

# Genetic association models are robust to common population kinship estimation biases

Zhuoran Hou<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



## Abstract

- Background:** commonly kinship estimators can have severe biases.
- Results:** kinship matrices of different bias types result in equal association statistics and performance in simulations and 1000 Genomes.
- Intercept and relatedness (PCs in Principal Component Analysis (PCA), random effect in Linear Mixed-effects Models (LMM)) coefficients compensate for the kinship bias.

## Model

$x_{ij} \in \{0, 1, 2\}$ : genotype of ind.  $j$ , biallelic SNP  $i$ .

$p_i$ : ancestral allele frequency.  $\varphi_{ij}$ : kinship coefficient

$$E[x_i] = 2p_i \mathbf{1}, \quad \text{Cov}(x_i) = 4p_i(1 - p_i)\Phi$$

where  $x_i = (x_{ij})$  is the length- $n$  column vector of genotypes at locus  $i$ ,  $\Phi = (\varphi_{ij})$  is the  $n \times n$  kinship matrix, and  $\mathbf{1}$  is a length- $n$  column vector of ones [1].

## Kinship estimators

### Standard estimator

Ratio-of-means (ROM) [1,2]:

$$\hat{p}_i = \frac{1}{2n} \sum_{j=1}^n x_{ij}, \quad \hat{\varphi}_{ij}^{\text{std-ROM}} = \frac{\sum_{k=1}^m (x_{ij} - 2\hat{p}_i)(x_{ik} - 2\hat{p}_i)}{\sum_{k=1}^m 4\hat{p}_i(1 - \hat{p}_i)} \xrightarrow[m \rightarrow \infty]{a.s.} \frac{\varphi_{ij} - \bar{\varphi}_j - \bar{\varphi}_k + \bar{\varphi}}{1 - \bar{\varphi}}$$

Mean-of-ratios (MOR)

$$\hat{\varphi}_{ij}^{\text{std-MOR}} = \frac{1}{m} \sum_{k=1}^m \frac{(x_{ij} - 2\hat{p}_i)(x_{ik} - 2\hat{p}_i)}{4\hat{p}_i(1 - \hat{p}_i)}$$

### Popkin estimator[1]

$$A_{ij} = \frac{1}{m} \sum_{k=1}^m w_i ((x_{ij} - 1)(x_{ik} - 1) - 1), \quad \hat{\varphi}_{ij}^{\text{popkin}} = 1 - \frac{A_{ij}}{\min_{j \neq k} A_{jk}},$$

$$w_i = 1 \text{ for ROM, } \hat{\varphi}_{ij}^{\text{popkin-ROM}} \xrightarrow[m \rightarrow \infty]{a.s.} \varphi_{ij},$$

$$w_i = (\hat{p}_i(1 - \hat{p}_i))^{-1} \text{ for MOR.}$$

## Association models

LMM [3,4]:  $y = \mathbf{1}\alpha + x_i\beta_i + s + \epsilon, \quad s \sim \text{Normal}(\mathbf{0}, 2\sigma^2\Phi)$

PCA [3,4]:  $y = \mathbf{1}\alpha + x_i\beta_i + U_d\gamma_d + \epsilon, \quad \Phi = U\Lambda U^T$

