# Modelling complex population structure using $F$-statistics and Principal Component Analysis

Benjamin M Peter

October 14, 2021

**Abstract**

Human genetic diversity is shaped by our complex history. Population genetic tools to understand this variation can broadly be classified into data-driven methods such as Principal Component Analysis (PCA), and model-based approaches such as $F$-statistics. Here, I show that these two perspectives are closely related, and I derive explicit connections between the two approaches. I show that $F$-statistics have a simple geometrical interpretation in the context of PCA, and that orthogonal projections are the key concept to establish this link. I illustrate my results on two examples, one of local, and one of global human diversity. In both examples, I find that population structure is sparse, and only a few components contribute to most statistics. Based on these results, I develop novel visualizations that allow for investigating specific hypotheses, checking the assumptions of more sophisticated models. My results extend $F$-statistics to non-discrete populations, moving towards more complete and less biased descriptions of human genetic variation.

# 1   Introduction

As most species, the genetic diversity of human populations is influenced by our history and environment over the last several hundred thousand years (e.g Cavalli-Sforza et al., 1994, **?**, **?**). Population genetic models use observed patterns of variation to investigate and reconstruct the demographic and evolutionary history of our species (Schraiber and Akey, 2015, **?**).

In particular, in isolated populations genetic drift will slowly change allele frequencies. As a result, isolated populations are expected to slowly differentiate (Wahlund, 1928, Cavalli-Sforza and Piazza, 1975). In humans, this may be caused because continental-scale geographic distances limit migration, causing a pattern known as isolation-by-distance (SLATKIN, 1985). However, isolation-by-distance patterns are usually not uniform, but shaped by geography, particularly barriers to migration such as mountain ranges, oceans or deserts (Cavalli-Sforza et al., 1994, **?**, Rosenberg et al., 2005, Bradburd et al., 2013, Peter et al., 2020). In addition, major historical population movements such as the out-of-Africa, Austronesian or Bantu expansions lead to more gradual patterns of genetic diversity over space (Cavalli-Sforza et al., 1994, Ramachandran et al., 2005, Novembre et al., 2008, Stoneking, 2016, Racimo et al., 2020). Local migration between neighboring populations will reduce differentiation, and long-distance migrations (Alves et al., 2016) and secondary contact between diverged populations, such as Neandertals and modern humans (Green et al., 2010) may lead to locally increased diversity (**?**).

The interplay of these demographic processes leads to the complex genetic structure observed in present-day human populations (The 1000 Genomes Project Consortium, 2015, **?**) with both discrete and continuous components (Pritchard et al., 2000, Rosenberg et al., 2002, Serre and Pääbo, 2004, Rosenberg et al., 2005, Bradburd et al., 2018, Reich, 2018, Peter et al., 2020) that we model in order to reconstruct human demographic history. This is challenging because particularly for large and

heterogeneous data sets we cannot expect to devise a single model that captures all processes. A commonly used analysis paradigm is thus to integrate tools based on different sets of assumptions. each emphasizing particular aspects of the data.

A typical analysis starts with data-driven, exploratory methods that summarize data making minimal assumptions (e.g. Schraiber and Akey, 2015). Examples are population trees (Cavalli-Sforza and Edwards, 1967, Felsenstein, 1973, Cavalli-Sforza and Piazza, 1975), Principal Component Analysis (PCA, Cavalli-Sforza et al., 1994, Patterson et al., 2006)) structure-like models (Pritchard et al., 2000, Alexander et al., 2009) or multidimensional scaling (MDS **?**)). However, these methods are not designed to answer specific research questions, and are limited in their ability to estimate biologically meaningful parameters. For this purpose, methods based on explicit demographic models are often used that aim to fit a specified or estimated model of divergence, migration and genetic drift to the data (Gutenkunst et al., 2009, Excoffier et al., 2013, Kamm et al., 2015). The drawback of these methods is that, to make inference mathematically feasible, we need to introduce strong modeling assumptions such as that populations are discrete, randomly mating, or at equilibrium. While in most cases these assumptions are violated to some extent and cannot be verified, but we hope that the resulting model fits provide a sufficiently accurate answer for specific research questions.

**F-stats**    However, when the number of populations exceeds a few dozen, even codifying reasonable population models can be prohibitively difficult. One approach is to pick a small set of "representative" samples, and restrict modeling to this subset (e.g. Gravel et al., 2011, Harney et al., 2021). However, this has the drawback that a large proportion of the data may be unused. An increasingly popular alternative approach, particularly in the analysis of human ancient DNA, is therefore to build up complex models from smaller building blocks based on the relationship between two, three or four populations.

The framework is based on a set of parameters called $F$-statistics *sensu* Patterson (Reich et al., 2009, Patterson et al., 2012, Peter, 2016). I save the formal definition for later; but the easiest way to motivate them assumes that populations are related as a tree, where the edge lengths measure how much genetic drift has occurred. (Figure 2; Semple and Steel, 2003, Peter, 2016).

In most applications, these $F$-statistics are estimated from data, and then used as tests of tree-ness. In particular, under the assumption of a tree, $F_3$ is restricted to be non-negative, and many $F_4$-statistics will be zero (Semple and Steel, 2003, Patterson et al., 2012), and data that violates these constraints is incompatible with a tree-like relationship between populations. The canonical alternative model is an admixture graph (or phylogenetic network) (Patterson et al., 2012, Huson et al., 2010), which is a tree which allows for additional edges reflecting gene flow (Figure 3A). However, admixture graphs are not the only plausible alternative model, and expected $F$-statistics can be calculated for a wide range of population genetic demographic models (Peter, 2016).

**F-stats and PCA**    The practical issue addressed in this study is how $F$-statistics can be reconciled with one of the most widely used data-driven techniques, PCA. One way PCA can be motivated is as generating a low-dimensional representation of the data, with each dimension (called principal components, PCs) retaining a maximum of the variance present in the data. In population genetics, the use of PCA has been pioneered by Cavalli-Sforza et al. (1964), who used allele-frequency data at a population level to visualize global genetic diversity (Cavalli-Sforza et al., 1994). Currently, PCA is most commonly performed on individual-level genotype data (e.g. Patterson et al., 2006, Novembre et al., 2008), making use of the hundreds of thousands of loci available in most genome-scale data sets. The PCA-decomposition has been studied for a number of population genetic models including trees (Cavalli-Sforza and Piazza, 1975), spatially continuous structure (Novembre and Stephens, 2008), the coalescent (McVean, 2009) and discrete population models (**?**). Here, in order to link PCA to $F$-statistics, I interpret both of them geometrically in *allele frequency space*, i.e. as functions of a high-dimensional Euclidean space. For $F$-statistics, this interpretation was recently

developed by Oteo-Garcia and Oteo (2021), and for PCA it follows naturally from the interpretation of approximating a high-dimensional space with a low-dimensional one.

In the next section, I will formally derive the connection between $F$-statistics and PCA, and show how $F$-statistics can be interpreted geometrically, with a particular emphasis on two-dimensional PCA plots. In the results section, I will then discuss how some of the most common applications of $F$-statistics manifest themselves on a PCA, and illustrate them on two example data sets.

# 2 Theory

In this section, I will introduce the mathematics and notations for $F$-statistics and PCA. A comprehensive treatise on PCA is given by e.g. Jolliffe (2013) a useful primer on the mathematics is Pachter (2014), and a useful guide to interpretation is Cavalli-Sforza et al. (1994). Readers unfamiliar with $F$-statistics may find Patterson et al. (2012), Peter (2016) or Oteo-Garcia and Oteo (2021) helpful.

## 2.1 Formal Definition of $F$-statistics

Let us assume we have a set of populations for which we have SNP allele frequency data from $S$ loci. Let $x_{il}$ denote the frequency at the $l$-th SNP in the $i$-th population; and let $X_i = (x_{i1}, x_{i2}, \ldots x_{iS})$ be a vector collecting all allele frequencies for population $i$. As $X_i$ will be the only data summary considered here for population $i$, I make no distinction between the population and the allele frequency vector used to represent it.

The three $F$-statistics can then be defined as

$$F_2(X_1, X_2) = \frac{1}{S} \sum_{l=1}^{S} (x_{1l} - x_{2l})^2 \tag{1a}$$

$$F_3(X_1; X_2, X_3) = \frac{1}{S} \sum_{l=1}^{S} (x_{1l} - x_{2l})(x_{1l} - x_{3l}) \tag{1b}$$

$$F_4(X_1, X_2; X_3, X_4) = \frac{1}{S} \sum_{l=1}^{S} (x_{1l} - x_{2l})(x_{3l} - x_{4l}), \tag{1c}$$

The normalization by the number of SNPs $S$ is assumed to be the same for all calculations and is thus omitted subsequently. Both $F_3$ and $F_4$ can be written as sums of $F_2$-statistics:

$$2F_3(X_1; X_2, X_3) = F_2(X_1, X_2) + F_2(X_1, X_3) - F_2(X_2, X_3) \tag{2a}$$
$$2F_4(X_1, X_2; X_3, X_4) = F_2(X_1, X_3) + F_2(X_2, X_4) - F_2(X_1, X_4) - F_2(X_2, X_3) \tag{2b}$$

As highlighted in the introduction, $F$-statistics have been primarily motivated in the context of trees and admixture graphs (Patterson et al., 2012). In a tree, the squared Euclidean distance $F_2(X_1, X_2)$ measures the length of the path between populations $X_1$ and $X_2$ (Figure 2A); $F_3$ represents the length of an external branch (Figure 2B) and $F_4$ the length of an internal branch between four nodes, respectively (Figure 2C). Crucially, for branches that do not exist in the tree (as in Figure 2D), $F_4$ will be zero. The length of each branch can be thought of in units of genetic drift, and is non-negative (Patterson et al., 2012).

Thinking of $F$-statistics as branch lengths is useful to understand a number of applications: In particular, one common task is to find the population most closely related to an unknown sample $X_U$ (Raghavan et al., 2014). One way to do that is using an *outgroup-$F_3$-statistic* $F_3(X_O; X_U, X_i)$, where $X_O$ denotes an outgroup, and the $X_i$ are a panel of populations that are candidates for the closest match. The highest values of $F_3$ indicate the population $X_i$ most closely related to $U$, using the outgroup $O$ to correct for differences in sample times. The intuition is given in Figure 1A; where

the outgroup-$F_3$-statistic $F_3(X_O; X_U, X_3)$ is highlighted. It represents the length of the branch from $X_O$ to the common node between the three samples in the statistic, and the closer this node is to $X_U$, the longer the branch and hence the larger the statistic. In contrast to a simple genetic distance, the sample time has no effect: The branch between $X_U$ and $X_2$, would be shorter than to one between $X_U$ and $X_1$, but the path to the shared junction and hence the $F_3$-statistic would be the same. Larger sets of $F_3$ and $F_4$-statistics are also frequently used for complex models, such as reconstructing admixture graphs (Patterson et al., 2012, Lipson et al., 2013) and estimating admixture proportions (Petr et al., 2019, Harney et al., 2021).

Most commonly however, $F_3$ and $F_4$ are used as tests of treeness (Patterson et al., 2012): Negative $F_3$-values correspond to a branch with negative genetic drift, which is a violation of treeness. Similarly if four populations are related as a tree, then at least one of the $F_4$ statistics between the populations will be zero (Patterson et al., 2012). The most widely considered alternative model is an admixture graph (Patterson et al., 2012), an example is given in Figure 3A, where (the typically unobserved) population $X_y$ is generated by a mixture of individuals from the ancestors of $X_2$ and $X_3$. Over time, genetic drift will change $X_y$ to $X_x$, which is the population we observe. This will result in $F_4$-statistics that are non-zero, and, in some cases, in negative $F_3$-statistics (exact conditions can be found in Peter, 2016).

