

Testing the effectiveness of principal components in adjusting for relatedness in genetic association studies

Yiqi Yao¹, 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

The genome-wide association study is widely used to identify the association between loci and trait. However, admixture population is involved in the current GWAS research which may result in the violation of Hardy-Weinberg Equilibrium (HWE). PCA can be used to correct for the bias result from admixture population. In this paper, we test the performance of PCA with desirable admixture simulation population. Since previous research about comparing PCA and LMM is limited, we also make comparison of performance among PCA and LMM in terms of both power and type one error controlling. According to the admixture simulation structure, three scenarios are considered in this case: large sample size, small sample size and large sample with complicated family structure. In order to fully compare the performance of PCA and LMM under different admixture structures, Root mean square deviation (RMSD) and areas under curve (AUC) to measure the performance of type one error controlling and power separately. According to the results of simulation, PCA fails in both type one error and power when the number of PCA is smaller than the true rank of genotype matrix. Additionally, PCA is robust to type one error once the number of PCs is larger than true rank of genotype matrix regardless of sample size, but the existence of family structure will impose a negative influence on it. In terms of power, PCA will be punished with excessive usage of PCs in the case of small sample size. Once enough PCs are used in PCA, it always outperforms LM in both type one error controlling and power. Regarding LMM, it has similar performance to PCA in type one error

controlling without complicated family structure, however, it has advantage over PCA when complex family structure exists. Considering power, LMM has better performance of power than PCA without complex family structure even enough PCs are used, whereas, the existence of complicated family structure will lead to PCA outperforms LMM in terms of power.

Contents

1	Introduction	4
2	Methods	5
2.1	Models for genetic association studies	5
2.1.1	The complex trait model and PCA approximation	6
2.1.2	Kinship model for genotypes	7
2.1.3	Estimation of principal components from genotype data	8
2.1.4	Linear mixed-effects model	8
2.2	Simulations	9
2.2.1	Genotype simulation from the admixture model	9
2.2.2	Genotype simulation from the family model	10
2.2.3	Trait Simulation	11
2.3	Evaluation of performance	12
2.3.1	RMSD_p : a measure of p-value uniformity	12
2.3.2	The area under the precision-recall curve	13
3	Results	14
4	Discussion	17

1 Introduction

The goal of a genome-wide association study (GWAS) is to identify loci whose genotypes are correlated significantly with a certain trait. An important assumption made by basic association tests is that genotypes at non-associated loci are drawn independently from a common allele frequency, so that they are in Hardy-Weinberg Equilibrium (HWE). However, HWE does not hold for structured populations, which includes multiethnic cohorts and admixed individuals, and for family data. When naive approaches are incorrectly applied to structured populations and/or family data, association statistics (such as χ^2) become inflated relative to the null expectation, resulting in greater numbers of false positives than expected (Devlin and Roeder, 1999; Voight and Pritchard, 2005; Astle and Balding, 2009).

Modern approaches for conducting genetic association studies with structured populations involve modeling the population structure via covariates. Such covariates may be inferred ancestry proportions (Pritchard et al., 2000) or transformations of these. Principal components analysis (PCA) represents the most common of these variants, in which the top eigenvectors of the kinship matrix are used to model the population structure (Price et al., 2006). These top eigenvectors are commonly referred to as Principal Components (PCs) in the genetics literature (the convention we adopt here; Patterson et al., 2006), but it is worth noting that in other fields the PCs would instead denote the projections of the data onto the eigenvectors. Various works have found that PCs map to ancestry, and PCs work as well as ancestry in GWAS and can be inferred more quickly (Patterson et al., 2006).

The other dominant approach for genetic association studies under population structure is the Linear Mixed-effect Model (LMM), in which population structure is a random effect drawn from a covariance model parametrized by the kinship matrix. LMM and PCA share deep connections that suggest that both models ought to perform similarly (Hoffman, 2013). However, many previous studies have found that LMM outperforms the PCA approach, although many evaluations are inconclusive or are limited to unrealistic population structures often with unrealistically low differentiation (Astle and Balding, 2009; Kang et al., 2010; Price et al., 2010; Wang et al., 2013). Moreover, various explanations for if and why LMM outperforms PCA are vague and have not been

tested directly (Price et al., 2010; Sul and Eskin, 2013; Price et al., 2013). Since LMMs tend to be considerably slower than the PCA approach, it is important to understand when the difference in performance between these two approaches is outweighed by their difference in runtime.

In this work, we study the performance of the PCA method for GWAS, comparing it to a leading LMM approach, characterizing its behavior under various numbers of PCs and varying sample sizes, under a reasonable admixture model and a model with admixture and family structure. Our evaluation is more thorough than previous ones, directly measuring the uniformity of null p-values (as required for accurate FDR control via q-values; Storey, 2003; Storey and Tibshirani, 2003) and orthogonally measuring predictive power by calculating the area under precision-recall curves. We find that the performance of the PCA approach is favorable when sample sizes are large (at least 1,000 individuals), matching the performance of LMMs as long as enough PCs are used. Remarkably, the approach is robust even when the number of PCs far exceeds the optimal number, suggesting that the degrees of freedom is not a concern in reasonably large studies. However, for smaller studies (100 individuals) there is a more pronounced loss of power when the number of PCs exceeds the optimal number. Moreover, LMMs outperform PCA in the presence of family structure, which is a well-known scenario where the problematic structure is not low-dimensional so PCA naturally cannot model it entirely (Patterson et al., 2006; Price et al., 2010). All together, our simulation studies provide clear criteria under which use of PCA results in acceptable performance compared to LMMs.

2 Methods

2.1 Models for genetic association studies

In this subsection we describe the complex trait model and kinship model that motivates both the PCA and LMM models for genetic association studies, followed by further details regarding the PCA and LMM approaches.

