

# Correcting principal component maps for effects of spatial autocorrelation in population genetic data

Eric Frichot, Sean D Schoville, Guillaume Bouchard and Olivier Francois

Journal Name:	Frontiers in Genetics
ISSN:	1664-8021
Article type:	Original Research Article
Received on:	21 Aug 2012
Accepted on:	29 Oct 2012
Provisional PDF published on:	29 Oct 2012
Frontiers website link:	www.frontiersin.org
Citation:	Frichot E, Schoville SD, Bouchard G and Francois O(2012) Correcting principal component maps for effects of spatial autocorrelation in population genetic data. 3:254. doi:10.3389/fgene.2012.00254
Article URL:	http://www.frontiersin.org/Journal/Abstract.aspx?s=1265& name=applied%20genetic%20epidemiology&ART_DOI=10.3389 /fgene.2012.00254
	(If clicking on the link doesn't work, try copying and pasting it into your browser.)
Copyright statement:	© 2012 Frichot, Schoville, Bouchard and Francois. This is an open-access article distributed under the terms of the <u>Creative</u> <u>Commons Attribution License</u> , which permits use, distribution and reproduction in other forums, provided the original authors and source are credited and subject to any copyright notices concerning any third-party graphics etc.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after rigorous peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.

1	Correcting principal component maps for effects of spatial autocorrelation in population genetic data
2	autocorrelation in population genetic data
3	Eric Frichot <sup>1</sup> , Sean Schoville <sup>1</sup> , Guillaume Bouchard <sup>2</sup> , Olivier François <sup>1</sup> *
4 5 6	<ol> <li>Université Joseph Fourier Grenoble, Centre National de la Recherche, TIMC-IMAG UMR 5525, Grenoble, France</li> <li>Xerox Research Center Europe, Meylan, France</li> </ol>
7	Abstract
8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24	In many species, spatial genetic variation displays patterns of "isolation- by-distance". Characterized by locally correlated allele frequencies, these pat- terns are known to create periodic shapes in geographic maps of principal components which confound signatures of specific migration events and in- fluence interpretations of principal component analyses (PCA). In this study, we introduced models combining probabilistic PCA and kriging models to in- fer population genetic structure from genetic data while correcting for effects generated by spatial autocorrelation. The corresponding algorithms are based on singular value decomposition and low rank approximation of the genotypic data. As their complexity is close to that of PCA, these algorithms scale with the dimension of the data. To illustrate the utility of these new models, we simulated isolation-by-distance patterns and broad-scale geographic variation using spatial coalescent models. Our methods remove the horseshoe patterns usually observed in PC maps and simplify interpretations of spatial genetic variation. We demonstrate our approach by analyzing single nucleotide poly- morphism data from the Human Genome Diversity Panel, and provide com- parisons with other recently introduced methods.

3700 words, 8 figures.

25

### 26 Correspondence:

- 27 Dr Olivier François
- 28 Grenoble INP
- 29 Laboratoire TIMC-IMAG
- 30 Faculty of Medicine,
- 31 F38706 La Tronche, France
- 32 olivier.francois@imag.fr
- 33 Running Title: Correcting principal component analysis
- 34 Keywords: principal component analysis, isolation-by-distance, spatial autocorrelation,
- 35 spatial factor analysis

#### 36 Introduction

The concept of "isolation-by-distance" (IBD) was introduced by S. Wright to describe 37 the accumulation of local genetic differences under spatially restricted dispersal (Wright, 38 1943). In species that are continuously distributed in geographic space and disperse over 39 short distances, the theory predicts that genetic differentation will increase with geographic 40 distance (Kimura and Weiss, 1964; Malécot, 1948). IBD can be described by spatial au-41 tocorrelation, a measure of the degree of dependency among observations in a geographic 42 space. Although studying IBD patterns could lead to useful estimates of gene dispersal 43 (Rousset, 1997), spatial autocorrelation derived from IBD often presents a problem for 44 population genetic analyses. More specifically, the presence of spatial autocorrelation pat-45 terns can increase the rate of false positive tests for hierarchical population structure or for 46 the detection of loci under selection (Meirmans, 2012). 47