## Empirical analysis using 1000 Genomes

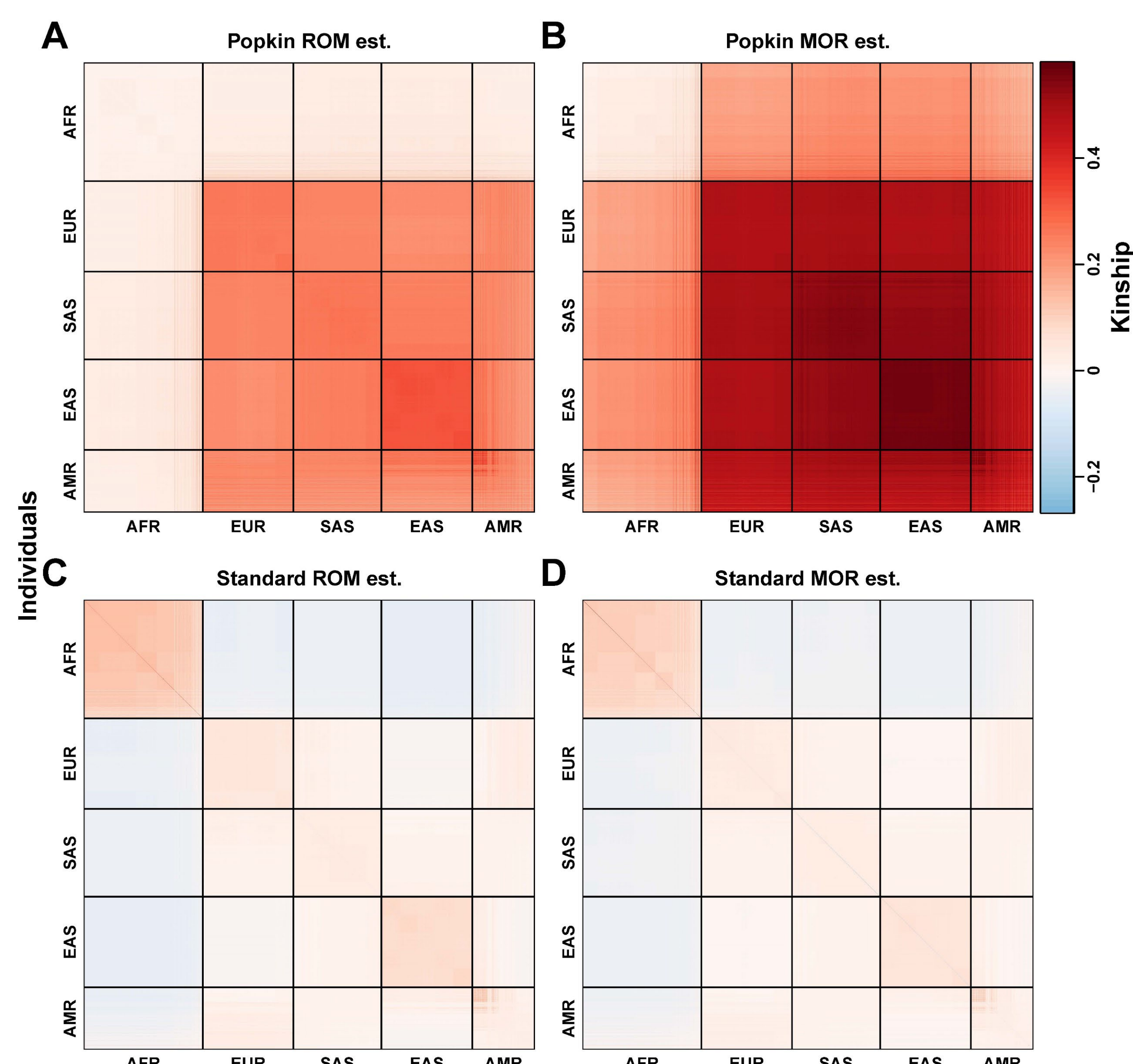


Figure 1: Kinship estimates on 1000 Genomes