2.1.1 The complex trait model and PCA approximation

Let $x_{ij} \in \{0, 1, 2\}$ be the genotype at 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 which represents trait value for each individual. The approaches we consider are based on the following additive linear model for a quantitative (continuous) trait:

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

where $\mathbf{1}$ is a length- n vector of ones, α is the scalar intercept coefficient, β is the length- m vector of locus effect sizes, and ϵ is a length- n vector of residuals. The residuals are assumed to follow a normal distribution: $\epsilon_j \sim \text{Normal}(0, \sigma^2)$ independently for each individual j , for some residual variance parameter σ^2 .

Typically the number of loci m is in the order of millions while the number of individuals n is in the thousands. Hence, the full model above cannot be fit in this typical $n \ll m$ case, as there are only n datapoints to fit (the trait vector) but there are $m + 1$ parameters to fit (α and the β vector). The PCA model with r PCs corresponds to the following approximation to the full model, corresponding to a model fit at a single locus i :

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i \beta_i + \mathbf{U}_r \gamma_r + \epsilon, \quad (2)$$

where \mathbf{x}_i is the length- n vector of genotypes at locus i only, β_i is the effect size coefficient for that locus, \mathbf{U}_r is an $n \times r$ matrix of PCs, and γ_r is the length- r vector of coefficients for the PCs. This approximation is explained by first noticing that the genotype matrix has the following singular value decomposition: $\mathbf{X}^\top = \mathbf{U} \mathbf{D} \mathbf{V}^\top$, where assuming $n < m$ we have that \mathbf{U} is an $n \times n$ matrix of the left singular vectors of \mathbf{X} , \mathbf{V} is an $m \times n$ matrix of its right singular vectors, and \mathbf{D} is an $n \times n$ diagonal matrix of its singular values. Thus, in the full model we have $\mathbf{X}^\top \beta = \mathbf{U} \gamma$, where $\gamma = \mathbf{D} \mathbf{V}^\top \beta$ is a length- n vector. The approximation consists solely of replacing $\mathbf{U} \gamma$ (the full set of n left singular vectors and their coefficients) with $\mathbf{U}_r \gamma_r$ (the top r singular vectors only, which constitutes the best approximation of rank r). Thus, the extra terms in the PCA approach approximate the polygenic

effect of the whole genome, and assumes that the locus i being tested does not contribute greatly to this signal.

The statistical significance of a given association test is performed as follows. The null hypothesis is $\beta_j = 0$ (no association). The null and alternative models are each fit (fitting the coefficients of the multiple regression, where β_j is excluded under the null while it is fit under the alternative). The resulting regression residuals are compared to each other using the F-test, which results in a two-sided p-value. Note that many common PCA implementations trade the more exact F-test for a χ^2 test, which is simpler to implement but only asymptotically accurate. As this is a multiple hypothesis test, there are a large number of loci (m) tested for association, so it is best to control the FDR rather than setting a fixed p-value threshold. We recommend estimating q-values and setting a threshold of $q < 0.05$ so that the FDR is controlled at the 5% level.

2.1.2 Kinship model for genotypes

In order to better motivate the most common estimation procedure of PCs for genotype data, and to connect PCA to LMMs, we shall review the kinship model for genotypes. The model states that genotypes are random variables with a mean and covariance structure given by

$$\mathbb{E}[x_{ij}] = 2p_i, \quad \text{Cov}(x_{ij}, x_{ik}) = 4p_i(1 - p_i)\varphi_{jk},$$

where p_i is the ancestral allele frequency at locus i and φ_{jk} is the kinship coefficient between individuals j and k (Malécot, 1948; Wright, 1951; Jacquard, 1970). Thus, if we standardize the genotype matrix as

$$\mathbf{X}_S = \left(\frac{x_{ij} - 2p_i}{\sqrt{4p_i(1 - p_i)}} \right),$$

then this results in a straightforward kinship matrix estimator:

$$\mathbb{E} \left[\frac{1}{m} \mathbf{X}_S^T \mathbf{X}_S \right] = \mathbf{\Phi},$$

where $\mathbf{\Phi} = (\varphi_{jk})$ is the $n \times n$ kinship matrix. Note that replacing the raw genotype matrix \mathbf{X} with the standardized matrix \mathbf{X}_S in the trait model of Eq. (1) results in an equivalent model, as this

covariate differs only by a linear transformation. Thus, under the standardized genotype model, the PCs of interest are equal in expectation to the top eigenvectors of the kinship matrix.

2.1.3 Estimation of principal components from genotype data

In practice, the matrix of principal components \mathbf{U}_r in Eq. (2) is determined from an estimate of the earlier standardized genotype matrix \mathbf{X}_S , namely

$$\hat{\mathbf{X}}_S = \left(\frac{x_{ij} - 2\hat{p}_i}{\sqrt{4\hat{p}_i(1 - \hat{p}_i)}} \right),$$

where the true ancestral allele frequency p_i is replaced by the estimate $\hat{p}_i = \frac{1}{2n} \sum_{j=1}^n x_{ij}$, and results in the kinship estimate $\hat{\Phi} = \frac{1}{m} \hat{\mathbf{X}}_S^\top \hat{\mathbf{X}}_S$. This kinship estimate and minor variants are also employed in LMMs (Yang et al., 2011). This estimator of the kinship matrix is biased, and this bias is different for every individual pair (Ochoa and Storey, 2016b; Ochoa and Storey, 2018). However, in the present context of PCA regression in genetic association studies, the existing approach performs as well as when the above estimate is replaced by the true kinship matrix (not shown). Thus, it appears that in combination with the intercept term ($\mathbf{1}\alpha$ in Eq. (2)), the rowspace of this kinship matrix estimate approximately equals that of the true kinship matrix.