Recently, it has been acknowledged that distortions caused by spatial autocorrelation 48 could also bias interpretations of population genetic structure as inferred from PCA or 49 50 from Bayesian clustering methods (François et al., 2010; Novembre and Stephens, 2008). Principal component analysis (PCA) is a method that searches for axes, called principal 51 components, along which projected individuals show the highest variance. As a result, the 52 first PCs are often used to explore the structure of variation in the sample. Characterized 53 by locally correlated allele frequencies, IBD patterns create periodic shapes in PC maps 54 that can confound signatures of migration events and influence interpretations of princi-55 pal component analyses (Novembre and Stephens, 2008). In scenarios where covariance 56 57 decays exponentially with geographic distance, PC plots are indeed expected to exhibit 58 horseshoe effects, an artifact in which the second axis is curved relative to the first axis. 59 These effects lead to counterintuitive representations of the data (Diaconis et al., 2008;

60 Legendre and Gallagher, 2001).

Several methods have been proposed to correct for the effects of spatial autocorrela-61 tion in exploratory data analyses. In particular, those methods include spatial Principal 62 Component Analysis (sPCA, Borcard and Legendre 2002; Borcard et al. 2004; Dray et al. 63 2006; Jombart et al. 2008), and sparse factor analysis (SFA, Engelhardt and Stephens 64 65 2010). Generally the methods share the objective of separating local and regional geographic scales in the data. In this study, we introduce a novel approach, based on latent 66 67 factors models, that addresses the separation of geographic scales more directly than the 68 two previous methods. The new method, Spatial factor analysis (spFA), combines prob-69 abilistic PCA (Tipping and Bishop, 1999) and kriging models (Cressie, 1993) to infer 70 population genetic structure from genetic data while correcting for errors introduced by spatial autocorrelation. While many approaches have been argued to improve interpreta-71 tions of the data, their outputs have not yet been compared to each other on the basis of 72 spatial simulations. To compare methods, we generated patterns of IBD and broad-scale 73 geographic variation using computer simulations of spatial coalescent models. We com-74 pared the outcomes of methods under population genetic models of isolation-by-distance, 75 and we argued that the methods provided insights on distinct aspects of the data. We report 76 that the new spFA method was able to remove the horseshoe effect observed in spatially 77 structured data, whereas this was not the case in PCA, sPCA and SFA analyses. We dis-78 cuss the significance of this result in an assessment of single nucleotide polymorphism 79 data from worldwide samples of the Human Genome Diversity Panel. 80

#### 81 Material and Methods

We considered single nucleotide polymorphism (SNP) data for *n* individuals genotyped at *L* loci. For these data, the genotypic matrix entries,  $(G_{i\ell})$ , record the number of derived alleles at locus  $\ell$  for individual *i*. For autosomal data,  $G_{i\ell}$  is thus equal to 0, 1 or 2, and corresponds to the genotype at locus  $\ell$ . The data were centered by substracting the mean value of each column of *G* and scaled by dividing by the standard deviation value of each column of *G*. In addition to the genotypic data, we assumed that geographical coordinates,  $(X_i)$ , were recorded for each individual.

We evaluated the effects of IBD patterns on inference of population genetic structure using 4 statistical methods: Principal Component Analysis (PCA, Jolliffe 1986; Patterson et al. 2006), spatial PCA (sPCA, Jombart et al. 2008), Sparse Factor Analysis (SFA, Engelhardt and Stephens 2010), and a new method called *spatial Factor Analysis* (spFA).

93 **Principal Component Analysis.** PCA is a popular method that searches for a set of *K* 94 orthogonal axes (the principal components), each of which is a linear combination of the 95 original axes, such that projections of the original data display maximal variance onto the 96 new axes (McVean, 2009). We computed the score matrix, *U* of dimension  $n \times K$ , and the 97 loading matrix, *V* of dimensions  $K \times L$ , using the rank *K* singular value decomposition 98 method implemented in the R function prcomp and in the computer program *SmartPCA* 99 (Patterson et al., 2006).

Moran eigenvectors and spatial PCA. Moran eigenvectors maps were proposed as an
alternative to trend surface analysis for incorporating spatial variation in population genetics models (Dray et al., 2006; Jombart et al., 2008). In Moran eigenvectors maps, there are

positive and negative eigenvalues. Eigenvectors associated with positive eigenvalues have
positive autocorrelation, and they describe global structures. Eigenvectors associated with
negative eigenvalues describe local structures. Implemented in an algorithm called spatial
PCA (sPCA), Moran's eigenvector maps (MEM) maximize Moran's spatial autocorrelation index, defined as follows

$$I(g) = \frac{\sum_{i,j} w_{ij} (g_i - \bar{g}) (g_j - \bar{g})}{\sum_{i,j} w_{ij} \sum_i (g_i - \bar{g})^2}$$

108 with respect to a spatial weighting matrix, W, deduced from geographical distances (Dray
109 et al., 2006). We implemented MEMs and sPCA using the R package adegenet using a
110 Delaunay weighting matrix (Jombart et al., 2008).