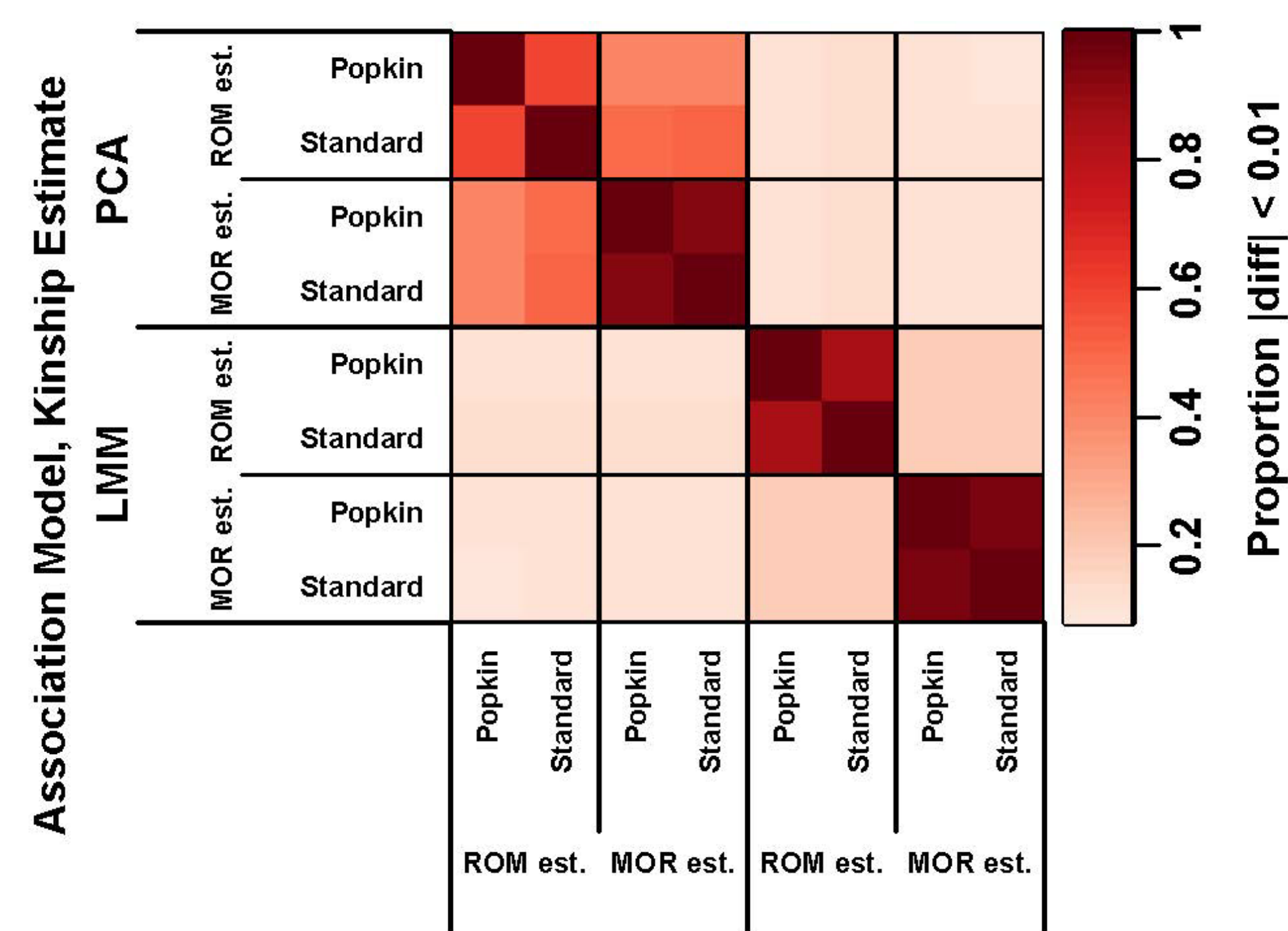


Figure 2: Approximate agreement between p-values on 1000 Genomes