### 2.1.1 Geometric interpretation of $F$-statistics

An implicit assumption in the development of $F$-statistics is that population lineages are mostly discrete, and that gene flow is rare. Recently, Oteo-Garcia and Oteo (2021) showed that these assumptions are not necessary by re-deriving $F$-statistics in a geometric framework. Specifically, they interpret the populations $X_i$ as points or vectors in the $S$-dimensional *allele frequency space* $\mathbb{R}^S$. In this case, the $F$-statistics can be thought of as inner (or dot) products, and that all properties and tests related to treeness can be derived from this larger space. In particular the $F$-statistics can be written as

$$F_2(X_1, X_2) \quad = \frac{1}{S}\sum_{l=1}^{S}(x_{1l} - x_{2l})^2 \qquad = \frac{1}{S}\langle X_1 - X_2, X_1 - X_2 \rangle = \frac{1}{S}\|X_1 - X_2\|^2 \quad (3a)$$

$$F_3(X_1; X_2, X_3) \quad = \frac{1}{S}\sum_{l=1}^{S}(x_{1l} - x_{2l})(x_{1l} - x_{3l}) \quad = \frac{1}{S}\langle X_1 - X_2, X_1 - X_3 \rangle \qquad (3b)$$

$$F_4(X_1, X_2; X_3, X_4) \quad = \frac{1}{S}\sum_{l=1}^{S}(x_{1l} - x_{2l})(x_{3l} - x_{4l}) \quad = \frac{1}{S}\langle X_1 - X_2, X_3 - X_4 \rangle, \qquad (3c)$$

where $\|\cdot\|$ denotes the Euclidean norm and $\langle\cdot,\cdot\rangle$ denotes the dot product. Some elementary properties of the dot product between vectors $a, b, c$ that I will use later are

$$\langle a, b \rangle = \sum_i a_i b_i \qquad (4a)$$

$$\langle a, b \rangle = \|a\|\,\|b\|\cos(\phi) \qquad (4b)$$

$$\langle a, a \rangle = \|a\|^2 \qquad (4c)$$

$$\langle a + c, b \rangle = \langle a, b \rangle + \langle b, c \rangle, \qquad (4d)$$

where $\phi$ is the angle between $a$ and $b$. The inner product is closely related to vector projections

$$proj_b a = \frac{\langle a, b \rangle}{\|b\|^2} b, \qquad (5)$$

which is a vector colinear to $b$ whose length measures how much vector $a$ points in the direction of $b$.

139 The drawback of the geometric approach of Oteo-Garcia and Oteo (2021) is that we have to deal
140 with an very high-dimensional space, as the number of SNPs is frequently in the millions. However,
141 it has been commonly observed that population structure is quite low-dimensional, and that the first
142 few PCs provide a good approximation of the covariance structure in the data (Patterson et al., 2006).
143 Therefore, we may hope that PCA could yield a reasonable approximation of the allele frequency
144 space, and that $F$-statistics as measures of population structure may likewise be well-approximated
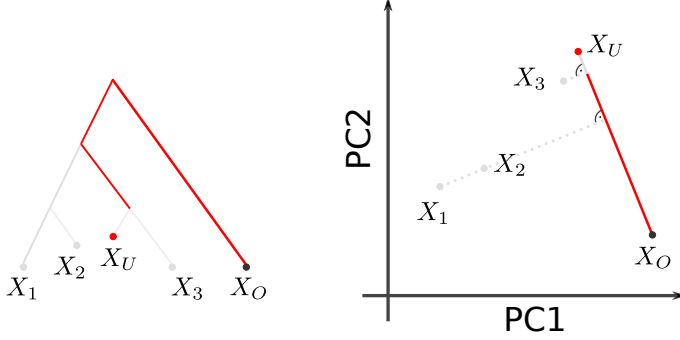145 by the first few PCs.



Figure 1: **Outgroup-$F_3$-statistics**

## 2.2 Formal Definition of PCA

147 PCA is a common way of summarizing genetic data, and so a large number of variations of PCA
148 exist, e.g. in how SNPs are standardized, how missing data is treated or whether we use individuals
149 or populations as units of analysis. The version of PCA I describe here is set up in a way that
150 the similarities to the $F$-statistics framework are maximized, and does *not* reflect how PCA is most
151 commonly applied to genome-scale human genetic variation data sets. In particular, I assume that a
152 PCA is performed on unscaled, estimated population allele frequencies, whereas many applications of
153 PCA are based on individual-level sample allele frequency, scaled by the estimated standard deviation
154 of each SNP (Patterson et al., 2006). The differences this causes will be addressed in the discussion.
155 Let us again assume we have allele frequency data as above, but let us now assume we aggregate
156 the allele frequency vectors $X_i$ in a matrix $\mathbf{X}$ whose entry $x_{il}$ reflects the allele frequency of the $i$-th
157 population at the $l$-th genotype. If we have $S$ SNPs and $n$ populations, $\mathbf{X}$ will have dimension $n \times S$.
158 Since the allele frequencies are between zero and one, we can interpret each Population $X_i$ of $\mathbf{X}$ as
159 a point in $[0,1]^S$, the allele frequency or *data space*, which is a subset of $\mathbb{R}^S$.
160 One way PCA can be motivated is that it aims to find a $K$-dimensional subspace of the data
161 space that retains most variation in the data. $K$ is at most $n - 1$, in which case the data is simply
162 rotated. However, the historical processes that generated genetic variation often result in *low-rank*
163 data (Engelhardt and Stephens, 2010), so that $K \ll n$ explains a substantial portion of the variation;
164 for visualization $K = 2$ is frequently used.

There are several algorithms that are used to perform PCAs, the most common one is based on
singular value decomposition (Jolliffe, 2013). In this approach, we first mean-center $\mathbf{X}$, obtaining a
centered matrix $\mathbf{Y}$

$$y_{il} = x_{il} - \mu_l$$

165 where $\mu_l$ is the mean allele frequency at the $l$-th locus.
166 PCA can then be written as

$$\mathbf{Y} = \mathbf{CX} = (\mathbf{U\Sigma})\mathbf{V}^T = \mathbf{PL}, \tag{6}$$

167 where $\mathbf{C} = \mathbf{I} - \frac{1}{n}\mathbf{1}$ is a centering matrix that subtracts row means, with $\mathbf{I}, \mathbf{1}$ the identity matrix
168 and a matrix of ones, respectively. For any matrix $\mathbf{Y}$, we can perform a singular value decomposition

5

$\mathbf{Y} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$ which, in the context of PCA, is interpreted as follows: The matrix of principal components $\mathbf{P} = \mathbf{U}\boldsymbol{\Sigma}$ has size $n \times n$ and contains information about population structure. The SNP loadings $\mathbf{L} = \mathbf{V}^T$ form an orthonormal basis of size $n \times S$, its rows give the contribution of each SNP to each PC. It is often used to look for outliers, which might be indicative of selection (e.g **?**). Alternatively, the PCs can also be obtained from an eigendecomposition of the covariance matrix $\mathbf{Y}\mathbf{Y}^T$. This can be motivated from (6):

$$\mathbf{Y}\mathbf{Y}^T = \mathbf{P}\mathbf{L}\mathbf{L}^T\mathbf{P}^T = \mathbf{P}\mathbf{P}^T, \tag{7}$$

since $\mathbf{L}\mathbf{L}^T = \mathbf{I}$.

## 2.3   Connection between PCA and $F$-statsitics

### 2.3.1   Principal components from $F$-statistics

PCA, as defined above, and $F$-statistics are closely related. In fact, the principal components can be directly calculated from $F$-statistics using multidimensional scaling, which, for squared Euclidean ($F_2$)-distances, leads to an identical decomposition to PCA (Gower, 1966). Suppose we calculate the pairwise $F_2(X_i, X_j)$ between all $n$ populations, and collect them in a matrix $\mathbf{F}_2$. We can obtain the principal components from this matrix by double-centering it, so that its row and column means are zero, and perform an eigendecomposition of the resulting matrix:

$$\mathbf{P}\mathbf{P}^T = -\frac{1}{2}\mathbf{C}\mathbf{F}_2\mathbf{C}. \tag{8}$$

### 2.3.2   $F$-statistics in PCA-space

By performing a PCA, we rotate our data to reveal the axes of highest variation. However, the dot product is invariant under rotation, and $F$-statistics can be thought of as dot products (Oteo-Garcia and Oteo, 2021). What this means is that we are free to calculate $F_2$ either on the uncentered data $\mathbf{X}$, the centered data $\mathbf{Y}$ or any other orthogonal basis such as the principal components $\mathbf{P}$. Formally,

$$
\begin{aligned}
F_2(X_i, X_j) &= \sum_{l=1}^{L} \left(x_{il} - x_{jl}\right)^2 \\
&= \sum_{l=1}^{L} \left((x_{il} - \mu_l) - (x_{jl} - \mu_l)\right)^2 = F_2(Y_i, Y_j) \\
&= \sum_{k=1}^{n} (p_{ik} - p_{jk})^2 = F_2(P_i, P_j), \tag{9}
\end{aligned}
$$

A derivation of this change-of-basis is given in Appendix A, Equation A1. As $F_3$ and $F_4$ can be written as sums of $F_2$-terms (Eqs. 2a, 2b), analogous relations apply.

In most applications, we do not use all PCs, but instead truncate to the first $K$ PCs, which explain most of the between-population genetic variation. Thus,

$$
\begin{aligned}
F_2(P_i, P_j) &= \sum_{k=1}^{K} (p_{ik} - p_{jk})^2 + \sum_{k=K+1}^{n} (p_{ik} - p_{jk})^2 \\
&= \hat{F}_2^{(K)}(P_i, P_j) + \epsilon^{(K)}(P_i, P_j) \qquad . \tag{10}
\end{aligned}
$$

In this notation, $\hat{F}_2^{(K)}$ is the approximation of $F_2$ with only the first $K$ PCs considered, and $\epsilon^{(K)}$ is the corresponding approximation error. I will omit the superscript of $\hat{F}_2$ when the exact number of

PCs is not relevant. If we sum up the squared approximation errors over all pairs of populations, we obtain

$$\sum_{i,j} \epsilon^{(K)}(P_i, P_j)^2 = \sum_{i,j} \left( \hat{F}_2^{(K)}(P_i, P_j) - F_2^{(K)}(P_i, P_j) \right)^2 = \left\| \mathbf{F}_2 - \hat{\mathbf{F}}_2 \right\|_F^2, \tag{11}$$

where the Frobenius-norm $\|\cdot\|_F^2$ of a matrix is defined as the square root of the sum-of-squares of all its elements. This is precisely the function that is minimized in MDS (Jolliffe, 2013). In that sense, $\hat{\mathbf{F}}_2^{(K)}$ is the optimal low-rank approximation of $\mathbf{F}_2$ for any $K$ in that it minimizes the sum of approximation errors of all $F_2$-statistics.