111 Spatial Factor Analysis. We introduced a new spatial factor analysis model (spFA)
112 which incorporates spatial information in factor analysis in an explicit way. In spFA, infer113 ence was performed in a matrix factorization model similar to probabilistic PCA (Tipping
114 and Bishop, 1999)

$$G_{i\ell} = U_i^T V_\ell + \epsilon_{i\ell} \,, \tag{1}$$

115 where  $\epsilon_{i\ell}$  are statistically dependent Gaussian variables with mean zero and with covari-116 ance matrix  $\Sigma_{\theta}$ . Similarly to Kriging approaches (Cressie, 1993), a radial basis covariance 117 matrix was chosen to model spatial autocorrelation patterns generated by IBD (see also 118 Durand et al. 2009). The covariance matrix  $\Sigma_{\theta}$  was defined as follows. For all pairs of 119 individuals, *i* and *j*, we have

$$\Sigma_{\theta}(i,j) = \exp(-d(X_i, X_j)/\theta), \quad \theta > 0,$$
(2)

120 where  $d(X_i, X_j)$  represents the squared Euclidean or great-circle distance between sites 121 with coordinate  $X_i$  and with coordinate  $X_j$ . To avoid colinearity issues, we assumed that 122 the individual geographical coordinates were distinct to each other (ties were broken by 123 adding small perturbations to the original spatial coordinates). The parameter  $\theta$  is a scale 124 parameter measured in units of average pairwise distance between geographic sites,  $\bar{d}$ . In 125 practice, spFA required that an array of  $\theta$  values (scale parameter) were explored, and  $\theta$ 126 was varied in the range (0,  $10\bar{d}$ ).

127 To solve the spFA model, we used a Cholesky decomposition,  $C^T C = \Sigma_{\theta}^{-1}$ , and we 128 established an equivalence with the following matrix factorization model

$$\tilde{G}_{i\ell} = \tilde{U}_i^T \tilde{V}_\ell + \tilde{\epsilon}_{i\ell} , \qquad (3)$$

where  $\tilde{G} = CG$ ,  $\tilde{U} = CU$ ,  $\tilde{V} = V$  and where  $\tilde{\epsilon}_{\ell}$  are statistically independent Gaussian 129 vectors of mean zero and covariance matrix equal to identity. The matrix  $\tilde{U}$  and  $\tilde{V}$  were 130 obtained by applying a singular value decomposition of rank K to the transformed data 131 matrix, CG. Then, U and V were obtained by applying a singular value decomposition 132 of rank K to  $C^{-1}\tilde{U}\tilde{V}$ . To avoid multiple solutions, the orthogonality condition  $VV^T = I_K$ , 133 where  $I_K$  is the identity matrix in K dimensions, was imposed to V (Figure 1). The time 134 needed to compute SpFA is the same order as the time needed to compute K scores and 135 loadings for a standard PCA (Patterson et al., 2006). For an example of implementation, 136 137 see our R code (http://membres-timc.imag.fr/Olivier.Francois/spfa.R).

138 [Figure\_1.TIF]

139 Sparse Factor Analysis. Sparse Factor Analysis (SFA) was introduced by Engelhardt
140 and Stephens (2010) as an alternative to admixture-based models, and this method can

recapitulate the results of PCA when population structure is influenced by IBD patterns.To give a description of SFA, we considered a regression model of the following form

$$G_{i\ell} = U_i^T V_\ell + \epsilon_{i\ell} \tag{4}$$

in which the residual errors are independent Gaussian random variables,  $\epsilon_{i,\ell} \sim N(0, 1/\psi_i)$ , and where the prior distribution on the precision parameter,  $\psi_i$ , is a Gamma distribution. In the SFA model, an automatic relevance determination prior is considered for the score vectors,  $U_{ik} \sim N(0, \sigma_{ik}^2)$ , where some  $\sigma_{ik}^2$  are constrained to be equal to zero. We implemented SFA using the code distributed in (Engelhardt and Stephens, 2010), and we used 1,000 iterations. Eigenvectors in spFA and in SFA were also referred to as *factors* or *axes*.