2.1.4 Linear mixed-effects model

The LMM is another approximation to the complex trait model in Eq. (1), given by

$$\mathbf{y} = \mathbf{1}\alpha + \mathbf{x}_i\beta_i + \mathbf{s} + \epsilon, \tag{3}$$

which is like the PCA model in Eq. (2) except that the PC terms $\mathbf{U}_r\gamma_r$ are replaced by the random effect \mathbf{s} , which is a length- n vector drawn from

$$\mathbf{s} \sim \text{Normal}(\mathbf{0}, \sigma_s^2 \Phi),$$

where Φ is the kinship matrix and σ_s^2 is a trait-specific variance scaling factor. This model is derived from treating the standardized genotype matrix \mathbf{X}_S as random rather than fixed, so that the standardized genetic effect $\mathbf{X}_S^\top \beta_S$ in Eq. (1) has mean zero and a covariance matrix of

$$\text{Cov}(\mathbf{X}_S^\top \beta_S) = \|\beta_S\|^2 \Phi.$$

The above random effect \mathbf{s} satisfies those equations, where the variance scale equals $\sigma_s^2 = \|\beta_S\|^2$. Thus, the PCA approach is the fixed model equivalent of the LMM under the additional approximation that only the top r eigenvectors are used in PCA whereas the LMM uses all eigenvectors.

A key advantage of LMM over PCA is that it has fewer degrees of freedom: ignoring the shared terms in Eq. (2) and Eq. (3), PCA has r parameters to fit (each PC coefficient in the γ vector), whereas LMMs only fit one additional parameter, namely σ_s^2 . Therefore, PCA is expected to overfit more substantially than LMM—and thus lose power—when r is very large, and especially when the sample size (the number of individuals n) is very small.

Due to its accuracy and speed, the LMM implementation that we chose for our evaluations is GCTA (Yang et al., 2011).

2.2 Simulations

2.2.1 Genotype simulation from the admixture model

We consider three simulation scenarios, referred to as (1) large sample size, (2) small sample size, and (3) family structure. All cases are based on the admixture model described previously (Ochoa and Storey, 2016a; Ochoa and Storey, 2016b), and which is implemented in the R package `bnpsd` available on GitHub and the Comprehensive R Archive Network (CRAN).

Here we consider scenarios where the number of individuals n varies: the large sample size and family structure scenarios have $n = 1,000$ whereas small sample size has $n = 100$. The number of loci in all cases is $m = 100,000$. Individuals are admixed from $K = 10$ intermediate subpopulations, where K is also the rank of the population structure; thus, after taking into account the intercept’s rank-1 contribution, the population structure can be fit with $r = K - 1$ PCs. Each subpopulation S_u

($u \in \{1, \dots, K\}$) has an inbreeding coefficient $f_{S_u} = u\tau$, individual-specific admixture proportions q_{ju} for individual j and intermediate subpopulation S_u arise from a random walk model for the intermediate subpopulations on a 1-dimensional geography with spread σ , where the free parameters τ and σ are fit to result in $F_{ST} = 0.1$ for the admixed individuals and a bias coefficient of $s = 0.5$, exactly as before (Ochoa and Storey, 2016b).

Random genotypes are drawn from this model, as follows. First, uniform ancestral allele frequencies p_i are drawn. The allele frequency $p_i^{S_u}$ at locus i of each intermediate subpopulation S_u is drawn from the Beta distribution with mean p_i and variance $p_i(1 - p_i)f_{S_u}$ (Balding and Nichols, 1995). The individual-specific allele frequency of individual j and locus i is given by $\pi_{ij} = \sum_{u=1}^K q_{ju}p_i^{S_u}$. Lastly, genotypes are drawn from $x_{ij} \sim \text{Binomial}(2, \pi_{ij})$. Loci that are fixed (where for some i we had $x_{ij} = 0$ for all j , or $x_{ij} = 2$ for all j) are drawn again from the model, starting from p_i , iterating until no loci are fixed.

2.2.2 Genotype simulation from the family model

Here we describe a simulation of a family structure with admixture that aims to be realistic by: (1) pairing all individuals in every generation, resulting in two children per couple; (2) strictly avoiding close relatives when pairing individuals; (3) strongly favoring pairs that are nearby in their 1-dimensional geography, which helps preserve the population structure across the generations by preferentially pairing individuals with more similar admixture proportions (a form of assortative mating); and (4) iterating for many generations so that a broad distribution of close and distant relatives is present in the data.

Generation 1 has individuals with genotypes drawn from the large sample size scenario described earlier, which features admixture. In subsequent generations, every individual is paired as follows. The local kinship matrix of individuals is stored and updated after every generation, which records the pedigree relatedness; in the first generation, everybody is locally unrelated. Also, individuals are ordered, initially by the 1-dimensional geography, and in subsequent generations paired individuals are grouped and reordered by their average coordinate, preserving the original order when there are ties. For every remaining unpaired individual, one is drawn randomly from the population, and

it is paired with the nearest individual that is not a second cousin or closer relative (local kinship must be $< 1/4^3$). Note that every individual is initially genderless, and after pairing one individual in the pair may be set to male and the other to female without giving rise to contradictions. If there are individuals that could not be paired (occurs if unpaired individuals are all close relatives), then the process of pairing individuals randomly is repeated entirely for this generation. If after 100 iterations no solution could be found randomly (there were always unpaired individuals), then the simulation restarts from the very first generation; this may occur for very small populations, but was not observed when $n = 1000$. Once individuals are paired, two children per pair have their genotypes drawn independently of each other. In particular, at every locus, one allele is drawn randomly from one of the parents and the other allele from the other parent. Loci are constructed independently of the rest (no linkage disequilibrium). The simulation continues for 20 generations. As this simulation is very computationally expensive, it was run only once (genotypes did not change as new random traits were constructed as described next).