### 2.3.3 $F$-statistics and projection on PCA

One of the easiest ways of dealing with missing data in PCA is to calculate the principal components (equation 6) only on a subset of the data with no missingness, and then to *project* the lower quality samples with high missingness onto this PCA. The simplest way to do this is to note that

$$\mathbf{Y}\mathbf{L}^T = \mathbf{P}\mathbf{L}\mathbf{L}^T = \mathbf{P},$$

and so a new (centered) population $Y_{\text{new}}$ can be projected onto an existing PCA simply by post-multiplying it with $\mathbf{L}^T$:

$$P_{\text{proj}} = Y_{\text{new}}\mathbf{L}^T;$$

the $k$-th entry of $P_{\text{proj}}$ gives the coordinates of the new sample on the $k$-th PC. However, it is likely that $Y_{\text{new}}$ lies outside the variation of the original sample. In this case, there is a projection error

$$\|Y_{\text{new}} - P_{\text{proj}}\mathbf{L}\|^2 = F_2(P_{\text{proj}}\mathbf{L}, Y_{\text{new}}).$$

If we project with missing data, a similar projection can be used where we remove the rows from $Y_{\text{new}}$ and $\mathbf{L}$ where data in $Y_{\text{new}}$ is missing (Patterson et al., 2006).

Thus, if we compare the $F$-statistic of a projected sample, we have

$$\begin{aligned} F_2(X_i, X_{\text{new}}) &= F_2(Y_i, Y_{\text{new}}) \\ &= F_2(P_i, P_{\text{proj}}) + F_2(P_{\text{proj}}\mathbf{L}, Y_{\text{new}}) \\ &= \hat{F}_2(P_i, P_j) + \epsilon(P_i, P_j) + F_2(P_{\text{proj}}\mathbf{L}, Y_{\text{new}}). \end{aligned} \tag{12}$$

The second row follows because the projection error and projection are orthogonal to each other. The main implication of equation 12 is that we do not need to worry about distinguishing the projection and approximation errors.

## 3 Material & Methods

The theory outlined in the previous section suggests that $F$-statistics have a geometric interpretation in PCA-space, which can be approximated on PCA plots. In the next section I explore this connection in detail, and illustrate it on two sample data sets that I briefly introduce here. Both are based on the analyses by Lazaridis et al. (2014). The data is from the Reich lab compendium data set (v44.3), downloaded from `https://reich.hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-` using data on the "Human Origins"-SNP set (597,573 SNPs). SNPs with missing data in any population are excluded. The code used to create all figures and analyses will be available on `https://github.com/BenjaminPeter/fstats_pca`.

<sub>218</sub> **"World" data set** This data set is a subset of the "World Foci" data set of Lazaridis et al. (2014),
<sub>219</sub> where I removed samples which are not permitted for free reuse. These populations span the globe and
<sub>220</sub> roughly represents global human genetic variation (638 individuals from 33 population) As adjacent
<sub>221</sub> sampling locations are often thousands of kilometers apart, I speculate that gene flow between these
<sub>222</sub> populations may not be particularly common; and their structure may therefore be well-approximated
<sub>223</sub> by an admixture graph. A file with all individuals used and their assigned population is given in
<sub>224</sub> **Supplementary File 1.**

<sub>225</sub> **Western Eurasian data set** This data set of 1,119 individuals from 62 populations contains
<sub>226</sub> present-day individuals from the Eastern Mediterranean, Caucasus and Europe. It is motivated
<sub>227</sub> by the analysis of Lazaridis et al. (2014), who used it as a basis of comparison for ancient genetic
<sub>228</sub> analyses of Western Eurasian individuals, and PCAs based on similar sets of samples have been used
<sub>229</sub> in many other ancient DNA studies (e.g. ?Haak et al., 2015). Genetic differentiation in this region is
<sub>230</sub> low and closely mirrors geography (Novembre et al., 2008). I thus speculate that gene flow between
<sub>231</sub> these populations is common (Ralph and Coop, 2013), and a discrete model such as a tree or an
<sub>232</sub> admixture graph might be a rather poor reflection of this data. A file with all individuals used and
<sub>233</sub> their assigned population is given in **Supplementary File 2.**

<sub>234</sub> **Computing $F$-statistics and PCA** I perform analyses at the level of populations to ease pre-
<sub>235</sub> sentation. It is an assumption of $F$-statistics that the genetic variation within sampled population
<sub>236</sub> is independent of the variation between samples (Patterson et al., 2012). All computations are per-
<sub>237</sub> formed in R. I use `admixtools 2.0.0` (`https://github.com/uqrmaie1/admixtools`) to compute
<sub>238</sub> $F$-statistics. To obtain a PC-decomposition, I first calculate all pairwise $F_2$-statistics, and then use
<sub>239</sub> equation 8 and the `eigen` function to obtain the PCs. A squared Euclidean distance matrix such as
<sub>240</sub> $\mathbf{F}_2$ has all negative or zero eigenvalues (i.e. $\mathbf{F}_2$ is negative-semidefinite). However as $F_2$-statistics are
<sub>241</sub> estimates, some eigenvalues might be slightly positive, which would lead to imaginary PCs. I avoid
<sub>242</sub> this by using the `nearPD`-function in R that ensures all eigenvalues have the correct sign.

# 4  Results

<sub>244</sub> The transformation from the previous section allows us to consider the geometry of $F$-statistics in
<sub>245</sub> PCA-space. The relationships we will discuss formally only hold if we use all PCs. However, the
<sub>246</sub> appeal of PCA is that frequently, only a very small number $K \ll n$ of PCS contain most information
<sub>247</sub> that is relevant for population structure (for visualization $K = 2$ is often used).

## 4.1  $F_2$ in PC-space

<sub>249</sub> The $F_2$-statistic is an estimate of the squared allele-frequency distance between two populations. On
<sub>250</sub> a tree (Figure 2A) this corresponds to the branch between two populations. In allele-frequency space,
<sub>251</sub> it corresponds to the squared Euclidean distance, and thus reflects the intuition that closely related
<sub>252</sub> populations will be close to each other on a PCA-plot, and have low pairwise $F_2$-statistics. However,
<sub>253</sub> since the right-hand side of equation 9 is a sum of squared (non-negative) terms, the $F_2$-distance on
<sub>254</sub> a PCA-plot will be an underestimate of the full distance. Thus, if we find two populations with high
<sub>255</sub> $F_2$-distance nearby on a PCA-plot this suggests that substantial variation is hidden on other PCs,
<sub>256</sub> and that these particular PCs are not suitable to understand and visualize the relationship between
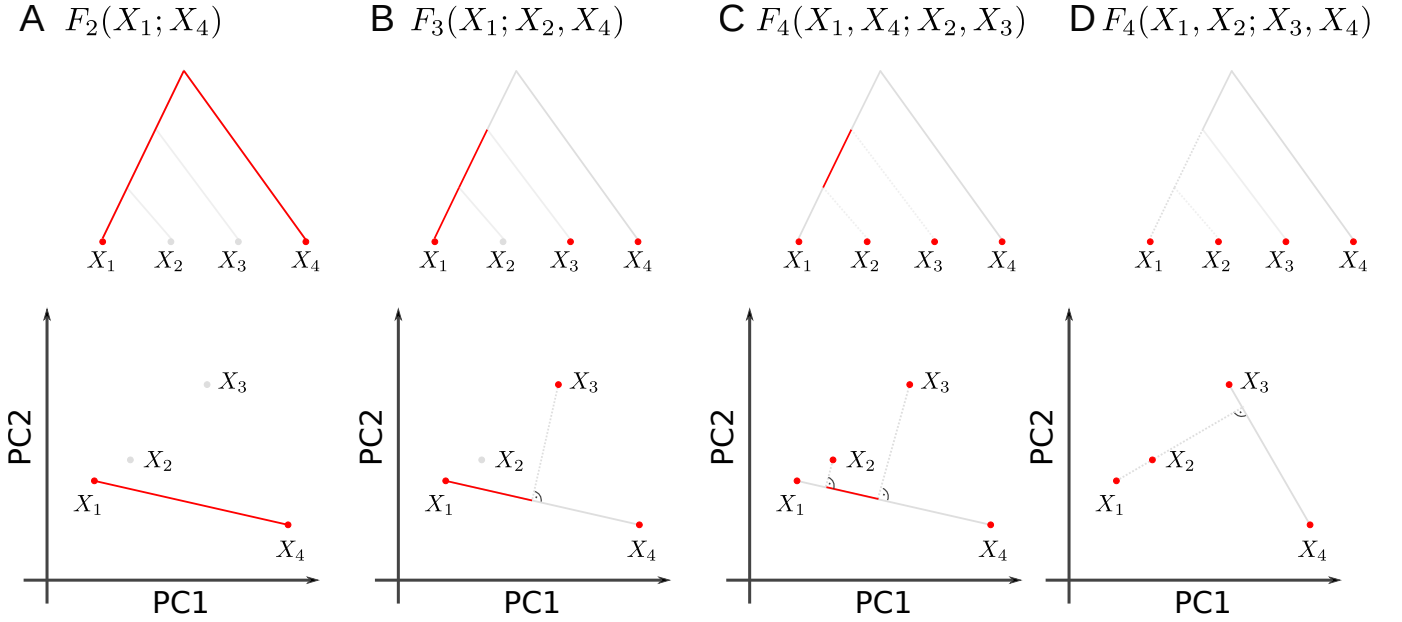<sub>257</sub> these particular populations.

Figure 2: **Representation of $F$-statistics on tree and 2D-PCA-plot.** The schematics show four populations and their representation using a tree (top row) or a 2D-PCA plot. A: $F_2$ represents the (squared) Euclidean distance between two tree leafs, and in PC-space. B: $F_3(X_1; X_3, X_4)$ corresponds to the external branch from $X_1$ to the internal node joining the populations, and is proportional to the orthoginal projection of $X_1 - X_3$ onto $X_1$-$X_4$. C: $F_4(X_1, X_4; X_2, X_3)$ corresponds to the internal branch in the tree, or the orthogonal projection of $X_2 - X_3$ on $X_1 - X4$. D: $F_4(X_1, X_2; X_3, X_4)$ The two paths from $X_1$ to $X_2$ and $X_3$ and $X_4$ are non-overlapping in the tree, which corresponds to orthogonal vectors in PCA-space.



Figure 3: **Admixture representation on 2D-PCA-plot.** The schematics show four populations and their representation using an admixture graph (A) or a 2D-PCA plot. A: Admixture graph, with population $X_y$ originating as an admixture of $X_2$ and $X_3$, with $X_2$ contributing proportion $\alpha$. Subsequent drift (red branch) will change allele frequency to sampled admixture population $X_x$. B: PCA representation of the scenario in A. $X_y$ originates on the segment connecting $X_2$ and $X_3$, and subsequent drift may move it in a random direction. C: $F_3(X_x; X_2, X_3)$ and negative region (light gray circle). $F_4(X_1, X_x; X_3, X_4)$ will no longer be zero (compare to Figure 2D).

## 4.2 When are admixture-$F_3$ statistics negative?

Consider the admixture scenario in Figure 3A, where population $X_y$ is the result of a mixture of $X_2$ and $X_3$, and subsequent drift changes the allele frequencies of the admixed population from $X_y$ to the sampled population $X_x$. How is such a scenario displayed on a PCA? Since the allele frequencies of $X_y$ are a linear combination of $X_2$ and $X_3$, it will lie on the line segment connecting these two populations (Figure 3B), at a location proportional to the admixture proportions. Subsequent drift

will change the allele frequency of $X_x$, and so in general it might fall on a different point on a PCA-plot. An exception occurs when $X_x$ (and no other populations related to $X_x$) are not part of the construction of the PCA, so that $F_2(X_x, X_y)$ is orthogonal to all PCs. In this case, $X_x$ and $X_y$ project to the same point, and the location on the PCA can be used to predict the admixture proportions (Brisbin et al., 2012, McVean, 2009, Oteo-Garcia and Oteo, 2021).