149 Simulated data. We generated simulated data for two diverging populations using coalescent models implemented in the computer program ms (Hudson, 2002). In these mod-150 els, each population was simulated according to a linear stepping-stone model with 50 151 demes. To reproduce the simulation settings of Novembre and Stephens (2008), the ef-152 153 fective migration rate between pairs of adjacent demes was set to the value 4Nm = 1. The divergence time  $\tau$  between the two populations was varied within the range of values 154  $\tau = (0, 100)$  measured in coalescent units. We sampled 100 individuals, one from each 155 deme both side of a (fictive) geographic barrier. For each simulation, we evaluated Wilks' 156  $\Lambda$ , a statistic used in multivariate analysis of variance to test whether there are differences 157 158 between the means of identified groups of individuals on the combination of genotypes (Mardia et al., 1979). 159

#### 160 **Results**

161 **Pure isolation-by-distance patterns.** In a first series of experiments, we used simula-162 tions of one-dimensional stepping-stone models reproducing the patterns of IBD described 163 in Novembre and Stephens (2008). In those simulated data, the divergence time between 164 the two populations was thus set to  $\tau = 0$ , and the populations were connected by recurrent gene flow (4Nm = 1). As expected from theoretical results for PCA and for other 165 ordination methods (Novembre and Stephens 2008; Ahmed et al. 1974; Dray et al. 2006), 166 the first PC maps displayed oscillating patterns. In addition, the frequency of oscillation 167 increased as we examined axes of higher orders (Figure 2A). When we used sPCA, the 168 169 first three positive components were almost identical to those obtained with PCA (not reported). 170

Running spFA with K = 3 and with 3 distinct values of the scale parameter ( $\theta/\bar{d} =$ 171 0.1, 0.2 and 0.3) led to different interpretations of the genetic data (Figure 2B-D). Gradu-172 ally varying  $\theta$  allowed us to evaluate the scales at which the IBD effects were apparent, and 173 also allowed us to remove those effects sequentially. For  $\theta/\bar{d} = 0.1$ , the maps correspond-174 ing to factor 1 and 2 displayed sinusoidal curves similar to PC1 and PC2, whereas the map 175 for factor 3 was flat as expected if the effect of IBD is removed (Figure 2B). For  $\theta/\bar{d} = 0.2$ , 176 the map corresponding to factor 1 remained similar to PC1, but the maps for factor 2 and 177 factor 3 were flat (Figure 2C). For  $\theta/\bar{d} = 0.3$ , the effects of isolation by distance were 178 corrected in all axes (Figure 2D). 179

180 When we ran SFA with K = 3 factors, the resulting maps also emphasized aspects of 181 the data different from the ones described by PC maps and spatial factor maps (Figure 3). 182 Maps for SFA are interpreted in terms of clusters, similar to those obtained in non-spatial 183 Bayesian assignment programs like STRUCTURE (Pritchard et al., 2000). Clusters created by 184 clustering programs under IBD models are often reported as being undesirable (François185 and Durand, 2010; Meirmans, 2012).

**Two diverging populations with IBD patterns.** In a second series of experiments, we used simulations of a two-population model, where each population consisted of a linear network of 50 demes. In these experiments, the two populations were separated by a geographic barrier to gene flow.

First the divergence time was set to  $\tau = 10$  coalescent units. Using PCA, the first 2 components displayed oscillating patterns, similar to those obtained with  $\tau = 0$  (pure IBD simulations) (Figure 4A). The PC1-PC2 plot exbihited a clear horseshoe pattern. Differentiation between the two populations was visible in the PC1 map, where a discontinuity was observed at the center of the habitat. This discontinuity corresponded to the localization of the geographic barrier. Results for the positive eigenvectors of sPCA strongly resembled those obtained for the first PCs (Figure 4B).

Turning to spFA, we argued for a particular choice of  $\theta/\bar{d}$  based on Wilks' A statistic, 199 200 a standard measure of separation of groups in discriminant analysis, and computed this statistic for  $\theta/\bar{d}$  ranging between 0.01 and 10. As spatial factor analysis provided differ-201 202 ent interpretations of the data depending on the scale at which the data were analyzed, 203 the choice of  $\theta$  was crucial to the method. Figure 5 reports the value of Wilks'  $\Lambda$  as a function of the logarithm of  $\theta/\bar{d}$ . Values of  $\theta/\bar{d}$  minimizing Wilks' statistic and providing 204 205 the best description of our data into clusters were about 0.32 (Figure 5). When spFA was 206 applied with K = 2, the first factor map grouped demes at the left and the right of the 207 geographic barrier in two main clusters, while simultaneously correcting for IBD patterns 208 within the two clusters (Figure 4C). The spFA Axis1-Axis2 plot removed the horsehoe 209 effect observed in PCA and sPCA plots. The resulting figure emphasized a discontinuous 210 population structure consisting of two differentiated genetic clusters. Running SFA with 211 K = 2 also led to a description of the data in two genetic clusters, located both sides of 212 the geographic barrier, but the method failed to describe the two clusters as discontinuous 213 entities (Figure 4D).

