

1 **An Eigenvalue Test for spatial Principal Component Analysis**

2 Montano V^{1*} and Jombart T²

3 ¹ School of Biology, University of St Andrews, Bute Building, St Andrews KY16 9TS, UK

4 ² MRC Centre for Outbreak Analysis and Modelling, Department of Infectious Disease

5 Epidemiology, Imperial College, St Mary's Campus, Norfolk Place, London W2 1PG, UK

6 *Corresponding author: mirainoshojo@gmail.com

7 **Abstract**

8 **Background**

9 The spatial Principal Component Analysis (sPCA, Jombart 2008) is designed to investigate
10 non-random spatial distributions of genetic variation. Unfortunately, the associated tests
11 used for assessing the existence of spatial patterns (*global and local test*; Jombart et al.
12 2008) lack statistical power and may fail to reveal existing spatial patterns. Here, we
13 present a non-parametric test for the significance of specific patterns recovered by sPCA.

14 **Results**

15 We compared the performance of this new test to the original *global* and *local* tests using
16 datasets simulated under classical population genetic models. Results show that our test
17 outperforms the original *global* and *local* tests, exhibiting improved statistical power while
18 retaining similar, and reliable type I errors. Moreover, by allowing to test various sets of
19 axes, it can be used to guide the selection of retained sPCA components.

20 **Conclusions**

21 As such, our test represents a valuable complement to the original analysis, and should
22 prove useful for the investigation of spatial genetic patterns.

23 **Keywords;** eigenvalues; sPCA; spatial genetic patterns; Monte-Carlo

24 INTRODUCTION

25 The principal component analysis (PCA; Pearson 1901; Hotelling 1933) is one of the most
26 common multivariate approaches in population genetics (Jombart et al 2009). Although
27 PCA is not explicitly accounting for spatial information, it has often been used for
28 investigating spatial genetic patterns (Novembre and Stephens 2008). As a complement to
29 PCA, the spatial principal component analysis (sPCA; Jombart et al. 2008) has been
30 introduced to explicitly include spatial information in the analysis of genetic variation, and
31 gain more power for investigating spatial genetic structures.

32

33 sPCA finds synthetic variables, the principal components (PCs), which maximise both the
34 genetic variance and the spatial autocorrelation as measured by Moran's I (Moran 1950).
35 As such, PCs can reveal two types of patterns: '*global*' structures, which correspond to
36 positive autocorrelation typically observed in the presence of patches or clines, and '*local*'
37 structures, which correspond to negative autocorrelation, whereby neighboring individuals
38 are more genetically distinct than expected at random (for a more detailed explanation on
39 the meaning of *global* and *local* structures see Jombart et al.. 2008). The *global* and *local*
40 tests have been developed for detecting the presence of global and local patterns,
41 respectively (Jombart et al. 2008). Unfortunately, while these tests have robust type I error,
42 they also typically lack power, and can therefore fail to identify existing spatial genetic
43 patterns (Jombart et al.. 2008). Moreover, they can only be used to diagnose the presence
44 or absence of spatial patterns, and are unable to test the significance of specific structures
45 revealed by sPCA axes.

46

47 In this paper, we introduce an alternative statistical test which addresses these issues.

48 This approach relies on computing the cumulative sum of a defined set of sPCA

49 eigenvalues as a test statistic, and uses a Monte-Carlo procedure to generate null
50 distributions of the test statistics and approximate p-values. After describing our approach,
51 we compare its performances to the global and local tests using simulated datasets,
52 investigating several standard spatial population genetics models. Our approach is
53 implemented as the function *sPCA_randtest* in the package *adegenet* (Jombart 2008;
54 Jombart and Ahmed 2011) for the R software (R Core Team 2017).

55 **METHODS**

56 ***Test statistic***

57 As in most multivariate analyses of genetic markers, our approach analyses a table of
58 centred allele frequencies (i.e. set to a mean frequency of zero), in which rows represent
59 individuals or populations, and columns correspond to alleles of various loci (Jombart et al
60 2008; Jombart et al 2009; Jombart et al 2010). We note X the resulting matrix, and n the
61 number of individuals analysed. In addition, the sPCA introduces spatial data in the form of
62 a n by n matrix of spatial weights L , in which the i^{th} row contains weights reflecting the
63 spatial proximity of all individuals to individual i . The PCs of sPCA are then found by the
64 eigen-analysis of the symmetric matrix (Jombart et al. 2008):

