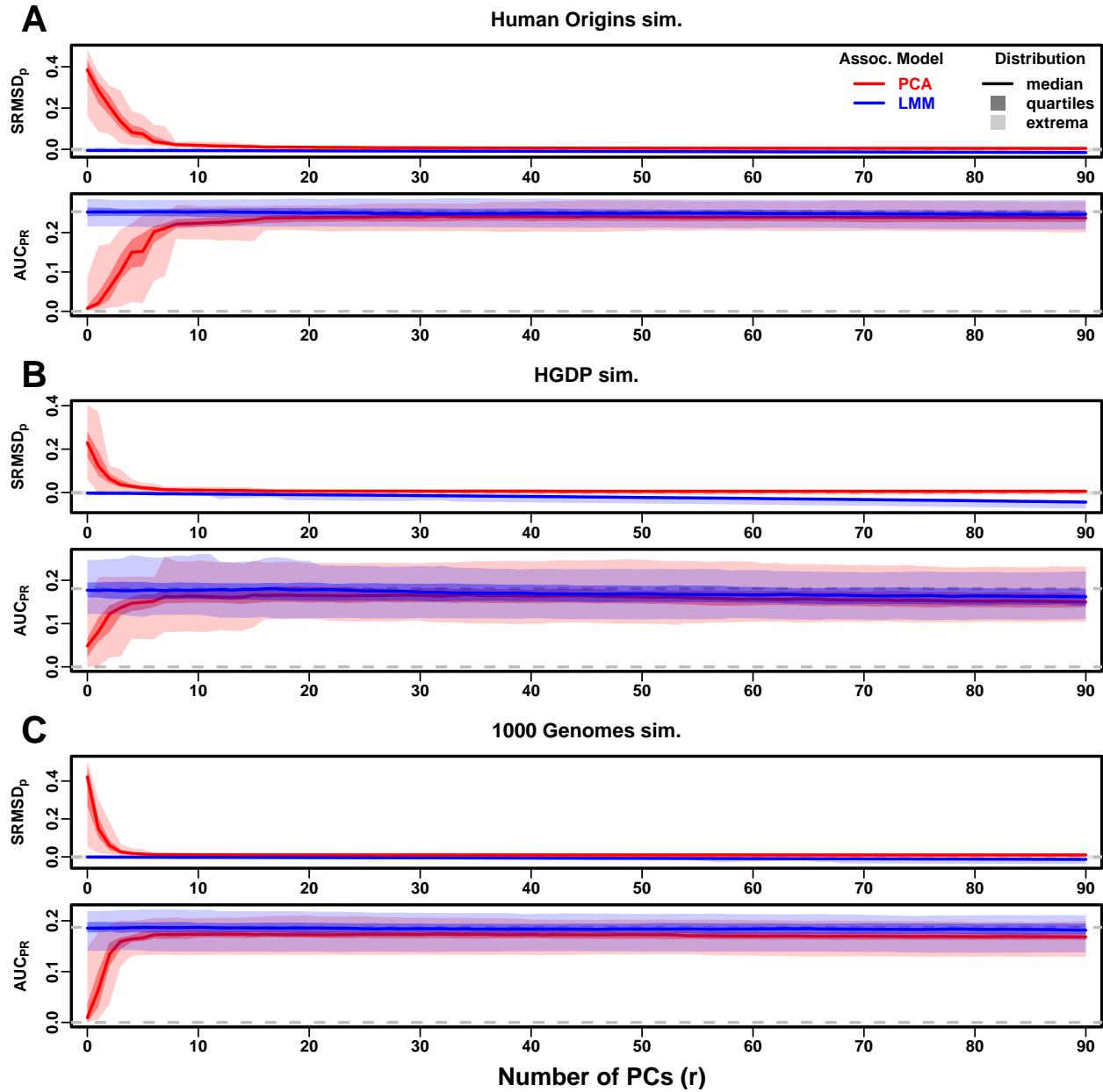


Figure S4: Evaluations in tree simulations fit to human data. Traits simulated from RC model, otherwise the same as Fig. 5.