214 Based on PC and factor plots, we next computed Wilks' A statistic for all methods, and 215 for divergence times  $\tau$  ranging between 0 and 100 (Figure 6). Lower values of  $\Lambda$  generally 216 indicated better discrimination of the 2 divergent populations in PC or factor plots. For 217 all methods, the  $\Lambda$  statistic decreased as the divergence time between the 2 populations increased (McVean 2009). In our spatially explicit framework, SFA (green curve) detected 218 219 the existence of diverging populations earlier than PCA (red curve) and than sPCA (not 220 shown, similar to PCA). SpFA was the most sensitive method, and provided an earlier detection of divergent clusters than SFA and PCA (blue curve). 221

- 222 [Figure\_4.TIF]
- 223 [Figure\_5.TIF]
- 224 [Figure\_6.TIF]

Human data analysis. Next we applied PCA, sPCA, SpFA and SFA to a worldwide sample of genomic DNA from 418 individuals in 27 Asian populations, from the Harvard Human Genome Diversity Project - Centre Etude Polymorphism Humain (Harvard HGDP-CEPH) (ftp://ftp.cephb.fr/hgdp\_v3/). In those data, each marker has been ascer-

tained in samples of Mongolian ancestry (referenced population HGDP01224). We selected all samples from Central and East-Asia at the exception of Xibe, who originated
in northeastern China, but migrated to northwestern China only recently (Powell et al.,
2007) (Figure 7A). The data set used a panel of 10,664 SNPs (see Patterson et al., 2012,
ftp://ftp.cephb.fr/hgdp\_supp10/).

234 In our analysis, samples from Central Asia, West to the Tibetan plateau, were rep-235 resented with red/orange colors, whereas populations from East-Asian were represented 236 with blue colors (Figure 7A). For those samples, the PC plot exhibited a horseshoe pattern, 237 which was a signature of the presence of IBD patterns in the data (Figure 7B). PCA led 238 to a continuum of samples without observable genetic discontinuities. Running spFA with K = 2 and setting  $\theta/\bar{d} = 10^{-2}$  on the basis of Wilks' statistic analysis, spFA corrected for 239 240 the effects of IBD in axes 1 and 2 (Figure 7C). The spFA method provided evidence of 241 a major discontinuity separating two clusters, one in Central Asia and one in East-Asia. In addition, Uyghur and Hazara population samples aligned with the two main clusters 242 and were placed in an intermediate position, suggesting genetic admixture from ancestral 243 244 Central Asian and East-Asian gene pools. Essentially the same patterns emerged when spFA was applied with K = 3 at the same scale (Figure 8C, 8D). 245

Using SFA with K = 2, factors 1 and 2 confirmed the main discontinuity, in a representation of clusters closer to Bayesian clustering methods than to PCA (Figure 7D). Uyghur and Hazara population samples were also placed between the main clusters. When we used SFA with K = 3, we obtained shapes without natural interpretations (Figure 8A, 8B). SFA detected additional discontinuities whereas the other methods suggested that continuous genetic variation in geographic space was predominant.

#### [Figure\_7.TIF]

252

#### 254 Discussion

Principal component analysis and related methods used to describe genomic variation 255 among large population samples are known to produce results that can be distorted by 256 257 IBD, and that may thus be difficult to interpret. The horseshoe effect is one of the distor-258 tions observed in PC plots that arises when covariance between allele frequencies decays 259 exponentially with geographic distance. In this case, there is an established mathematical 260 correspondence between the eigenvectors of the covariance matrix and the columns of a 261 discrete cosine-transform (Ahmed et al., 1974; Diaconis et al., 2008). In this study, we 262 used this correspondence to propose a new approach based on spatial models for the co-263 variance structure of residual errors in factor analysis. In spFA, IBD effects were modeled through the introduction of a covariance matrix that accounts for the geographic distance 264 265 between individuals explicitly.