$$65 \quad 1/(2n) X^T(L^T + L)X \quad (1)$$

66 We note λ the corresponding non-zero eigenvalues. We differentiate the r positive
67 eigenvalues λ^+ , corresponding to global structures, and the 's' negative eigenvalues λ^- ,
68 corresponding to local structures, so that $\lambda = \{\lambda^+, \lambda^-\}$. Without loss of generality, we
69 assume both sets of eigenvalues are ordered by decreasing absolute value, so that $\lambda_1^+ >$
70 $\lambda_2^+ > \dots > \lambda_r^+$ and $|\lambda_1^-| > |\lambda_2^-| > \dots > |\lambda_s^-|$. Simply put, each eigenvalue quantifies the
71 magnitude of the spatial genetic patterns in the corresponding PC: larger absolute values
72 indicate stronger global (respectively local) structures. We note $V^+ = \{v_1^+, \dots, v_r^+\}$ and $V^- =$
73 $\{v_1^-, \dots, v_s^-\}$ the sets of corresponding PCs. The most natural choice of test statistic to
74 assess whether a given PC contains significant structure would seem to be the
75 corresponding eigenvalue. This would, however, not account for the dependence on
76 previous PCs: v_j^+ (respectively v_j^-) can only be significant if all previous PCs $\{v_1^+, \dots, v_{j-1}^+\}$
77 are also significant. To account for this, we define the test statistic for v_j^+ as:

$$78 \quad f_j^+ = \sum_{i=1, \dots, j} \lambda_i^+$$

79 and as:

80
$$f_i^- = \sum_{i=1, \dots, j} |\lambda_i^-|$$

81 for v_j^- .

82

83 ***Permutation procedure***

84 f_i^+ and f_i^- become larger in the presence of strong global or local structures in the first i^{th}

85 global / local PCs. Therefore, they can be used as test statistics against the null

86 hypotheses of absence of global or local structures in these PCs. The expected

87 distribution of f_i^+ and f_i^- in the absence of spatial structure is not known analytically.

88 Fortunately, it can be approximated using a Monte-Carlo procedure, in which at each

89 permutation individual genotypes are shuffled to be assigned to a different pair of

90 coordinates than in the observed original dataset and f_i^+ and f_i^- are computed. Note that the

91 original values of the test statistic are also included in these distributions, as the initial

92 spatial configuration is by definition a possible random outcome. The p -values are then

93 computed as the relative frequencies of permuted statistics equal to or greater than the

94 initial value of f_i^+ or f_i^- .

95

96 To guide the selection of global and local PCs to retain, the simulated values of each

97 eigenvalue (from most positive to most negative), which make up the f_i^+ and f_i^- statistics,

98 are also recorded during the permutation procedure. In this way, if global or local

99 structures are detected to be significant, an observed p -value for each observed

100 eigenvalue can be estimated by comparison with its simulated eigenvalue distribution.

101 Note that the number of eigenvalues produced by an sPCA does not change between the

102 observed and permuted datasets, so each observed eigenvalue can be compared with

103 the distribution of the corresponding simulated one. This testing procedure can be used

104 with increasing numbers of retained axes. Because each test is conditional on the previous

105 tests, incremental Bonferroni correction is used to avoid the inflation of type I error, so that
106 the significance level for the i^{th} PC will be α / i , where α is the target type I error. Hence, the
107 correction implies that if the most positive (or negative) eigenvalue is significant in regards
108 with the chosen p -value threshold, the second eigenvalue is tested for a p -value threshold
109 that is the half of the previous and so on. The entire testing procedure is implemented in
110 the function *sPCA_randtest* in the package *adeigenet* (Jombart 2008; Jombart and Ahmed
111 2011) for R (R Core Team 2017). A flow chart of the test procedure is shown in Figure 1.

112

113 **Simulation study**

114 To assess the performance of our test, we simulated genetic data under three migration
115 models: island (IS) and stepping stone (SS), using the software GenomePop 2.7 (Carvajal-
116 Rodríguez 2008), and isolation by distance (IBD), using *IBDSimV2.0* (Leblois 2009). We
117 simulated the IS and SS models with 4 populations, each with 25 individuals, and a single
118 population under IBD with 100 individuals. 200 unlinked biallelic diploid loci (or single
119 nucleotide polymorphisms; SNPs) were simulated. Populations evolved under constant
120 effective population size $\theta = 20$, and interchanged migrants at three different symmetric
121 and homogeneous rates (0.005, 0.01, and 0.1). We performed 100 independent runs for
122 each of the three migration rates, for a total of 300 simulated dataset per migration model.