However, if $X_x$, is included in the construction of the PCA, or if the population structure is sufficiently complex that gene flow occurred between $X_x$ and any of the populations used to construct the PCA, $X_x$ and $X_y$ will project on different spots (Figure 3B). Thus, a natural question to ask is given two source populations $X_2$, $X_3$, can we use PCA to predict which populations might be considered admixed between them?

One way to address this question is to consider the space for which $F_3$ is negative, i.e.

$$
\begin{aligned}
2F_3(X_x; X_2, X_3) &= 2\langle X_x - X_2, X_x - X_3 \rangle \\
&= \|X_x - X_2\|^2 + \|X_x - X_3\|^2 - \|X_2 - X_3\|^2 < 0.
\end{aligned} \tag{13}
$$

By the Pythagorean theorem, $F_3 = 0$ if and only if $X_2, X_3$ and $X_x$ form a right-angled triangle. The associated region where $F_3 = 0$ is a $n$-sphere (or a circle in two dimensions) with diameter $\overline{X_2 X_3}$ (The overline denotes a line segment). $F_3$ is negative when the triangle is obtuse, i.e. $X_x$ could be considered admixed if it lies inside the $n$-ball with diameter $\overline{X_2 X_3}$ (Figure 2B, Equation A2).

$F_3$ **on a 2D PCA-plot.** If we project this $n$-ball on a two-dimensional plot, $\overline{X_2 X_3}$ will usually not align with the PCs; thus the ball may be somewhat larger than it appears on the plot. This geometry is perhaps easiest visualized on a globe. If we look at the globe from a view point parallel to the equator, both the north and south poles are visible at the very edge of the circle. But if we look at it from above the north pole, the north- and south-poles will be at the very same point.

Thus if $\hat{F}_3 \ll F_3$, the true circle will be bigger than would be predicted from a 2D-plot. In this case, substantial relevant genetic differentiation is "hidden" in the higher PCs, and populations that appear inside the circle on a PCA-plot may, in fact, have positive $F_3$-statistics. This is because they are outside the $n$-ball in higher dimensions. The converse interpretation is more strict: if a population lies outside the circle on *any* 2D-projection, $F_3$ is guaranteed to be bigger than 0 (see Equation A4 in the Appendix).

**Example** As an example, I visualize the admixture statistic $F_3(X; \text{Basque}, \text{Turkish})$, on the first two PCs of the Western Eurasian data set (Figure 4A). In this case, the projected $n$-ball (light gray) and circle based on 2D (dark gray) align relatively closely, but several populations inside the ball (e.g. Sardinian, Finnish) have, in fact, positive $F_3$-values. This reveals that the first two PCs do not capture all the genetic variation relevant for Southern European population structure. Consequently, approximating $F_3$ by the first two or ten PCs (Figure 4B) only gives a coarse approximation of $F_3$, and from Figure 4C we see that many higher PCs contribute to $F_3$ statistics.

However, many populations, particularly from Western Asia and the Caucasus, fall outside the circle. This allows us to immediately conclude that their $F_3$-statistics must be positive; and we should not consider them as a mixture between Basques and Turks.

## 4.3 $F_3$-statistics as projections

Outgorup-$F_3$-statistic motivate a comparison of $F_3$-statistics to projections (Equation 5), consider again the case displayed in Figure 1A, where the goal is to find the population $X_i$ that is closest to an unknown population $X_U$, with respect to an outgroup $X_O$, which we can do using the statistic $F_3(X_O; X_U, X_i)$. On a PCA-plot, we can visualize this $F_3$-statistic as the projection of the vector
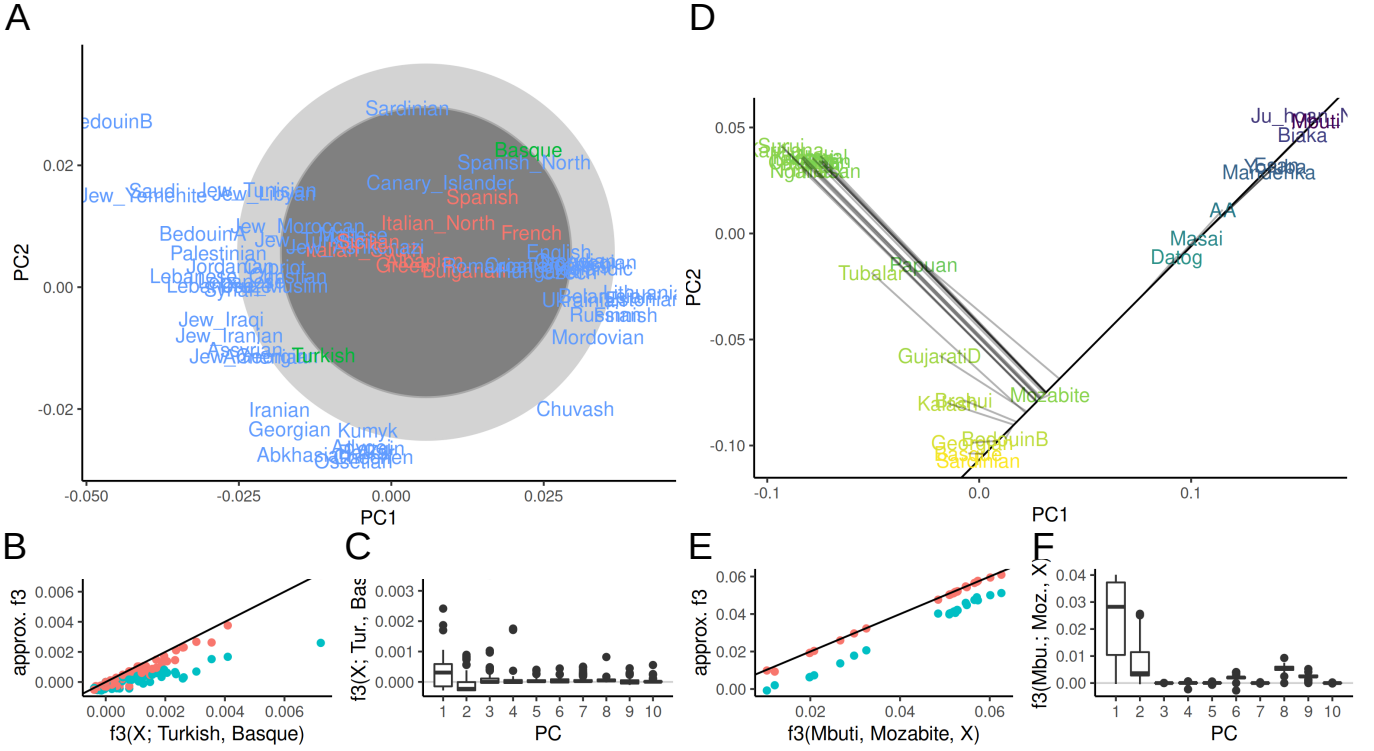
Figure 4: **PCA and $F_3$-statistics** A: PCA of Western Eurasian data; the circle denotes the region for which $F_3(X; \text{Basque}, \text{Turkish})$ may be negative. Populations for which $F_3$ is negative are colored in red. B, E: $F_3$ approximated with two (blue) and ten (red) PCs versus the full spectrum. C,F: Contributions of PCs 1-10 to each $F_3$-statistic. D: PCA of World data set, color indicates value of $F_3(\text{Mbuti}; \text{Mozabite}, X)$. The black line shows the projection axis Mbuti-Mozabite, the gray lines indicates the projected position of each population.

$X_i - X_O$ onto $X_U - X_O$:

$$proj_{X_U - X_O} X_i - X_O = \frac{F_3(X_O; X_U, X_i)}{F_2(X_O; X_U)}(X_U - X_O).$$

Of the right-hand-side terms, $X_U - X_O$ gives the direction of the resulting vector, and the $F_2$-term in the denominator is a normalizing constant. Neiter of these terms depend on $X_i$, which justifies the argument that the $F_3$-statistic and length of the projected vector are proportional to each other, and can thus be interpreted similarly. Thus, the outgroup-$F_3$-statistic is larger for whichever $X_i$ projects furthest along the axis from the outgroup to the unknown population; in the example in Figure 1B this is $X_3$.

**Example**  Again, these projections will be orthogonal when using the full data, and may only be approximately orthogonal when approximated using the first few PCs. In Figure 4D, I visualize the outgroup-$F_3$-statistic $F_3(\text{Mbuti}; \text{Mozabite}, X_i)$, i.e. a statistic that aims to find the population most closely related to Mozabite (a Berber ethnic group from the northern Sahara), assuming the Mbuti are an outgroup. On a PCA, we can interpret this $F_3$ statistic as the projection of the line segment from Mbuti to population $X_i$ onto the line through Mbuti and Mozabite (black line). For each population, the projection is indicated with a grey line. In the full data space, this line is always orthogonal to the segment Mbuti-Mozabite, but on the plot (i.e. the subspace spanned by the first two PCs), this is not necessarily the case. The coloring is based on the $F_3$-statistic calculated from all the data, with brighter values indicating higher $F_3$-statistics. In this case, the first two PCs approximate the $F_3$-statistic very well: Particularly the samples from East Asia, Siberia and the

11

Americas project very close to orthogonally, suggesting that most of the genetic variation relevant for this analysis is captured by these first two PCs. We can quantify this and find that the first two PCs slightly underestimate the absolute value of $F_3$ (Figure 4E), but keep the relative ordering. I also find that many PCs, e.g. PCs 3-5, 7 and 10 have almost zero contribution to all $F_3$-statistics (Figure 4F), and PCs 6, 8 and 9 having a similar non-zero contribution for almost all statistics, likely because these PCs explain within-African variation.

## 4.4 $F_4$-statistics as angles

The interpretation of $F_4$ in PCA is similar to that of $F_3$ as a projection of one vector onto another, with the difference that now all four points may be distinct. $F_4$-statistics that correspond to a branch in a tree (as in Figure 2C), can be interpreted as being proportional to the length of a projected segment on a PCA plot (Figure 2G), again with the caveat that we need to scale it by a constant. If the $F_4$-statistic corresponds to a branch that does not exist in the tree, i.e it is a test statistic (Figure 2D), then, from the tree-interpretation, we expect $F_4(X_1, X_2; X_3, X_4) = 0$ implies that the vectors $X_1 - X_2$ and $X_3 - X_4$ are orthogonal to each other, or that the two populations map to the same point (Figure 2H). In the case of an admixture graph, this is no longer the case: Population $X_x$ in Figure 3D does *not* map to the same point as $X_1$ or $X_2$ do, implying that statistics of the form $F_4(X_1, X_x; X_3, X_4) \neq 0$.
$F_4(X_3, X_3'; X_4, X_4') = 0$.

Since $F_4$ is a covariance, its magnitude lacks an interpretation. Therefore, commonly correlation coefficients are used, as there, zero means independence and one means maximum correlation. For $F_4$, we can write

$$\text{Cor}(X_1 - X_2, X_3 - X_4) = \frac{\langle X_1 - X_2, X_3 - X_4 \rangle}{\|X_1 - X_2\| \|X_3 - X_3\|} = \cos(\phi), \tag{14}$$

where $\phi$ is the angle between $X_1 - X_2$ and $X_3 - X_4$. Thus, independent drift events lead to $\cos(\phi) = 0$, so that the angle is 90 degrees, whereas an angle close to zero means $\cos(\phi) \approx 1$, which means most of the genetic drift on this branch is shared.