266 We compared spFA to PCA and to two recent methods that also attempt to correct for IBD effects: spatial Principal Component Analysis (sPCA, Jombart et al. 2008) and sparse 267 factor analysis (SFA, Engelhardt and Stephens 2010). When we applied PCA to simulated 268 data from spatial coalescent models, PC maps displayed sinusoidal curves as observed in 269 previous studies (Novembre and Stephens, 2008). We observed that sPCA, which includes 270 several distance matrices within Moran eigenvector maps of genetic data, produced results 271 similar to those of PCA, and did not correct for IBD effects. When we applied SFA to 272 spatial coalescent simulations, the algorithm clustered individuals in several small groups 273 depending on the number of latent factors used in the method. SFA factor maps actually 274 displayed outcomes closer to discrete clusters than to continuous variation. After adjusting 275

276 for the spatial scale in the covariance model, spFA was able to remove the oscillating277 shapes observed in the first PCs sequentially.

278 When PCA was applied to spatially explicit simulations of two diverging populations, PC maps failed to firmly identify genetic discontinuities between populations. Despite a 279 280 relatively long period of isolation in simulations, the populations were not strongly sep-281 arated in PC maps due to the horseshoe effect. Compared to PCA and sPCA, the spFA method had increased power to identify genetic discontinuities where they were masked 282 283 by spurious autocorrelation effects. When we applied SFA, we found that, up to normal-284 ization of outputs, the results were similar to those generated by clustering algorithms like 285 STRUCTURE. For simulations of two diverging populations, SFA detected a main separation 286 between two differentiated populations, but this approach did not correct for IBD effects 287 within the main genetic clusters. Similarly to STRUCTURE, the results of SFA were influ-288 enced by the presence of IBD patterns in the samples. We found that spFA alleviated this 289 issue, and that it produced results more robust to the choice of the number of factors than 290 SFA.

291 The methods used in this study provided quite distinct descriptions of the data when they were applied to human population samples from Central and East Asia, and they un-292 293 derlined several aspects of the data. With PCA, a typical horseshoe pattern was observed, but no obvious genetic discontinuities were observed. In contrast, SFA provided evidence 294 for two main clusters which were also confirmed by spFA. When we used SFA with K =295 296 3, we obtained shapes without natural interpretations (Figure 8). SFA detected additional 297 discontinuities whereas the other methods suggested that continuous genetic variation in 298 geographic space was predominant. We observed that SFA behaves like clustering algo-299 rithms and did not correct for spurious clusters created by IBD patterns. This issue makes

the SFA results difficult to interpret in terms of admixture and ancestral populations. The 300 spFA method corrected for the horseshoe pattern observed in PC plots by removing au-301 tocorrelation effects from the second and third axes. The method suggested that Asian 302 population structure is strongly influenced by IBD patterns. In the spFA plot, Hazara of 303 304 Pakistan and Uygur of northwestern China grouped together, and were placed between 305 Pakistani and East Asian populations (Rosenberg et al., 2002). A way to interpret those results is as a support for admixed genomes in Hazara and Uygur populations, or as fa-306 307 voring the hypothesis of a central Asian migration route of modern humans in East Asia 308 (Zhang et al., 2007). The public availability of data sets other than the HGDP will en-309 able us to make further assessment of the interest of the method for the analysis of human 310 genetic data in the future.

311 A potential limitation of the spFA approach is to be sensitive to the choice of the scale parameter,  $\theta$ . The  $\theta$  parameter actually determines the scale of the spatial effects that 312 could be removed by spFA. Note that spFA is essentially performing a standard principal 313 component analysis when it is applied with small values of the scale parameter. In this 314 315 study, we suggested to explore a grid of  $\theta$  values so that IBD effects could be removed at distinct scales sequentially. The choice of the number of factors, K, in spFA is also 316 tied to the particular value of  $\theta$  implemented in the model. One way to determine K is by 317 using Tracy-Widom tests on the matrix of genotypes,  $\tilde{G}$  (Patterson et al., 2006). Gradually 318 increasing the value of  $\theta$  enabled a fine grain analysis of genetic discontinuities in human 319 data, and allowed us to study IBD patterns within genetic clusters. The computational 320 321 complexity of spFA is linear in function of the number of markers. Since it is equivalent 322 to the computation of a low rank approximation of the genotypic matrix (lower than a standard PCA, a few seconds on standard computer systems), applying spFA at multiple 323

324 scales was not highly time-consuming.