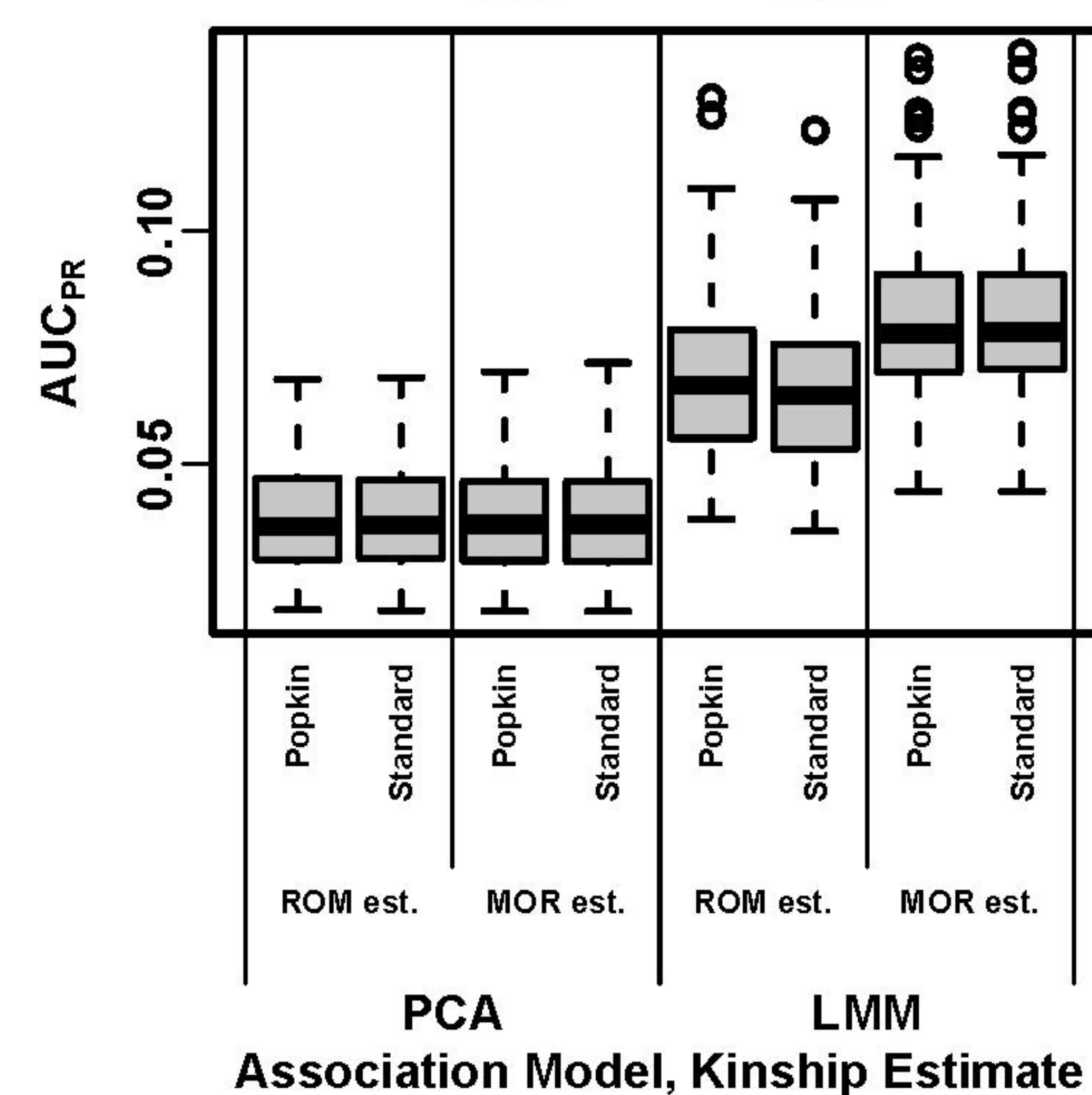


Figure 3: Distributions of Area Under the Precision-Recall Curve (AUC<sub>PR</sub>) on 1000 Genomes