2.2.3 Trait Simulation

For a given genotype matrix (simulated or real), a simulated complex trait that follows the additive quantitative trait model in Eq. (1) is constructed as follows. In all cases we set the heritability of the trait to be $h^2 = 0.8$. We varied the number of causal loci (m_1) together with the number of individuals (n) so power would remain balanced: for the $n = 1,000$ cases we set $m_1 = 100$, whereas the $n = 100$ simulation had $m_1 = 10$.

Each simulation replicate consists of different causal loci with different effect sizes, as follows. The non-genetic effects are drawn from $\epsilon_j \sim \text{Normal}(0, 1 - h^2)$ independently for each individual j . A subset of size m_1 of loci was selected at random from the genotype matrix to be causal loci. The effect size β_i at each causal locus i is drawn initially from a Standard Normal distribution. At non-causal loci i we have $\beta_i = 0$. Under the kinship model, the resulting genetic variance component is given by

$$\sigma_0^2 = \sum_{i=1}^m 2p_i(1 - p_i)\beta_i^2,$$

where p_i is the true ancestral allele frequency at locus i , which is known in our simulations. The

desired genetic variance of h^2 is therefore obtained by multiplying every β_i by $\frac{h}{\sigma_0}$. Lastly, the intercept coefficient in Eq. (1) is set to $\alpha = -\sum_{i=1}^m 2p_i\beta_i$, so the trait expectation is zero. This trait simulation procedure is implemented in the `simtrait` R package, available at <https://github.com/OchoaLab/simtrait>.

2.3 Evaluation of performance

All of the approaches considered here are evaluated in two orthogonal dimensions. The first one—the RMSD_p statistic below—quantifies the extent to which null p-values are uniform, which is a prerequisite for accurate control of the type-I error and successful FDR control via q-values. The second measure—the area under the precision-recall curve—quantifies the predictive power of each method, which makes it possible to qualitatively compare the statistical power of each method without having to select a single threshold, and most importantly, overcoming the problem that methods may not have accurate p-values.

2.3.1 RMSD_p : a measure of p-value uniformity

From their definition, correct p-values (for continuous test statistics) have a uniform distribution when the null hypothesis holds. This fact is crucial for accurate control of the type-I error, and is a prerequisite for the most common approaches that control the FDR, such as q-values (Storey, 2003; Storey and Tibshirani, 2003). We use the Root Mean Square Deviation (RMSD) to measure the disagreement between the observed p-value quantiles and the expected uniform quantiles:

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

where $m_0 = m - m_1$ is the number of null loci ($\beta_i = 0$ cases only), here i indexes null loci only, $p_{(i)}$ is the i th ordered null p-value, and $u_i = (i - 0.5)/m_0$ is its expectation. Thus, $\text{RMSD}_p = 0$ corresponds to the best performance in this test, and larger RMSD_p values correspond to worse performance.

In previous evaluations, test statistic inflation has been used to measure the success of corrections

for population structure (Astle and Balding, 2009; Price et al., 2010). The inflation factor λ is defined as the median χ^2 association statistic divided by theoretical median under the null hypothesis (Devlin and Roeder, 1999). Hence, when null test statistics have their expected distribution, we get $\lambda = 1$ (same as $\text{RMSD}_p = 0$ above). However, any other null test statistic distribution with the same median results in $\lambda = 1$ as well, which is a flaw of this test that RMSD_p overcomes ($\text{RMSD}_p = 0$ if and only if null test statistics have their expected distribution). The $\lambda > 1$ case (gives $\text{RMSD}_p > 0$) corresponds to inflated statistics, which occurs when residual population structure is present. $\lambda < 1$ is not expected for genetic association studies (also gives $\text{RMSD}_p > 0$). Note that λ only use the median of the null distribution, whereas the RMSD_p makes use of the complete p-value distribution to evaluate its uniformity, which is more accurate.

2.3.2 The area under the precision-recall curve

Precision and recall are two common measures for evaluating binary classifiers. Let c_i be the true classification of locus i , where $c_i = 1$ for truly causal loci (if the true $\beta_i \neq 0$, where the alternative hypothesis holds), and $c_i = 0$ otherwise (null cases). For a given method and some threshold t on its per-locus test statistics, the method predicts a classification $\hat{c}_i(t)$ (for example, if t_i is the test statistic, the prediction could be $\hat{c}_i(t) = 1$ if $t_i \geq t$, and $\hat{c}_i(t) = 0$ otherwise). Across all loci, the number of true positives (TP), false positives (FP) and false negatives (FN) at the given threshold t is given by

$$\begin{aligned}\text{TP}(t) &= \sum_{i=1}^m c_i \hat{c}_i(t), \\ \text{FP}(t) &= \sum_{i=1}^m (1 - c_i) \hat{c}_i(t), \\ \text{FN}(t) &= \sum_{i=1}^m c_i (1 - \hat{c}_i(t)).\end{aligned}$$

Precision and recall at this threshold are given by

$$\begin{aligned}\text{Precision}(t) &= \frac{\text{TP}(t)}{\text{TP}(t) + \text{FP}(t)} = \frac{\sum_{i=1}^m c_i \hat{c}_i(t)}{\sum_{i=1}^m \hat{c}_i(t)}, \\ \text{Recall}(t) &= \frac{\text{TP}(t)}{\text{TP}(t) + \text{FN}(t)} = \frac{\sum_{i=1}^m c_i \hat{c}_i(t)}{\sum_{i=1}^m c_i}.\end{aligned}$$

The precision-recall curve results from calculating the above two values at every threshold t , tracing a curve as recall goes from zero (everything is classified as null) to one (everything is classified as alternative), and the area under this curve is our final measure AUC_{PR} . A method obtains the maximum $\text{AUC}_{\text{PR}} = 1$ if there is some threshold that classifies all loci perfectly. In contrast, a method that classifies at random (for example, $\hat{c}_i(t) \sim \text{Bernoulli}(p)$ for any p) has an expected precision (= AUC_{PR}) approximately equal to the overall proportion of alternative cases: $\pi_1 = \frac{m_1}{m} = \frac{1}{m} \sum_{i=1}^m c_i$. The AUC_{PR} was calculated using the R package `PRROC`, which computes the area by integrating the correct non-linear piecewise function when interpolating between points (Grau et al., 2015).