**Example** To illustrate the angle interpretation I again use the Western Eurasian data. The PCA-biplot shows two roughly parallel clines (Figure 4A), a European gradient (from Sardinian to Chuvash), and a Asian cline (from Arab to Caucasus populations). This is quantified in Figure 5A, where I plot the angle corresponding to $F_4(X, \text{Sardinian}; \text{Saudi}, \text{Georgian})$. For most European populations, using two PCs (green points) gives an angle close to zero, corresponding to a correlation coefficient between the two clines of $r > 0.9$. Just adding PC3 (blue), however, shows that the parallelism of the clines is spurious: Using three PCs or the full data (red) shows that most correlations are low. I arrive at a similar interpretation from the spectrum of these statistics (Figure 5B), which has high loadings for the first three PCs, with minimal contributions from the higher ones.

## 4.5 Other projections

So far, I used eq. 9 to interpret $F$-statistics on a PC-plot, but the argument holds for *any* orthonormal projection of the data space. This is useful in particular for estimates of admixture proportions, which are often done in a small reference space (Patterson et al., 2012, Petr et al., 2019, Harney et al., 2021, Oteo-Garcia and Oteo, 2021).

The simplest approach is the $F_4$-ratio to infer the admixture sources of population $X$ as

$$\alpha = \frac{F_4(R_1, R_2; X, A)}{F_4(R_1, R_2; B, A)} = \frac{proj_{R_1 - R_2} X - A}{proj_{R_1 - R_2} B - A}, \tag{15}$$

which can be interpreted as projecting $X - A$ and $B - A$ onto $R_1 - R_2$ and measuring their relative proportions (Oteo-Garcia and Oteo, 2021). `qpAdm` extends this approach to a higher-dimensional reference space and multiple potential source populations. One open practical question in many applications is which reference and putative source populations to use (Harney et al., 2021). The theory developed here suggests some possible visualizations that may address this issue.

### 4.5.1  Example

In the PCA on the world overview data set, I found a gradient from Africans to Europeans (Figure 4D). I focus on this cline using an alternative projection by using $F$-statistics of the form $F_4(X, Y; \text{Sardinian}, \text{Yoruba}))$, which might e.g. be used in an $F_4$-ratio. These $F_4$-statistics are very well-approximated by the first two PCs, with a 99.2% correlation between $F_4$ and its approximation using the first two PCs (Figure 5C).

In Figure 5D, I show the projection $\langle X; \text{Sardinian}, \text{Yoruba} \rangle$ on the $X$-axis, which means that the horizontal difference between any pair of population is proportional to their $F_4$-statistic relative to Sardinians and Yorubans. We can also ask what variation is not represented by performinc a PCA on the residual of this projection, the first two residual PCs are given on the Y-axis and in the coloring. This visualization reveals that variation within Africans (with Mbuti, Biaka and Ju|'hoansi, top right) and the variation in East Asians and Americans are largely orthogonal to this projection axis, and so Sardinians and Yoruba would be poor references if we were interested in studying East Asian genetic variation.

The percentage of between-population variance explained by the Sardinia-Yoruba axis (24%) is much lower than that of the first PC (40%, Figure 5E). However, the cumulative variance explained by the first two axes is similar, with (52%) explained when adding residual PC1 to the projection, compared to 55% for the first two PCs. The advantage of specifying one axis is that it displays the orthogonal components more explicitly, reveals distinct structure in Africans and non-Africans and thus can be used to test assumptions of more complex models.

## 5  Discussion

Particularly for the analysis of human genetic variation, $F$-statistics are a powerful tool to describe population genetic diversity. Here, I show that the geometry of $F$-statistics (Oteo-Garcia and Oteo, 2021) leads to a number of simple interpretations of $F$-statistics on a PCA plot. This allows for direct and quantitative comparisons between $F$-statistic-based results and PCA biplots. As PCA is often ran in an early step in data analysis, this also aids in generation of hypotheses that can be more directly evaluated using generative models, (e.g using a lower number of populations). It also allows reconciling apparent contradictions between $F$-statistics and PCA-plots; differences between the two data summaries are explained solely by higher PCs, and so whenever such contradictions arise, higher PCs will be informative for population structure. Previous interpretation of PCA in the context of population genetic models have primarily focused on the PCs, which can be derived analytically for trees (Cavalli-Sforza and Piazza, 1975) and homogeneous spatial models (Novembre and Stephens, 2008). My interpretation here is different in that it puts more emphasis on the geometry itself, rather than directly interpreting the PCs. One consequence is that the results here are not impacted by sample ascertainment and sample sizes (McVean, 2009, Novembre and Stephens, 2008), which are common concerns in the interpretation of PCA. However, a very skewed sampling distribution will increases the likelihood that more or different PCs will have to be included in the analysis. From this perspective, one could envision a framework where $F$-statistics are used to decide which samples should be included to obtain a low-dimensional PCA-plot "representative" of the data

As $F$-statistics are motivated by trees, they assume that populations are discrete, related as a graph, and that gene flow between populations is rare (Patterson et al., 2012, Harney et al., 2021).
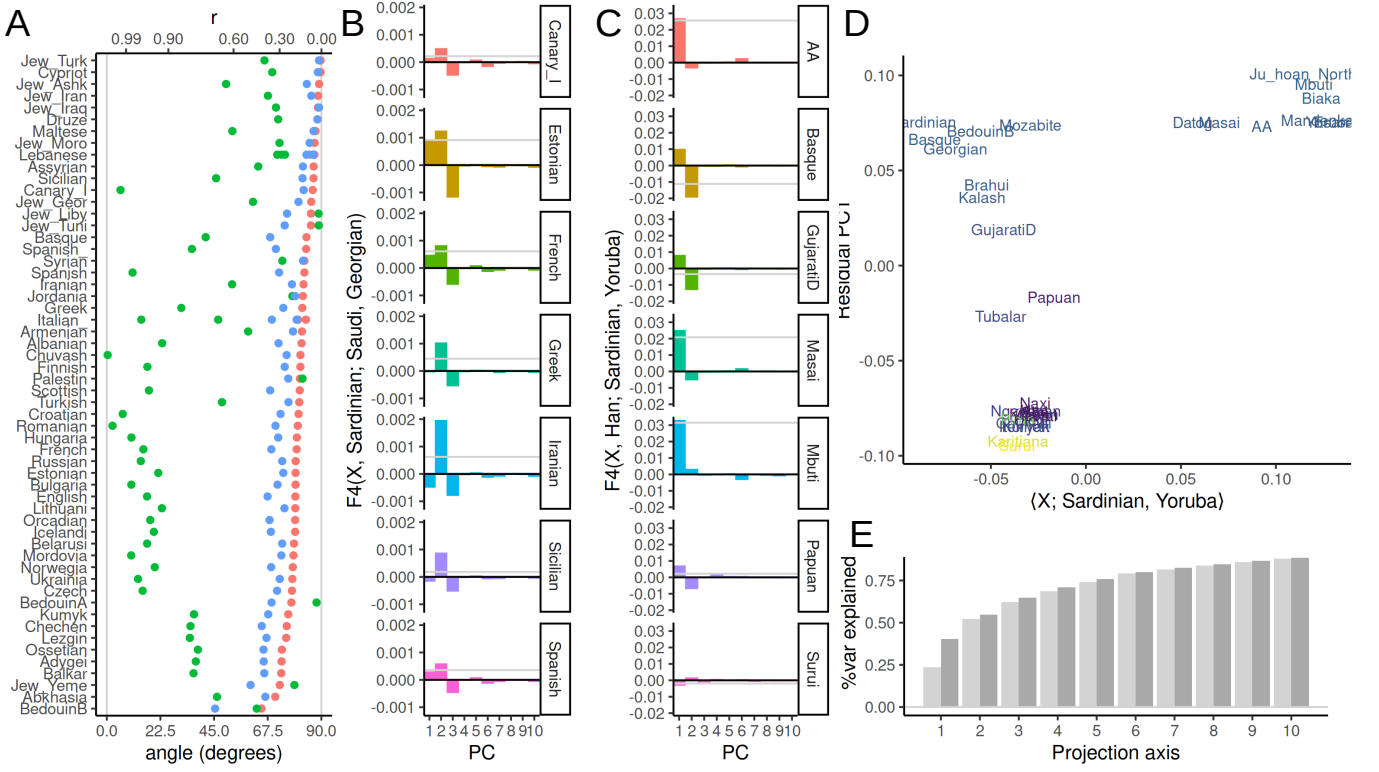
Figure 5: **PCA and $F_4$-statistics** A: Projection angle and correlation coefficient $r$ representation of $F_4(X, Sardinian; Saudi, Georgian)$ (red) in the Western Eurasian data set, and approximations using two (green) and three (blue) PCs. B: Spectrum of select $F_4$-statistics in the Western Eurasian data set. C: Spectrum of $F-4$-statistics in World data set. D: Scatterplot of $F_3$-projection on Sardinian-Yoruba axis and residual PC1. E: Percent variance explained from of the projection based on $F_3$ in panel D and first nine residual PCs (light gray), compared with percent variance explained by first ten PCs (dark gray).

However, in many regions, all humans populations are admixed to some degree (Pickrell and Reich, 2014), and in regions such as Europe, genetic diversity is distributed continuously (Novembre et al., 2008, Novembre and Stephens, 2008). This provides a challenge for interpretation; as many $F_3$ and $F_4$ statistics may indicate gene flow. In my example (Figure 4A), most Southern European populations are "admixed" between Basques and Turkish, but a more accurate model might be one of continuous variation where Basque and Turkish lie on one of multiple gradients; which is more directly visualized with PCA. There are a number of tools that have been developed that use multiple $F$-statistics to build complex models, such as `qpGraph` (Lazaridis et al., 2014) and `qpAdm` (Harney et al., 2021). One issue with these approaches is that they are usually restricted to at most a few dozen populations. As ancient DNA data sets now commonly include thousands of individuals, analysts are faced with the challenge of which data to include. A common approach is to sample a large number of distinct models, and retain the ones that are compatible with the data. However, as both `qpGraph` and `qpAdm` assume that gene flow is rare and discrete, selecting sets of populations that did experience little gene flow will provide good fits. One example of this is the world foci data set used here, which contains only 33 populations from across the world, and which is well-approximated by two PCs. However, this ascertainment misses a large amount of variation; a more dense sampling would show that in many places human genetic diversity is very gradual and multi-layered (Lazaridis et al., 2014, Peter et al., 2020). The PCA-based interpretation offers an alternative that trades interpretability for robustness. Particularly interpreting a (normalized) $F_4$-statistic as a correlation coefficient translates to generalized models of gene flow. Separating $F$-statistics in a sum of model and residuals, and performing a PCA on the latter (such as in Figure 5D) is another way how we

14

can visualize $F$-statistics and evaluate the model fit.

To make this link directly applicable to data analysis, there are a number of – primarily statistical – concerns that will need to be addressed. Fist, PCA is most frequently run on individuals, whereas $F$-statistics are often calculated on populations. This is largely because in most workflows, PCA is run much earlier than $F$-statistics; it is a standard assumption of $F$-statistics that there is no population substructure (Patterson et al., 2012), and an easy way to test that is ensure that all individuals cluster tightly on a PCA.