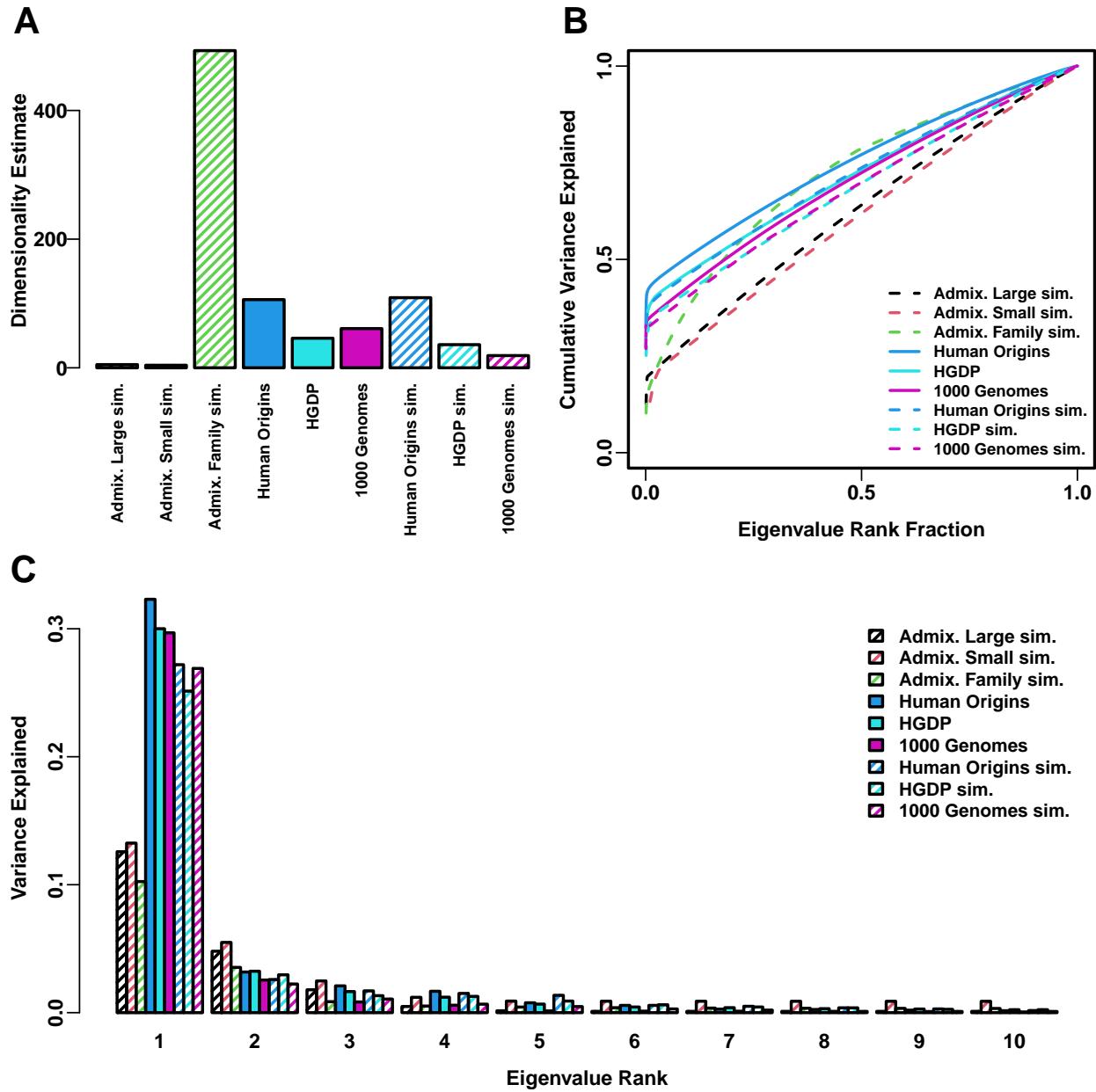


Figure S5: Estimated dimensionality of datasets. **A.** Kinship dimensionalities estimated with the Tracy-Widom test with  $p < 0.01$ . **B.** Cumulative variance explained versus eigenvalue rank fraction. **C.** Variance explained by first 10 eigenvalues.

## S2 Supplementary tables

Table S1: **Dataset sizes after 4th degree relative filter.**

Dataset	Loci ( $m$ )	Ind. ( $n$ )	Ind. removed (%)
Human Origins	189,722	2636	9.8
HGDP	905,838	842	9.4
1000 Genomes	1,097,415	2390	4.6