3 Results

We simulate genotype matrices and traits to go with the genotypes, in order to control important features of the population structure and to test all methods in an ideal setting where the true causal loci are known. Our simulations permit exact identification of true positives, false positives, and false negatives, ultimately yielding two measures of interest: RMSD_p measure null p-value uniformity and relates to the accuracy of type-I error control (smaller is better), while AUC_{PR} measures predictive power (higher is better) and serves as a proxy for statistical power when $\text{RMSD}_p \approx 0$. However, the simulation of genotypes followed by simulation of the trait leads to a considerable amount of variance in the final measured RMSD_p and AUC_{PR} , which are random variables. For that reason, every evaluation was replicated at least 10 times (varies by scenario), resulting in a distribution of RMSD_p and AUC_{PR} values per method. Except when noted, each replicate consisted of a new genotype matrix drawn from the same structure model of the scenario, followed by a new simulated trait based on this genotype matrix, which included selecting new causal loci with new effect sizes.

All scenarios are based on an admixture simulation from $K = 10$ subpopulations and a resulting generalized $F_{ST} = 0.1$, which establishes the population structure. We vary the sample size (number of individuals) in order to test the extent to which PCA overfits the population structure as the number of PCs increases ($r \in \{0, \dots, 90\}$), particularly in comparison to the LMM. Keep in mind that the ideal choice for the number of PCs in this simulation is $r = K - 1 = 9$ (the rank of the data minus the rank of the intercept). Lastly, to push all methods to their limits, we evaluate them in a scenario with both admixture and a complex family structure.

First we evaluate all methods in the large sample size scenario, which has a reasonable number of individuals ($n = 1,000$) typical for genetic association studies. In this scenario we find a clear transition around the ideal number of PCs of $r = 9$, below of which performance is poor and above of which performance is satisfactory (Fig. 1). In particular, when $r < 9$ we find the largest RMSD_p values, which indicate that p-values are highly non-uniform and would therefore result in inaccurate type-I error control. The smallest AUC_{PR} values also occur for $r < 9$, showing that not enough PCs results in loss of predictive power as well. As expected, $r = 9$ has the best performance in terms of both RMSD_p and AUC_{PR} . Remarkably, as r is increased up to $r = 90$, there is no noticeable change in the RMSD_p distribution, and only a small decrease in AUC_{PR} compared to the optimal $r = 9$ case. The LMM performs about as well as PCA with $r = 9$ here, with small RMSD_p values (though somewhat larger than those of $r = 9$) and larger AUC_{PR} values than PCA's with $r = 9$. Thus, in this common scenario where sample sizes are large enough, the PCA approach with enough PCs performs as well as LMM.

The previous observation, that PCA continues to perform well when the number of PCs is 10 times greater than its optimum value ($r = 90$ vs $r = 9$), propelled us to find a scenario where this is no longer the case. We expect the PCA approach to begin overfitting as the number of PCs r approaches the sample size n . Increasing r beyond 90 does not make sense, as this would never be done in practice. Instead, we reduced n to 100, a number of individuals that is small for typical association studies, but which may occur in studies of rare diseases, or be due to low budgets or other constraints. To compensate for the loss of power that results from reducing the sample size, we also reduced the number of causal loci from 100 before to $m_1 = 10$, which increases the magnitude