A second difference is that frequently, rare SNPs are weighted higher in PCA, whereas all SNPs are weighted the same for $F$-statistics (Patterson et al., 2006, 2012). This is a difference of convention (Cavalli-Sforza and Piazza, 1975); since $F$-statistics are summed over SNPs with the same expectation, $F$-statistics could also be calculated using the same weighting. The close connection between the two approaches developed here suggest that for most analyses, users might want to be consistent and use the same weighting for both types of analyses.

The third and perhaps biggest gap are statistical issues. The treatment here focuses on the mean estimated $F$-statistic, but many applications of $F$-statistics are based on hypothesis tests (Patterson et al., 2012). This requires estimating accurate standard errors for these statistics, which is difficult since nearby SNPs will be correlated due to recombination (Hahn, 2018). In contrast, PCA jointly models the covariance matrix due to population structure and sampling, so if hypothesis tests are desired this will need to be incorporated.

An advantage of calculating $F$-statistics based on PCs is that they yield consistent estimtates. For both data sets I investigated here, the matrix $\mathbf{F}_2$ of $F$-statistics obtained using admixtools2 is not a proper squared Euclidean distance matrix, i.e. it is not negative semidefinite and has imaginary PCs. A model-based framework based on probabilistic PCA (Hastie et al., 2015, Meisner et al., 2021, Agrawal et al., 2020) would likely be able to generate consistent $F$-statistics and PCs, while incorporating sampling error and missing data.

15

# References

Aman Agrawal, Alec M. Chiu, Minh Le, Eran Halperin, and Sriram Sankararaman. Scalable probabilistic PCA for large-scale genetic variation data. *PLOS Genetics*, 16(5):e1008773, 2020. ISSN 1553-7404. URL https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1008773.

David H. Alexander, John Novembre, and Kenneth Lange. Fast model-based estimation of ancestry in unrelated individuals. *Genome research*, 19(9):1655–1664, 2009. URL http://genome.cshlp.org/content/19/9/1655.short.

Isabel Alves, Miguel Arenas, Mathias Currat, Anna Sramkova Hanulova, Vitor C. Sousa, Nicolas Ray, and Laurent Excoffier. Long-distance dispersal shaped patterns of human genetic diversity in Eurasia. *Molecular biology and evolution*, 33(4):946–958, 2016.

Gideon S. Bradburd, Peter L. Ralph, and Graham M. Coop. Disentangling the Effects of Geographic and Ecological Isolation on Genetic Differentiation. *Evolution*, 67(11):3258–3273, 2013. ISSN 1558-5646. URL http://onlinelibrary.wiley.com/doi/10.1111/evo.12193/abstract.

Gideon S. Bradburd, Graham M. Coop, and Peter L. Ralph. Inferring continuous and discrete population genetic structure across space. *Genetics*, 210(1):33–52, 2018.

Abra Brisbin, Katarzyna Bryc, Jake Byrnes, Fouad Zakharia, Larsson Omberg, Jeremiah Degenhardt, Andrew Reynolds, Harry Ostrer, Jason G. Mezey, and Carlos D. Bustamante. PCAdmix: Principal Components-Based Assignment of Ancestry along Each Chromosome in Individuals with Admixed Ancestry from Two or More Populations. *Human biology*, 84(4):343–364, August 2012. ISSN 0018-7143. URL https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3740525/.

L. L. Cavalli-Sforza and A. W. F. Edwards. Phylogenetic Analysis: Models and Estimation Procedures. *Evolution*, 21(3):550–570, 1967. ISSN 0014-3820. URL http://www.jstor.org/stable/2406616.

L. L. Cavalli-Sforza and A. Piazza. Analysis of evolution: Evolutionary rates, independence and treeness. *Theoretical Population Biology*, 8(2):127–165, October 1975. ISSN 0040-5809. URL http://www.sciencedirect.com/science/article/pii/0040580975900295.

L. L Cavalli-Sforza, I. Barrai, and A. W. F Edwards. Analysis of Human Evolution Under Random Genetic Drift. *Cold Spring Harbor Symposia on Quantitative Biology*, 29:9–20, January 1964. ISSN 0091-7451, 1943-4456. URL http://symposium.cshlp.org/content/29/9.

L. L. Cavalli-Sforza, P. Menozzi, and A. Piazza. *The history and geography of human genes*. Princeton university press, 1994.

Barbara E. Engelhardt and Matthew Stephens. Analysis of Population Structure: A Unifying Framework and Novel Methods Based on Sparse Factor Analysis. *PLoS Genet*, 6(9):e1001117, September 2010. URL http://dx.doi.org/10.1371/journal.pgen.1001117.

Laurent Excoffier, Isabelle Dupanloup, Emilia Huerta-Sánchez, Vitor C. Sousa, and Matthieu Foll. Robust Demographic Inference from Genomic and SNP Data. *PLOS Genetics*, 9(10):e1003905, October 2013. ISSN 1553-7404. URL https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1003905.

J Felsenstein. Maximum-likelihood estimation of evolutionary trees from continuous characters. *American Journal of Human Genetics*, 25(5):471–492, September 1973. ISSN 0002-9297. URL http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1762641/.

J. C. Gower. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika*, 53(3-4):325–338, December 1966. ISSN 0006-3444. URL `https://doi.org/10.1093/biomet/53.3-4.325`.

Simon Gravel, Brenna M. Henn, Ryan N. Gutenkunst, Amit R. Indap, Gabor T. Marth, Andrew G. Clark, Fuli Yu, Richard A. Gibbs, Carlos D. Bustamante, David L. Altshuler, Richard M. Durbin, Gonçalo R. Abecasis, David R. Bentley, Aravinda Chakravarti, Andrew G. Clark, Francis S. Collins, Francisco M. De La Vega, Peter Donnelly, Michael Egholm, Paul Flicek, Stacey B. Gabriel, Richard A. Gibbs, Bartha M. Knoppers, Eric S. Lander, Hans Lehrach, Elaine R. Mardis, Gil A. McVean, Debbie A. Nickerson, Leena Peltonen, Alan J. Schafer, Stephen T. Sherry, Jun Wang, Richard K. Wilson, Richard A. Gibbs, David Deiros, Mike Metzker, Donna Muzny, Jeff Reid, David Wheeler, Jun Wang, Jingxiang Li, Min Jian, Guoqing Li, Ruiqiang Li, Huiqing Liang, Geng Tian, Bo Wang, Jian Wang, Wei Wang, Huanming Yang, Xiuqing Zhang, Huisong Zheng, Eric S. Lander, David L. Altshuler, Lauren Ambrogio, Toby Bloom, Kristian Cibulskis, Tim J. Fennell, Stacey B. Gabriel, David B. Jaffe, Erica Shefler, Carrie L. Sougnez, David R. Bentley, Niall Gormley, Sean Humphray, Zoya Kingsbury, Paula Koko-Gonzales, Jennifer Stone, Kevin J. McKernan, Gina L. Costa, Jeffry K. Ichikawa, Clarence C. Lee, Ralf Sudbrak, Hans Lehrach, Tatiana A. Borodina, Andreas Dahl, Alexey N. Davydov, Peter Marquardt, Florian Mertes, Wilfried Nietfeld, Philip Rosenstiel, Stefan Schreiber, Aleksey V. Soldatov, Bernd Timmermann, Marius Tolzmann, Michael Egholm, Jason Affourtit, Dana Ashworth, Said Attiya, Melissa Bachorski, Eli Buglione, Adam Burke, Amanda Caprio, Christopher Celone, Shauna Clark, David Conners, Brian Desany, Lisa Gu, Lorri Guccione, Kalvin Kao, Andrew Kebbel, Jennifer Knowlton, Matthew Labrecque, Louise McDade, Craig Mealmaker, Melissa Minderman, Anne Nawrocki, Faheem Niazi, Kristen Pareja, Ravi Ramenani, David Riches, Wanmin Song, Cynthia Turcotte, Shally Wang, Elaine R. Mardis, Richard K. Wilson, David Dooling, Lucinda Fulton, Robert Fulton, George Weinstock, Richard M. Durbin, John Burton, David M. Carter, Carol Churcher, Alison Coffey, Anthony Cox, Aarno Palotie, Michael Quail, Tom Skelly, James Stalker, Harold P. Swerdlow, Daniel Turner, Anniek De Witte, Shane Giles, Richard A. Gibbs, David Wheeler, Matthew Bainbridge, Danny Challis, Aniko Sabo, Fuli Yu, Jin Yu, Jun Wang, Xiaodong Fang, Xiaosen Guo, Ruiqiang Li, Yingrui Li, Ruibang Luo, Shuaishuai Tai, Honglong Wu, Hancheng Zheng, Xiaole Zheng, Yan Zhou, Guoqing Li, Jian Wang, Huanming Yang, Gabor T. Marth, Erik P. Garrison, Weichun Huang, Amit Indap, Deniz Kural, Wan-Ping Lee, Wen Fung Leong, Aaron R. Quinlan, Chip Stewart, Michael P. Stromberg, Alistair N. Ward, Jiantao Wu, Charles Lee, Ryan E. Mills, Xinghua Shi, Mark J. Daly, Mark A. DePristo, David L. Altshuler, Aaron D. Ball, Eric Banks, Toby Bloom, Brian L. Browning, Kristian Cibulskis, Tim J. Fennell, Kiran V. Garimella, Sharon R. Grossman, Robert E. Handsaker, Matt Hanna, Chris Hartl, David B. Jaffe, Andrew M. Kernytsky, Joshua M. Korn, Heng Li, Jared R. Maguire, Steven A. McCarroll, Aaron McKenna, James C. Nemesh, Anthony A. Philippakis, Ryan E. Poplin, Alkes Price, Manuel A. Rivas, Pardis C. Sabeti, Stephen F. Schaffner, Erica Shefler, Ilya A. Shlyakhter, David N. Cooper, Edward V. Ball, Matthew Mort, Andrew D. Phillips, Peter D. Stenson, Jonathan Sebat, Vladimir Makarov, Kenny Ye, Seungtai C. Yoon, Carlos D. Bustamante, Andrew G. Clark, Adam Boyko, Jeremiah Degenhardt, Simon Gravel, Ryan N. Gutenkunst, Mark Kaganovich, Alon Keinan, Phil Lacroute, Xin Ma, Andy Reynolds, Laura Clarke, Paul Flicek, Fiona Cunningham, Javier Herrero, Stephen Keenen, Eugene Kulesha, Rasko Leinonen, William M. McLaren, Rajesh Radhakrishnan, Richard E. Smith, Vadim Zalunin, Xiangqun Zheng-Bradley, Jan O. Korbel, Adrian M. Stütz, Sean Humphray, Markus Bauer, R. Keira Cheetham, Tony Cox, Michael Eberle, Terena James, Scott Kahn, Lisa Murray, Aravinda Chakravarti, Kai Ye, Francisco M. De La Vega, Yutao Fu, Fiona C. L. Hyland, Jonathan M. Manning, Stephen F. McLaughlin, Heather E. Peckham, Onur Sakarya, Yongming A. Sun, Eric F. Tsung, Mark A. Batzer, Miriam K. Konkel, Jerilyn A. Walker, Ralf Sudbrak, Marcus W. Albrecht, Vyacheslav S. Amstislavskiy, Ralf Herwig, Dimitri V. Parkhomchuk, Stephen T.