325 **Conclusion.** This study provided a comparison of existing methods that attempt to correct for IBD effects in population genetic analyses, and showed that each of studied 326 327 approaches provided different insights on the data. Under equilibrium IBD, PCA was 328 confounded by continuous variation and main genetic discontinuities may be missed or 329 misinterpreted. For the same data, SFA over-estimated the number of clusters in the genetic data, creating spurious clusters from continuous patterns. In the presence of IBD 330 331 patterns, spatial factor analysis provided clearer interpretations of the data than PCA and SFA. In a spatially explicit framework, we found that spFA identified genetic discontinu-332 ities more efficiently than did PCA or SFA when these discontinuities are blurred by noise 333 334 from IBD patterns in the genetic data.

#### 335 Acknowledgments

We thank Nicolas Duforet-Frebourg for his help with the software ms. This work was
supported by a grant from la Région Rhône-Alpes to Eric Frichot and Olivier François,
and by an NSF grant to Sean Schoville (OISE-0965038). Olivier François acknowledges
support from Grenoble INP.

#### 340 Legends

- 341 Figure 1: Algorithm for SpFA. For a genotypic matrix G with individual geographic coor-
- 342 dinates ( $X_i$ ), and for scale parameter  $\theta > 0$ , the spFA steps summarize as follows:

343

344 Figure 2: PC and SpFA factor maps for data simulated under an IBD model. A) PC 345 maps, B) SpFA factor maps for  $\theta/\bar{d} = 0.1$ , C) SpFA factor maps for  $\theta/\bar{d} = 0.2$ , SpFA 346 factor maps for  $\theta/\bar{d} = 0.3$ .

347

348 Figure 3: SFA factor maps for data simulated under an IBD model. Plots of the first349 three Factor maps for SFA.

350

351 Figure 4: Two discrete populations under equilibrium IBD. Plots of the first 2 maps for A)
352 PCA, B) sPCA, C) spFA, D) SFA.

353

354 Figure 5: Wilks'  $\Lambda$  statistic as a function of the scale parameter  $\theta/\bar{d}$  in spFA.

355

356 Figure 6: Wilks'  $\Lambda$  statistic as a function of the divergence time,  $\tau$ , ranging between 1 357 and 100.

358

359 Figure 7: A) Asia map with geographic locations of HGDP populations. PC and fac-360 tor plots for B) PCA C) SpFA, D) SFA.

361

362 Figure 8: Factor plots for A), B) SFA and C), D) SpFA with K = 3 clusters.

## 363 **References**

- 364 Ahmed, N., Natarajan, T., and Rao, K. R. (1974). Discrete cosine transfom. IEEE Trans-
- 365 *actions on Computers*, C-23(1):90–93.
- 366 Borcard, D. and Legendre, P. (2002). All-scale spatial analysis of ecological data by means
- 367 of principal coordinates of neighbour matrices. *Ecological Modelling*, 153(1-2):51–68.

- Borcard, D., Legendre, P., Avois-Jacquet, C., and Tuomisto, H. (2004). Dissecting the
  spatial structure of ecological data at multiple scales. *Ecology*, 85(7):1826–1832.
- 370 Cressie, N. A. C. (1993). Statistics for spatial data. Revised ed. Wiley, New York.
- 371 Diaconis, P., Goel, S., and Holmes, S. (2008). Horseshoes in multidimensional scaling
- and local kernel methods. *The Annals of Applied Statistics*, 2(3):777–807.
- 373 Dray, S., Legendre, P., and Peresneto, P. (2006). Spatial modelling: a comprehensive
- 374 framework for principal coordinate analysis of neighbour matrices (PCNM). *Ecological*
- 375 *Modelling*, 196(3-4):483–493.
- 376 Durand, E., Jay, F., Gaggiotti, O. E., and François, O. (2009). Spatial inference of ad-
- mixture proportions and secondary contact zones. *Molecular Biology and Evolution*,
  26:1963–1973.
- Engelhardt, B. E. and Stephens, M. (2010). Analysis of population structure: A unifying
  framework and novel methods based on sparse factor analysis. *PLoS Genetics*, 6:12.
- 381 François, O., Currat, M., Ray, N., Han, E., Excoffier, L., and Novembre, J. (2010). Prin-
- 382 cipal component analysis under population genetic models of range expansion and ad-
- mixture. *Molecular Biology and Evolution*, 27(6):1257–68.
- François, O. and Durand, E. (2010). Spatially explicit bayesian clustering models in population genetics. *Molecular Ecology Resources*, 10:773–784.
- 386 Hudson, R. R. (2002). Generating samples under a Wright-Fisher neutral model of genetic
- 387 variation. *Bioinformatics*, 18(2):337–338.
- 388 Jolliffe, I. T. (1986). Principal Component Analysis. Springer Verlag New York.