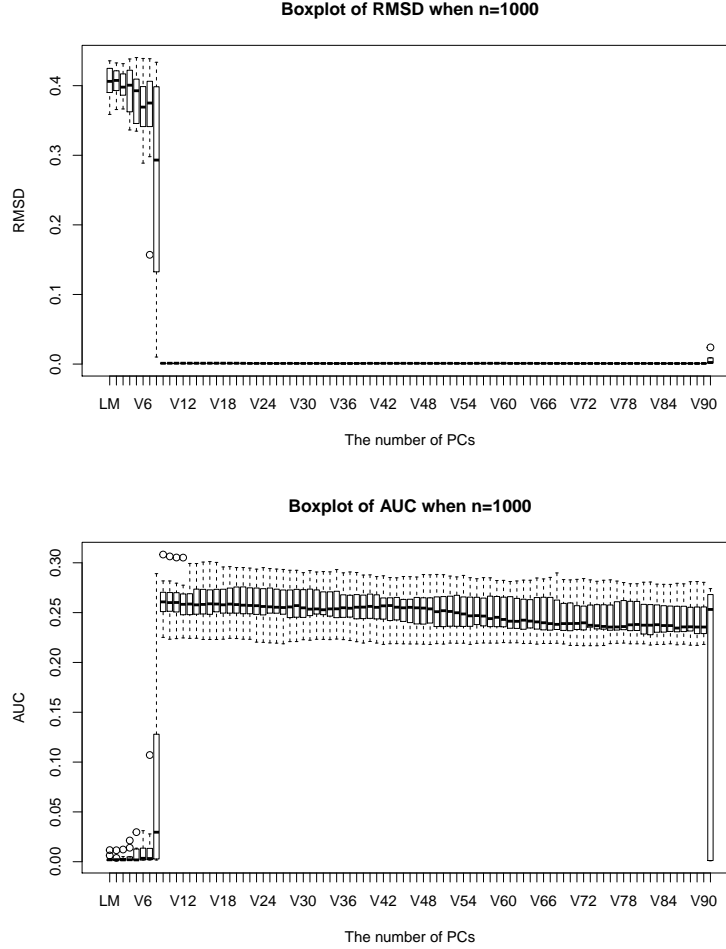


Figure 1: **Evaluation in large sample size admixture scenario.** Here there are $n = 1,000$ individuals in the simulation. The PCA approach is tested under varying number of PCs ($r \in \{0, \dots, 90\}$), alongside the LMM approach (x-axis), with boxplots for 10 replicates (y-axis) for the distributions of RMSD_p (top panel) and AUC_{PR} (bottom panel). Small RMSD_p and large AUC_{PR} correspond with better performance. The ideal number of PCs is $r = K - 1 = 9$, where K is the number of subpopulations prior to admixture, which results in near zero RMSD_p and peak AUC_{PR} , and performs nearly as well as the LMM. PCA with $r < 9$ has incorrect p-values ($\text{RMSD}_p \gg 0$ cases) and lowest predictive power (small AUC_{PR}). Remarkably, PCA remains robust even in extreme $r > 9$ cases, with RMSD_p near zero up to $r = 90$ and minimal loss of power as r increases to 90.

of the effect sizes. Note that this reduction in the number of causal loci results in more discreteness in AUC_{PR} values in Fig. 2. Interestingly, we find that the relationship between RMSD_p and r is similar under small and large sample sizes, with ideal near-zero RMSD_p distributions for $r \geq 9$. On the other hand, we do see a more severe overfitting effect here that results in decreased predictive power: AUC_{PR} peaks at $r = 9$ as expected, but drops more rapidly as r increases, with performance around $r = 50$ that is as bad as for $r = 0$, and practically zero AUC_{PR} at $r = 90$ (Fig. 2). Another notable difference from the large sample size scenario is that here LMM outperforms PCA with $r = 9$ by a sizable margin.

Previous work has shown that PCA performs poorly in the presence of family structure. Here we aim to characterize PCA’s behavior in a much more complex structure than before, by simulating a family of admixed founders for 20 generations, so that we may observe numerous siblings, first cousins, etc. In this case $r = 9$ is not the optimal choice, as the rank of the genotype matrix is much greater due to the family structure. We find that, although RMSD_p decreases monotonically as r increases, this distribution does not go to zero, instead converging to around 0.05 (Fig. 3). Additionally, the AUC_{PR} increases until $r = 4$ is reached (as opposed to $r = 9$ as before), then plateaus with marginal decreases in performance as r goes to 90. In contrast, the LMM does achieve a near-zero RMSD_p , although the AUC_{PR} distribution is much wider than the best performing PCA cases.

4 Discussion

Right now, both LMM and PCA have become standard approaches to correct for admixture population. In current PCA GWAS research, the number of PCs used in PCA is usually assumed to be 10. The default in EIGENSTRAT is to use 10 PCs (Price et al., 2006). For instance, based on the simulation of Hoffman (2013), 10 PCs are randomly selected from the first 30 PCs. Furthermore, Wojcik et al. (2019) also performed PCA GWAS for 26 traits with first 10 PCs. According to the result of Wojcik et al. (2019), the correlation plot between SNP genotype and PC1–PC10 illustrates different populations has different correlations over some PCs. Here, we further investigate about the optimal choice of number of PCA under different population structures. In most

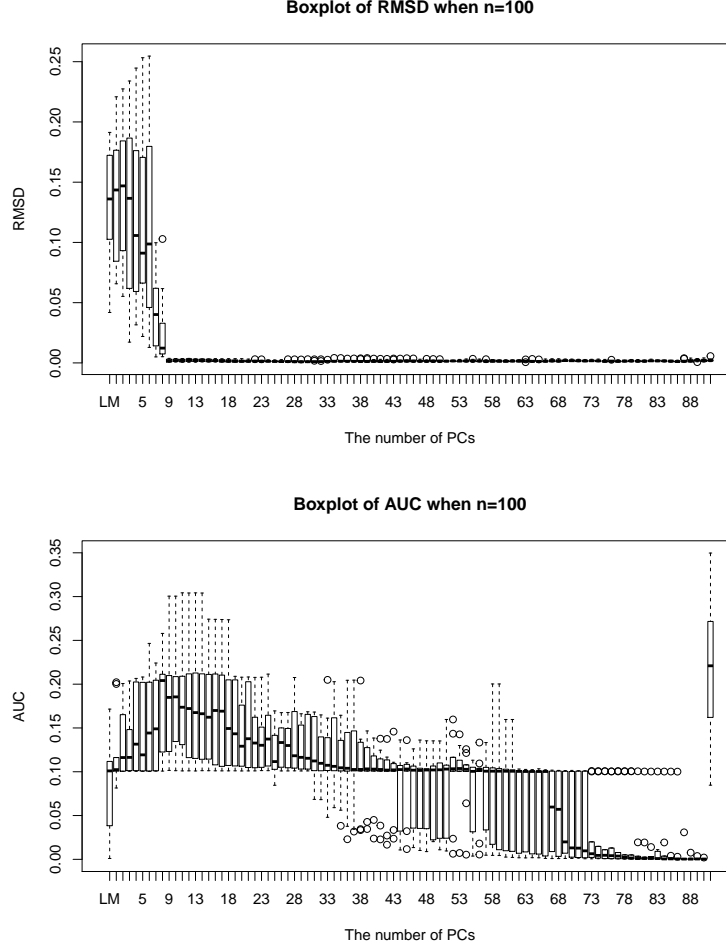


Figure 2: **Evaluation in small sample size admixture scenario.** Here there are $n = 100$ individuals in the simulation, otherwise the simulation and figure layout is the same as in Fig. 1. The pattern for RMSD_p in the top panel is similar to the previous figure. However, here there is a more pronounced drop in AUC_{PR} values as the number of PCs r increases from $r = 9$ to $r = 90$. Here LMM outperforms PCA with $r = 9$ in terms of AUC_{PR} by a greater margin.