Sherry, Richa Agarwala, Hoda M. Khouri, Aleksandr O. Morgulis, Justin E. Paschall, Lon D. Phan, Kirill E. Rotmistrovsky, Robert D. Sanders, Martin F. Shumway, Chunlin Xiao, Gil A. McVean, Adam Auton, Zamin Iqbal, Gerton Lunter, Jonathan L. Marchini, Loukas Moutsianas, Simon Myers, Afidalina Tumian, Brian Desany, James Knight, Roger Winer, David W. Craig, Steve M. Beckstrom-Sternberg, Alexis Christoforides, Ahmet A. Kurdoglu, John V. Pearson, Shripad A. Sinari, Waibhav D. Tembe, David Haussler, Angie S. Hinrichs, Sol J. Katzman, Andrew Kern, Robert M. Kuhn, Molly Przeworski, Ryan D. Hernandez, Bryan Howie, Joanna L. Kelley, S. Cord Melton, Gonçalo R. Abecasis, Yun Li, Paul Anderson, Tom Blackwell, Wei Chen, William O. Cookson, Jun Ding, Hyun Min Kang, Mark Lathrop, Liming Liang, Miriam F. Moffatt, Paul Scheet, Carlo Sidore, Matthew Snyder, Xiaowei Zhan, Sebastian Zöllner, Philip Awadalla, Ferran Casals, Youssef Idaghdour, John Keebler, Eric A. Stone, Martine Zilversmit, Lynn Jorde, Jinchuan Xing, Evan E. Eichler, Gozde Aksay, Can Alkan, Iman Hajirasouliha, Fereydoun Hormozdiari, Jeffrey M. Kidd, S. Cenk Sahinalp, Peter H. Sudmant, Elaine R. Mardis, Ken Chen, Asif Chinwalla, Li Ding, Daniel C. Koboldt, Mike D. McLellan, David Dooling, George Weinstock, John W. Wallis, Michael C. Wendl, Qunyuan Zhang, Richard M. Durbin, Cornelis A. Albers, Qasim Ayub, Senduran Balasubramaniam, Jeffrey C. Barrett, David M. Carter, Yuan Chen, Donald F. Conrad, Petr Danecek, Emmanouil T. Dermitzakis, Min Hu, Ni Huang, Matt E. Hurles, Hanjun Jin, Luke Jostins, Thomas M. Keane, Si Quang Le, Sarah Lindsay, Quan Long, Daniel G. MacArthur, Stephen B. Montgomery, Leopold Parts, James Stalker, Chris Tyler-Smith, Klaudia Walter, Yujun Zhang, Mark B. Gerstein, Michael Snyder, Alexej Abyzov, Suganthi Balasubramanian, Robert Bjornson, Jiang Du, Fabian Grubert, Lukas Habegger, Rajini Haraksingh, Justin Jee, Ekta Khurana, Hugo Y. K. Lam, Jing Leng, Xinmeng Jasmine Mu, Alexander E. Urban, Zhengdong Zhang, Yingrui Li, Ruibang Luo, Gabor T. Marth, Erik P. Garrison, Deniz Kural, Aaron R. Quinlan, Chip Stewart, Michael P. Stromberg, Alistair N. Ward, Jiantao Wu, Charles Lee, Ryan E. Mills, Xinghua Shi, Steven A. McCarroll, Eric Banks, Mark A. DePristo, Robert E. Handsaker, Chris Hartl, Joshua M. Korn, Heng Li, James C. Nemesh, Jonathan Sebat, Vladimir Makarov, Kenny Ye, Seungtai C. Yoon, Jeremiah Degenhardt, Mark Kaganovich, Laura Clarke, Richard E. Smith, Xiangqun Zheng-Bradley, Jan O. Korbel, Sean Humphray, R. Keira Cheetham, Michael Eberle, Scott Kahn, Lisa Murray, Kai Ye, Francisco M. De La Vega, Yutao Fu, Heather E. Peckham, Yongming A. Sun, Mark A. Batzer, Miriam K. Konkel, Jerilyn A. Walker, Chunlin Xiao, Zamin Iqbal, Brian Desany, Tom Blackwell, Matthew Snyder, Jinchuan Xing, Evan E. Eichler, Gozde Aksay, Can Alkan, Iman Hajirasouliha, Fereydoun Hormozdiari, Jeffrey M. Kidd, Ken Chen, Asif Chinwalla, Li Ding, Mike D. McLellan, John W. Wallis, Matt E. Hurles, Donald F. Conrad, Klaudia Walter, Yujun Zhang, Mark B. Gerstein, Michael Snyder, Alexej Abyzov, Jiang Du, Fabian Grubert, Rajini Haraksingh, Justin Jee, Ekta Khurana, Hugo Y. K. Lam, Jing Leng, Xinmeng Jasmine Mu, Alexander E. Urban, Zhengdong Zhang, Richard A. Gibbs, Matthew Bainbridge, Danny Challis, Cristian Coafra, Huyen Dinh, Christie Kovar, Sandy Lee, Donna Muzny, Lynne Nazareth, Jeff Reid, Aniko Sabo, Fuli Yu, Jin Yu, Gabor T. Marth, Erik P. Garrison, Amit Indap, Wen Fung Leong, Aaron R. Quinlan, Chip Stewart, Alistair N. Ward, Jiantao Wu, Kristian Cibulskis, Tim J. Fennell, Stacey B. Gabriel, Kiran V. Garimella, Chris Hartl, Erica Shefler, Carrie L. Sougnez, Jane Wilkinson, Andrew G. Clark, Simon Gravel, Fabian Grubert, Laura Clarke, Paul Flicek, Richard E. Smith, Xiangqun Zheng-Bradley, Stephen T. Sherry, Hoda M. Khouri, Justin E. Paschall, Martin F. Shumway, Chunlin Xiao, Gil A. McVean, Sol J. Katzman, Gonçalo R. Abecasis, Tom Blackwell, Elaine R. Mardis, David Dooling, Lucinda Fulton, Robert Fulton, Daniel C. Koboldt, Richard M. Durbin, Senduran Balasubramaniam, Allison Coffey, Thomas M. Keane, Daniel G. MacArthur, Aarno Palotie, Carol Scott, James Stalker, Chris Tyler-Smith, Mark B. Gerstein, Suganthi Balasubramanian, Aravinda Chakravarti, Bartha M. Knoppers, Gonçalo R. Abecasis, Carlos D. Bustamante, Neda Gharani, Richard A. Gibbs, Lynn Jorde, Jane S. Kaye, Alastair Kent, Taosha Li, Amy L. McGuire, Gil A. McVean, Pilar N. Ossorio, Charles N. Rotimi, Yeyang Su, Lorraine H. Toji, Chris TylerSmith, Lisa D. Brooks, Adam L. Felsenfeld, Jean E. McEwen, As-

sya Abdallah, Christopher R. Juenger, Nicholas C. Clemm, Francis S. Collins, Audrey Duncanson, Eric D. Green, Mark S. Guyer, Jane L. Peterson, Alan J. Schafer, Gonçalo R. Abecasis, David L. Altshuler, Adam Auton, Lisa D. Brooks, Richard M. Durbin, Richard A. Gibbs, Matt E. Hurles, and Gil A. McVean. Demographic history and rare allele sharing among human populations. *Proceedings of the National Academy of Sciences*, page 201019276, July 2011. ISSN 0027-8424, 1091-6490. URL `http://www.pnas.org/content/early/2011/06/30/1019276108`.

R.E. Green, J. Krause, A.W. Briggs, T. Maricic, U. Stenzel, M. Kircher, N. Patterson, H. Li, W. Zhai, M.H.Y. Fritz, et al. A draft sequence of the Neandertal genome. *science*, 328(5979):710, 2010.

Ryan N. Gutenkunst, Ryan D. Hernandez, Scott H. Williamson, and Carlos D. Bustamante. Inferring the Joint Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. *PLoS Genet*, 5(10):e1000695, October 2009. URL `http://dx.doi.org/10.1371/journal.pgen.1000695`.

Wolfgang Haak, Iosif Lazaridis, Nick Patterson, Nadin Rohland, Swapan Mallick, Bastien Llamas, Guido Brandt, Susanne Nordenfelt, Eadaoin Harney, Kristin Stewardson, Qiaomei Fu, Alissa Mittnik, Eszter Bánffy, Christos Economou, Michael Francken, Susanne Friederich, Rafael Garrido Pena, Fredrik Hallgren, Valery Khartanovich, Aleksandr Khokhlov, Michael Kunst, Pavel Kuznetsov, Harald Meller, Oleg Mochalov, Vayacheslav Moiseyev, Nicole Nicklisch, Sandra L. Pichler, Roberto Risch, Manuel A. Rojo Guerra, Christina Roth, Anna Szécsényi-Nagy, Joachim Wahl, Matthias Meyer, Johannes Krause, Dorcas Brown, David Anthony, Alan Cooper, Kurt Werner Alt, and David Reich. Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature*, 522(7555):207–211, June 2015. ISSN 0028-0836. URL `http://www.nature.com/nature/journal/v522/n7555/full/nature14317.html`.

Matthew Hahn. *Molecular Population Genetics*. Oxford University Press, Oxford, New York, August 2018. ISBN 978-0-87893-965-7.

Eadaoin Harney, Nick Patterson, David Reich, and John Wakeley. Assessing the performance of qpAdm: a statistical tool for studying population admixture. *Genetics*, 217(4), April 2021. ISSN 1943-2631. URL `https://doi.org/10.1093/genetics/iyaa045`.

Trevor Hastie, Rahul Mazumder, Jason D. Lee, and Reza Zadeh. Matrix completion and low-rank SVD via fast alternating least squares. *The Journal of Machine Learning Research*, 16(1):3367–3402, January 2015. ISSN 1532-4435.

Daniel H. Huson, Regula Rupp, and Celine Scornavacca. *Phylogenetic networks: concepts, algorithms and applications*. Cambridge University Press, 2010. URL `https://books.google.com/books?hl=en&lr=&id=0rB5I5GxveAC&oi=fnd&pg=PR5&dq=huson+phylogenetic+networks&ots=BaKyTHg9EO&sig=HrZB-uEusSsveNCDJEedODh7UHg`.

I. T. Jolliffe. *Principal Component Analysis*. Springer Science & Business Media, March 2013. ISBN 978-1-4757-1904-8.

John A. Kamm, Jonathan Terhorst, and Yun S. Song. Efficient computation of the joint sample frequency spectra for multiple populations. *arXiv:1503.01133 [math, q-bio]*, March 2015. URL `http://arxiv.org/abs/1503.01133`.

Iosif Lazaridis, Nick Patterson, Alissa Mittnik, Gabriel Renaud, Swapan Mallick, Karola Kirsanow, Peter H. Sudmant, Joshua G. Schraiber, Sergi Castellano, Mark Lipson, and others. Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature*, 513(7518): 409–413, 2014. URL `http://www.nature.com/nature/journal/v513/n7518/abs/nature13673.html`.

Mark Lipson, Po-Ru Loh, Alex Levin, David Reich, Nick Patterson, and Bonnie Berger. Efficient Moment-Based Inference of Admixture Parameters and Sources of Gene Flow. *Molecular Biology and Evolution*, 30(8):1788–1802, August 2013. ISSN 0737-4038, 1537-1719. URL http://mbe.oxfordjournals.org/content/30/8/1788.

Gil McVean. A genealogical interpretation of principal components analysis. *PLoS genetics*, 5(10): e1000686, October 2009. ISSN 1553-7404.

Jonas Meisner, Siyang Liu, Mingxi Huang, and Anders Albrechtsen. Large-scale Inference of Population Structure in Presence of Missingness using PCA. *Bioinformatics (Oxford, England)*, page btab027, January 2021. ISSN 1367-4811.