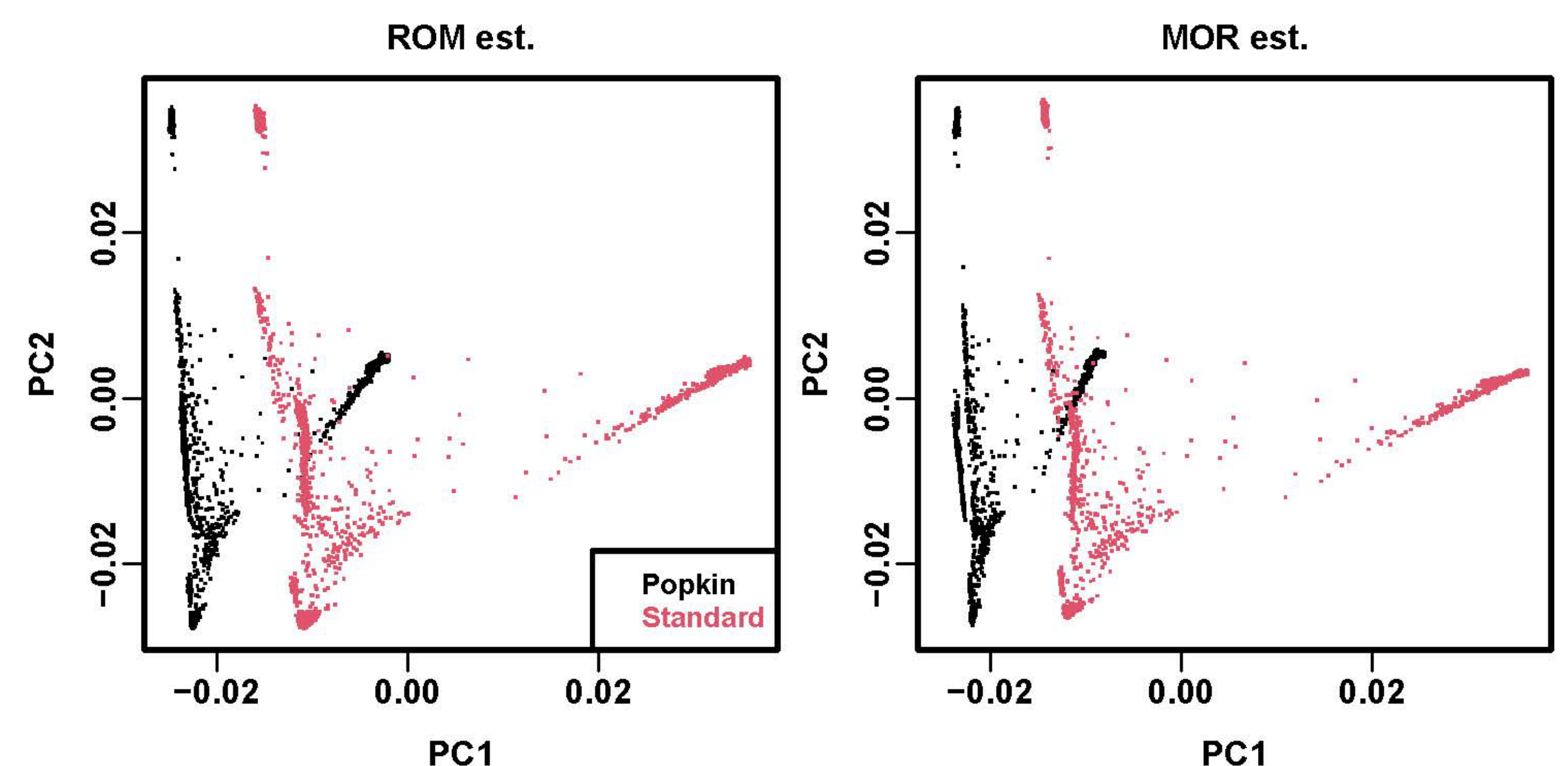


Figure 4: Visualization of PC shift due to kinship biases

## Proof of association invariability to common kinship biases

For standard kinship estimator:

$$\Phi' = \frac{1}{1 - \bar{\varphi}} \mathbf{C}\Phi\mathbf{C}, \quad \text{where } \mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \text{ is center matrix.}$$

In LMM model, kinship bias compensated by intercept:

$$y = \mathbf{1}\alpha' + x_i\beta_i + s' + \epsilon, \quad s' = \mathbf{C}s \sim \text{Normal}(\mathbf{0}, 2\sigma'^2\Phi'),$$

$$\sigma'^2 = (1 - \bar{\varphi})\sigma^2, \quad s' = s - \mathbf{1}\bar{s}, \quad \alpha' = \alpha + \bar{s}, \quad \bar{s} \sim \text{Normal}(0, \sigma^2\bar{\varphi})$$

Similarly, in PCA model:

$$U'_d \approx \mathbf{C}U_d$$

$$y = \mathbf{1}\alpha' + x_i\beta_i + U'_d\gamma'_d + \epsilon,$$

$$\gamma'_d = \gamma_d, \quad U'_d\gamma'_d = U_d\gamma_d - \mathbf{1}\bar{U}_d\gamma_d, \quad \alpha' = \alpha + \bar{U}_d\gamma_d$$

## Reference

- Ochoa, Alejandro and John D. Storey (2021). "Estimating FST and kinship for arbitrary population structures". PLoS Genet 17(1), e1009241.
- Bhatia, Gaurav et al. (2013). "Estimating and interpreting FST: the impact of rare variants". Genome Res. 23(9), pp. 1514–1521.
- 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.
- Yao, Yiqi and Alejandro Ochoa (2022). Limitations of principal components in quantitative genetic association models for human studies. bioRxiv, p. 2022.03.25.485885.