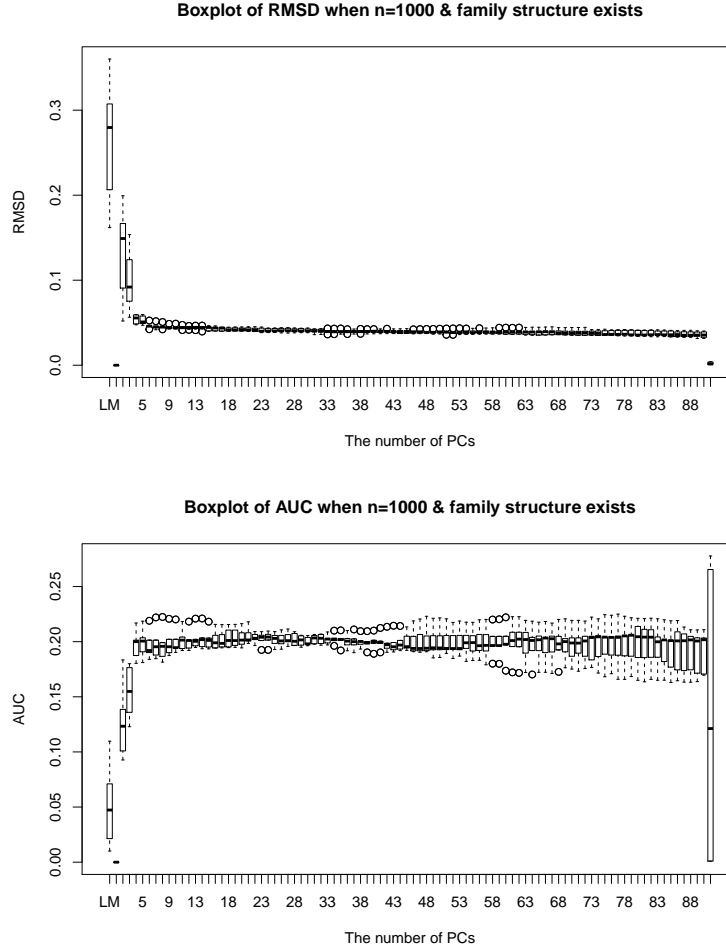


Figure 3: **Evaluation in family structure admixture scenario.** Here there are $n = 1,000$ individuals from a family structure simulation with admixed founders and large numbers of pairs of sibling, first cousins, second cousins, etc, from a realistic random pedigree that spans 20 generations. Unlike previous figures, here RMSD_p (top panel) for PCA does not go down to zero as r increases. For this complex relatedness structure $r = 9$ is not the optimal number of PCs, although performance remains steady for all $r \geq 9$ values tested. TODO: LMM looks weird and $r = 1$ is missing!

current research, the number of subpopulation is smaller 10 and thus, based on results of previous three scenarios, the number of PCs used in PCA is enough which certifies the good performance of PCA in current research. For instance, based on the scatter plot of first two principal components with HapMap3 dataset which contains 11 populations, it can be divided into three subpopulations (Abraham and Inouye, 2014). On the contrary, the lack of PCs will lead to the failure of PCA regardless of the sample size or family structure in both type one error and power. Lee et al. (2012) state that considering the large number of single-nucleotide polymorphisms (SNP) used in GWAS to infer structure, it is necessary to remove SNPs that have negligible loadings in PCA. Whereas, our simulation indicates that if there is no enough PCs used in PCA GWAS, its performance can be bad as SNPs that have significant loadings are vanished from the analysis.

PCA GWAS is still robust even though PCs are excessively used for type one error controlling. However, for power, it will receive punishment if the number of PCs are excessive and the sample size is small. The small sample size in our research project only has 100 individuals in total. In a GWAS study, this is an extreme situation which it not realistic in research. The AUC boxplot of small sample size indicates that the peak of AUC is close to AUC value for large sample size. This may result from that the number of causal loci in small sample size is reduced from 100 to 10. However, the boxplot of AUC has a bell shap with downward-sloping line on each side of the peak. It demonstrates that in small sample size, punishment of excessive use of PCs will come occur immedieately. Considering in the study of GWAS, we will expect to use thousnds of SNPs and number of individuals are also much larger than 100, the situation of small sample size may not be common (Lee et al., 2012) Hence, although PCA will fail when the number of PCs is much larger than the true rank of genotype matrix, it is still encouraged to use more PCs in PCA GWAS.

(Price et al., 2010) point out that GWAS may fail in the case of dataset contains family structure. The result of our simulation also supports this argument. Fig. 3 shows that PCA has worse performance in type one error controlling when family strcuture exists, compared with Fig. 1. Without complicated family structure, the vlaue of RMSD will converges to 0 when the number of PCs is large enough. However, in the case of complex family structure, it can be seen that RMSD converges to 0.05 which indicates that we do not have strong evidence to claim that type one error is

excellently controlled. Concerning power, it can be seen that in both situations, AUC will converge to 0.2. Whereas, the range of AUC in Fig. 3 will increase when excessive PCs are used in PCA. It illustrates that excessive use of PCs will result in the extra variance. Hence, in the case of complex family structure, we need to be cautious about using PCs to avoid unnecessary variance.