- 389 Jombart, T., Devillard, S., Dufour, A.-B., and Pontier, D. (2008). Revealing cryptic spatial
- patterns in genetic variability by a new multivariate method. *Heredity*, 101(1):92–103.
- 391 Kimura, M. and Weiss, G. H. (1964). The stepping stone model of population structure
- and the decrease of genetic correlation with distance. *Genetics*, 49(4):561–576.
- 393 Legendre, P. and Gallagher, E. (2001). Ecologically meaningful transformations for ordi-
- nation of species data. *Oecologia*, 129(2):271–280.
- 395 Malécot, G. (1948). Les Mathématiques de l'Hérédité. Masson, Paris.
- 396 Mardia, K. V., Kent, J. T., and Bibby, J. M. (1979). *Multivariate Analysis*. Academic
  397 Press, London.
- McVean, G. (2009). A genealogical interpretation of principal components analysis. *PLoS Genetics*, 5(10):10.
- 400 Meirmans, P. G. (2012). The trouble with isolation by distance. *Molecular Ecology*,
  401 21(12):2839–46.
- 402 Novembre, J. and Stephens, M. (2008). Interpreting principal component analyses of
  403 spatial population genetic variation. *Nature Genetics*, 40(5):646–649.
- 404 Patterson, N., Price, A. L., and Reich, D. (2006). Population structure and eigenanalysis.
  405 *PLoS Genetics*, 2:20.
- 406 Patterson, N. J., Moorjani, P., Luo, Y., Mallick, S., Rohland, N., Zhan, Y., Genschoreck,
- 407 T., Webster, T., and Reich, D. (2012). Ancient admixture in human history. *Genetics*,
- 408 *doi:10.1534/genetics.112.145037*.

- 409 Powell, G. T., Yang, H., Tyler-Smith, C., and Xue, Y. (2007). The population history of
- 410 the Xibe in northern China: a comparison of autosomal, mtDNA and Y-chromosomal
- 411 analyses of migration and gene flow. *Forensic Science International Genetics*, 1(2):115–
- 412 119.
- 413 Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure
- 414 using multilocus genotype data. *Genetics*, 155(2):945–959.
- 415 Rosenberg, N., Pritchard, J., Weber, J., Cann, H., Kidd, K., Zhivotovsky, L., and Feldman,
- 416 M. (2002). Genetic structure of human populations. *Science*, pages 2381–2385.
- 417 Rousset, F. (1997). Genetic differentiation and estimation of gene flow from F-statistics
  418 under isolation by distance. *Genetics*, 145(4):1219–1228.
- 419 Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis.
  420 *BMC Bioinformatics*, 61(3):611–622.
- 421 Wright, S. (1943). Isolation by distance. *Genetics*, 28(2):114–138.
- 422 Zhang, F., Su, B., Zhang, Y., and Jin, L. (2007). Genetic studies of human diversity in East
- 423 Asia. Philosophical Transactions of the Royal Society of London Series B: Biological
- 424 *Sciences*, 362(1482):987–996.

#### 425 Legends

- 426 Figure 1: Algorithm for SpFA. For a genotypic matrix G with individual geographic coor-
- 427 dinates ( $X_i$ ), and for scale parameter  $\theta > 0$ , the spFA steps summarize as follows:
- 428
- 429 Figure 2: PC and SpFA factor maps for data simulated under an IBD model. A) PC

maps, B) SpFA factor maps for  $\theta/\bar{d} = 0.1$ , C) SpFA factor maps for  $\theta/\bar{d} = 0.2$ , SpFA 430 factor maps for  $\theta/\bar{d} = 0.3$ . 431

432

Figure 3: SFA factor maps for data simulated under an IBD model. Plots of the first 433 three Factor maps for SFA. 434

435

437

Figure 4: Two discrete populations under equilibrium IBD. Plots of the first 2 maps for A) 436 PCA, B) sPCA, C) spFA, D) SFA.

438

Figure 5: Wilks' A statistic as a function of the scale parameter  $\theta/\bar{d}$  in spFA. 439

440

441 Figure 6: Wilks'  $\Lambda$  statistic as a function of the divergence time,  $\tau$ , ranging between 1 442 and 100.

443

Figure 7: A) Asia map with geographic locations of HGDP populations. PC and fac-444 tor plots for B) PCA C) SpFA, D) SFA. 445

446

447 Figure 8: Factor plots for A), B) SFA and C), D) SpFA with K = 3 clusters.



