J. Novembre and M. Stephens. Interpreting principal component analyses of spatial population genetic variation. *Nature genetics*, 40(5):646–649, 2008. URL http://www.nature.com/ng/journal/v40/n5/abs/ng.139.html.

John Novembre, Toby Johnson, Katarzyna Bryc, Zoltán Kutalik, Adam R Boyko, Adam Auton, Amit Indap, Karen S King, Sven Bergmann, Matthew R Nelson, Matthew Stephens, and Carlos D Bustamante. Genes mirror geography within Europe. *Nature*, 456(7218):98–101, 2008. URL http://www.ncbi.nlm.nih.gov/pubmed/18758442.

Gonzalo Oteo-Garcia and Jose-Angel Oteo. A geometrical framework for f-statistics. *Bulletin of Mathematical Biology*, 83(2):1–22, 2021.

Lior Pachter. What is principal component analysis?, May 2014. URL https://liorpachter.wordpress.com/2014/05/26/what-is-principal-component-analysis/.

Nick Patterson, Daniel J. Richter, Sante Gnerre, Eric S. Lander, and David Reich. Genetic evidence for complex speciation of humans and chimpanzees. *Nature*, 441(7097):1103–1108, June 2006. ISSN 0028-0836. URL http://www.nature.com/nature/journal/v441/n7097/abs/nature04789.html.

Nick J. Patterson, Priya Moorjani, Yontao Luo, Swapan Mallick, Nadin Rohland, Yiping Zhan, Teri Genschoreck, Teresa Webster, and David Reich. Ancient Admixture in Human History. *Genetics*, page genetics.112.145037, September 2012. ISSN 0016-6731, 1943-2631. URL http://www.genetics.org/content/early/2012/09/06/genetics.112.145037.

Benjamin M. Peter. Admixture, Population Structure and F-Statistics. *Genetics*, page genetics.115.183913, January 2016. ISSN 0016-6731, 1943-2631. URL http://www.genetics.org/content/early/2016/02/03/genetics.115.183913.

Benjamin M. Peter, Desislava Petkova, and John Novembre. Genetic landscapes reveal how human genetic diversity aligns with geography. *Molecular biology and evolution*, 37(4):943–951, 2020.

Martin Petr, Svante Pääbo, Janet Kelso, and Benjamin Vernot. Limits of long-term selection against Neandertal introgression. *Proceedings of the National Academy of Sciences*, 116(5):1639–1644, January 2019.

Joseph K. Pickrell and David Reich. Toward a new history and geography of human genes informed by ancient DNA. *Trends in Genetics*, 30(9):377–389, September 2014. ISSN 0168-9525. URL http://www.sciencedirect.com/science/article/pii/S0168952514001206.

Jonathan K Pritchard, Matthew Stephens, and Peter Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000. URL http://www.ncbi.nlm.nih.gov/pubmed/10835412.

Fernando Racimo, Jessie Woodbridge, Ralph M. Fyfe, Martin Sikora, Karl-Göran Sjögren, Kristian Kristiansen, and Marc Vander Linden. The spatiotemporal spread of human migrations during the European Holocene. *Proceedings of the National Academy of Sciences*, 117(16):8989–9000, April 2020.

Maanasa Raghavan, Pontus Skoglund, Kelly E. Graf, Mait Metspalu, Anders Albrechtsen, Ida Moltke, Simon Rasmussen, Thomas W. Stafford Jr, Ludovic Orlando, Ene Metspalu, and others. Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans. *Nature*, 505(7481): 87–91, 2014. URL `http://www.nature.com/nature/journal/v505/n7481/abs/nature12736.html`.

Peter Ralph and Graham Coop. The Geography of Recent Genetic Ancestry across Europe. *PLoS Biol*, 11(5):e1001555, May 2013. URL `http://dx.doi.org/10.1371/journal.pbio.1001555`.

Sohini Ramachandran, Omkar Deshpande, Charles C Roseman, Noah A Rosenberg, Marcus W Feldman, and L. Luca Cavalli-Sforza. Support from the relationship of genetic and geographic distance in human populations for a serial founder effect originating in Africa. *Proceedings of the National Academy of Sciences of the United States of America*, 102(44):15942–15947, 2005. ISSN 0027-8424, 1091-6490. URL `http://www.pnas.org/content/102/44/15942`.

D. Reich, K. Thangaraj, N. Patterson, A. L. Price, and L. Singh. Reconstructing Indian population history. *Nature*, 461(7263):489–494, 2009.

David Reich. *Who We Are and How We Got Here: Alte DNA und die neue Wissenschaft der menschlichen Vergangenheit.* Pantheon, New York, illustrated edition edition, 2018. ISBN 978-1-101-87032-7.

Noah A. Rosenberg, Jonathan K. Pritchard, James L. Weber, Howard M. Cann, Kenneth K. Kidd, Lev A. Zhivotovsky, and Marcus W. Feldman. Genetic structure of human populations. *Science (New York, N.Y.)*, 298(5602):2381–2385, December 2002. ISSN 1095-9203.

Noah A Rosenberg, Saurabh Mahajan, Sohini Ramachandran, Chengfeng Zhao, Jonathan K Pritchard, and Marcus W Feldman. Clines, Clusters, and the Effect of Study Design on the Inference of Human Population Structure. *PLoS Genet*, 1(6):e70, December 2005. URL `http://dx.plos.org/10.1371/journal.pgen.0010070`.

Joshua G. Schraiber and Joshua M. Akey. Methods and models for unravelling human evolutionary history. *Nature Reviews Genetics*, 2015. URL `http://www.nature.com/nrg/journal/vaop/ncurrent/full/nrg4005.html`.

Charles Semple and M. A. Steel. *Phylogenetics.* Oxford University Press, 2003. ISBN 978-0-19-850942-4.

David Serre and Svante Pääbo. Evidence for Gradients of Human Genetic Diversity Within and Among Continents. *Genome Research*, 14(9):1679–1685, September 2004. ISSN 1088-9051, 1549-5469. URL `https://genome.cshlp.org/content/14/9/1679`.

M SLATKIN. GENE FLOW IN NATURAL-POPULATIONS. *Annual Review of Ecology and Systematics*, 16:393–430, 1985. ISSN 0066-4162.

Mark Stoneking. *An Introduction to Molecular Anthropology.* John Wiley & Sons, December 2016. ISBN 978-1-118-06162-6.

716 The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*, 526
717 (7571):68–74, October 2015. ISSN 0028-0836. URL `http://www.nature.com.proxy.uchicago.`
718 `edu/nature/journal/v526/n7571/full/nature15393.html`.

719 Sten Wahlund. Zusammensetzung Von Populationen Und Korrelationserscheinungen Vom Stand-
720 punkt Der Vererbungslehre Aus Betrachtet. *Hereditas*, 11(1):65–106, May 1928. ISSN 1601-
721 5223. URL `http://onlinelibrary.wiley.com/doi/10.1111/j.1601-5223.1928.tb02483.x/`
722 `abstract`.

# A  Derivations

Depending on a readers' background in linear algebra, these results may appear elementary; I include them here for reference and because they were not obvious to me at the onset of this project.

**$F$-statistics are invariant under a change-of-basis**

$$
\begin{aligned}
F_2(X_i, X_j) &= \sum_{l=1}^{S} \big( (x_{il} - \mu_l) - (x_{jl} - \mu_l) \big)^2 = F_2(Y_i, Y_j) \\
&= \sum_{l=1}^{S} \Big( \sum_k L_{kl} P_{ik} - \sum_k L_{kl} P_{jk} \Big)^2 \\
&= \sum_{l=1}^{S} \left( \sum_k L_{kl}(P_{ik} - P_{jk}) \right)^2 \\
&= \sum_{l=1}^{S} \left( \sum_k L_{kl}^2 (P_{ik} - P_{jk})^2 + 2 \sum_{k \neq k'} L_{kl} L_{k'l} (P_{ik} - P_{jk'})^2 \right) \\
&= \sum_k \underbrace{\left( \sum_{l=1}^{L} L_{kl}^2 \right)}_{1} (P_{ik} - P_{jk})^2 + 2 \sum_{k \neq k'} \underbrace{\left( \sum_{l=1}^{S} L_{kl} L_{k'l} \right)}_{0} (P_{ik} - P_{jk'})^2 \\
&= \sum_k (P_{ik} - P_{jk})^2 \tag{A1}
\end{aligned}
$$

In summary, the first row shows that $F_2$ on the centered data will give the same results (as distances are invariant to translations), in the second row we apply the PC-decomposition. The third row is obtained from factoring out $L_{lk}$. Row four is obtained by multiplying out the sum inside the square term for a particular $l$. We have $k$ terms when for $\binom{k}{2}$ terms for different $k$'s. Row five is obtained by expanding the outer sum and grouping terms by $k$. The final line is obtained by recognizing that $\mathbf{L}$ is an orthonormal basis; where dot products of different vectors have lengths zero.

Note that if we estimate $F_2$, unbiased estimators are obtained by subtracting the population-heterozygosities $H_i, H_j$ from the statistic. As these are scalars, they do not change above calculation.

**The region of negative $F_3$-statistics is a $n$-ball**  Without loss of generality, assume that $X_1 = (r, 0, 0, \dots)$ and $X_2 = (-r, 0, 0, \dots)$, and let us assume that $X_x$ has coordinates $(x_1, x_2, \dots, x_S)$ Assuming $F_3(X_x; X_1, X_2) = 0$, equation 13 becomes

$$
\begin{aligned}
2F_3(X_x; X_1, X_2) &= \|X_x - X_1\|^2 + \|X_x - X_2\|^2 - \|X_1 - X_2\|^2 = 0 \\
&= \left[ (x_1 - r)^2 + \sum_{i=2}^{S} x_i^2 \right] + \left[ (x_1 + r)^2 + \sum_{i=2}^{S} x_i^2 \right] - 4r^2 \\
&= 2 \left[ \sum_{i=1}^{S} x_i^2 + r^2 + x_1 r - x_1 r \right] - 4r^2 \\
F_3(X_x; X_1, X_2) &= -r^2 + \sum_{i=1}^{S} x_i^2 = -r^2 + \|X_x\|^2 = 0, \tag{A2}
\end{aligned}
$$

which is the equation of a $n$-sphere with radius $r$ and center at the origin, as assumed from the placing of $X_1$ and $X_2$. Now, assume that $F_3$ is negative, i.e. $F_3(X_x; X_1, X_2) = -k < 0$. Moving $r^2$ to the left we obtain

$$
r^2 - k = \|X_x\|^2, \tag{A3}
$$

which is another $n$-sphere with a smaller radius, showing that all points inside the $n$-sphere will have negative $F_3$-values.

**If a population lies outside the circle of this $n$-Sphere in any 2D-projection, $F_3$ is positive**

Assume the center of the $n$-sphere $C = \frac{X_1+X_2}{2} = (c_1, c_2, \ldots c_S)$, and $X_x = (x_1, x_2, \ldots x_S)$. Then,

$$F_3(X_x; X_1, X_2) = \|X_x - C\|^2 - r^2$$

$$= \underbrace{(x_1 - c_1)^2 + (x_2 - c_2)^2}_{>r^2} + \underbrace{\sum_{i=3}^{S}(x_i - c_i)^2}_{\geq 0} - r^2$$

$$> 0. \tag{A4}$$

The condition $(x_1 - c_1)^2 + (x_2 - c_2)^2 > r^2$ is satisfied whenever $X_x$ is outside the circle obtained from projecting the $n$-sphere on the first two dimensions. An analogous argument applies for any low-dimensional representation.