Wang et al. (2013) argues that mixed effects model is preferred in the case of existence cryptic relatedness but not population stratification. In their paper, only first four PCs are used in PCA and performance of PCA may be underestimated. Furthermore, EMMAX which is a kind of linear mixed model is claimed to be better than PCA (TODO review: Gengxin and Hongjiang, 2013). Based on the result of our simulations, it can be seen that in both large sample size and small sample size, LMM are slightly better than PCA in terms of power. When sample size is large such as the scenario in Fig. 1B, it can be seen that although both two methods's are not good. However, the maximum AUC value of PCA is around 0.25 and LMM's AUC value is slightly larger than PCA. In addition to this, in the case of small sample size, advantage of LMM is more obvious. As mentioned before, PCA will receive punishment when the number of PCs is much larger than the true rank of genotype matrix. The maximum of AUC value of PCA in this scenario decreases to 0.2 and the average AUC value of LMM is around 0.24. However, when complicated family structure exists, PCA outperforms LMM in terms of power.

(Tucker et al., 2014)

TODO: This is a review: Gengxin and Hongjiang (2013) demonstrate that efficient mixed-model association eXpedited (EMMAX) which is based on linear mixed model outperforms PCA in both the population cohort study and case-control study.

References

- Abraham, Gad and Michael Inouye (9, 2014). "Fast Principal Component Analysis of Large-Scale Genome-Wide Data". *PLOS ONE* 9(4), e93766.
- Astle, William and David J. Balding (2009). "Population Structure and Cryptic Relatedness in Genetic Association Studies". *Statist. Sci.* 24(4). Mathematical Reviews number (MathSciNet): MR2779337, pp. 451–471.

- 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), pp. 3–12.
- Devlin, B. and Kathryn Roeder (1, 1999). “Genomic Control for Association Studies”. *Biometrics* 55(4), pp. 997–1004.
- Grau, Jan, Ivo Grosse, and Jens Keilwagen (1, 2015). “PRROC: computing and visualizing precision-recall and receiver operating characteristic curves in R”. *Bioinformatics* 31(15), pp. 2595–2597.
- 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. (2010). “Variance component model to account for sample structure in genome-wide association studies”. *Nat. Genet.* 42(4), pp. 348–354.
- Lee, Seokho et al. (2012). “Sparse Principal Component Analysis for Identifying Ancestry-Informative Markers in Genome-Wide Association Studies”. *Genetic Epidemiology* 36(4), pp. 293–302.
- Malécot, Gustave (1948). *Mathématiques de l’hérédité*. Masson et Cie.
- Ochoa, Alejandro and John D. Storey (2016a). “ F_{ST} and kinship for arbitrary population structures I: Generalized definitions”. Submitted, preprint at <http://biorxiv.org/content/early/2016/10/27/083915>.
- (2016b). “ F_{ST} and kinship for arbitrary population structures II: Method of moments estimators”. Submitted, preprint at <http://biorxiv.org/content/early/2016/10/27/083923>.
- (2018). “New kinship and F_{ST} estimates reveal higher levels of differentiation in the world-wide human population”. Submitted, preprint at <http://biorxiv.org/content/early/...>
- Patterson, Nick, Alkes L Price, and David Reich (22, 2006). “Population Structure and Eigenanalysis”. *PLoS Genet* 2(12), e190.
- Price, Alkes L. et al. (2006). “Principal components analysis corrects for stratification in genome-wide association studies”. *Nat. Genet.* 38(8), pp. 904–909.
- Price, Alkes L. et al. (2010). “New approaches to population stratification in genome-wide association studies”. *Nature Reviews Genetics* 11(7), pp. 459–463.

- Price, Alkes L. et al. (2013). “Response to Sul and Eskin”. *Nature Reviews Genetics* 14(4), p. 300.
- Pritchard, Jonathan K. et al. (2000). “Association Mapping in Structured Populations”. *The American Journal of Human Genetics* 67(1), pp. 170–181.
- Storey, John D. (2003). “The positive false discovery rate: a Bayesian interpretation and the q-value”. *Ann. Statist.* 31(6). Mathematical Reviews number (MathSciNet): MR2036398; Zentralblatt MATH identifier: 02067675, pp. 2013–2035.
- Storey, John D. and Robert Tibshirani (2003). “Statistical significance for genomewide studies”. *Proceedings of the National Academy of Sciences of the United States of America* 100(16), pp. 9440–9445.
- Sul, Jae Hoon and Eleazar Eskin (2013). “Mixed models can correct for population structure for genomic regions under selection”. *Nature Reviews Genetics* 14(4), p. 300.
- Tucker, George, Alkes L. Price, and Bonnie Berger (1, 2014). “Improving the Power of GWAS and Avoiding Confounding from Population Stratification with PC-Select”. *Genetics* 197(3), pp. 1045–1049.
- Voight, Benjamin F. and Jonathan K. Pritchard (2, 2005). “Confounding from Cryptic Relatedness in Case-Control Association Studies”. *PLOS Genetics* 1(3), e32.
- Wang, Kai, Xijian Hu, and Yingwei Peng (2013). “An Analytical Comparison of the Principal Component Method and the Mixed Effects Model for Association Studies in the Presence of Cryptic Relatedness and Population Stratification”. *HHE* 76(1), pp. 1–9.
- Wojcik, Genevieve L. et al. (2019). “Genetic analyses of diverse populations improves discovery for complex traits”. *Nature* 570(7762), pp. 514–518.
- Wright, S. (1951). “The genetical structure of populations”. *Ann Eugen* 15(4), pp. 323–354.
- Yang, Jian et al. (7, 2011). “GCTA: a tool for genome-wide complex trait analysis”. *Am. J. Hum. Genet.* 88(1), pp. 76–82.