123

124 To quantify type I error rates for the *sPCA_randtest*, *global* and *local tests*, we extracted
125 100 random coordinates from 10 square 2D grids, using the function *spsample* from the
126 *spdep* package (Bivand et al. 2013). In order to evaluate the rate of false negatives for
127 global patterns, we manually generated 10 sets of 100 pairs of coordinates simulating
128 gradients and/or patches from 2D grids. An example of simulated global patterns is
129 presented in Figure 2. To test for the rate of false negatives for local patterns, we perform

130 a principal component analysis on 10 random datasets simulated under the SS model with
131 0.005 migration rate. We used the coordinates of the individuals on the first principal
132 component and set the second coordinate to zero for all individuals (1D). With the
133 coordinates so produced, we used the function *chooseCN* in *adegenet* to obtain 10
134 neighbouring graphs where the most genetically distinct individuals (falling in the upper
135 quartile of the pairwise genetic distances) are considered as neighbors, while the others
136 are non-neighbors.

137

138 We tested 100 simulations each for all the 30 sets of geographic coordinates (random,
139 positive and negative), for each of the three migration rates (0.005, 0.01 and 0.1), for each
140 of the three migration models (IS, SS, IBD; total of 9,000 tests per migration model). We
141 repeated all tests using a subset of 40 SNPs per individual, for a total of 18,000 tests in the
142 absence of spatial structures, and 36,000 tests in the presence of global or local
143 structures.

144 **RESULTS**

145 ***Statistical power of the *spca_randtest****

146 We compared the performances of the *spca_randtest* with the *global* and *local* tests in
147 three settings: in the absence of spatial structure, and in the presence of global, and local
148 structures. The results obtained in the absence of spatial structure show that all tests have
149 reliable type I errors (Table 1 and 2). The *spca_randtest* exhibited consistently better
150 performances for detecting existing structures in the data than both *global* and *local tests*
151 (Table 1 and 2). Although our simulated local spatial patterns turned out more difficult to
152 detect than global patterns, the *spca_randtest* is twice to five times more effective than the
153 *local test* (Table 1 and 2). Generally, the underlying migration model, the migration rate
154 and the number of loci affect the ability of all tests to detect non-random spatial patterns.
155 Both *spca_randtest* and *global* and *local tests* have in fact a lower sensitivity in presence
156 of island migratory schemes, while results for stepping stone and isolation by distance
157 models are more satisfying (Table 1 and 2). Increasing migration rates lead to a higher
158 rates of false negatives for all tests, which can be overcome using more loci (Table 1 and
159 2).

160

161 Significant eigenvalues are assessed using a hierarchical Bonferroni correction which
162 accounts for non-independence of eigenvalues and multiple testing (Figure 2). Strong
163 patterns (e.g. IBD) tend to produce a higher number of significant components than weak
164 patterns (e.g. island models with high migration rates), which are otherwise captured by
165 fewer to no components.

166

167 ***Application to real data***

168 We have run the sPCA to compare the new *spca_randtest* and previous tests to a real

169 dataset of human mitochondrial DNA (mtDNA). We used a dataset of 85 populations from
170 Central-Western Africa that spans a big portion of the African continent (from Gabon to
171 Senegal; Montano et al 2013). Previous analysis on these data detected a clear genetic
172 structure from West to Central Africa with ongoing stepping stone migration movements.
173 We therefore expected that this spatial distribution of genetic variation would be detected
174 as significant. In the sPCA, populations were treated as units of the analysis, for which
175 allele frequencies of mtDNA polymorphisms are calculated per population. The same
176 approach was used in Montano et al 2013 to run a discriminant analysis of principal
177 components (DAPC; Jombart et al 2010) and detect population genetic structure. The
178 sPCA analysis is found non significant by *global* and *local* tests after 1e4 permutations (*p*-
179 value > 0.5), while the *spca_randtest* detects a significant global pattern already with 500
180 permutations, and with 1e4 permutations the *p*-value for global patterns is 0.005. The
181 second step of the test on single eigenvalues finds the three most positive components to
182 be significant after Bonferroni correction (Table 3). Significant axes can thus be plotted
183 against the spatial network to give a biological interpretation to the results (Figure 3).

184 **DISCUSSION**

185 We introduced a new statistical test associated to the sPCA to evaluate the statistical
186 significance of global and local spatial patterns. Using simulated data, we show that this
187 new approach outperforms previously implemented tests, having greater statistical power
188 (lower type II errors) whilst retaining consistent type I errors. Our simulations also suggest
189 that demographic settings and migratory models can substantially impact the ability to
190 detect spatial patterns. Indeed, high migration rates, non-hierarchical migration models,
191 such as island model, and low amount of loci can hamper or worsen the performance of
192 the test, preventing the detection of actual spatial patterns. In lack of previous information
193 on the demographic history and/or the movement ecology of the population under study, it
194 is certainly useful to exploit all the available genetic information. In this regards, our
195 simulations show how an increased number of loci does improve the ability of the test to
196 provide meaningful results.

197

198 The impact of specific factors such as the effective population size or the number of
199 individuals sampled per population remain to be investigated. A more extensive simulation
200 study, possibly comparing different non-model based methods such as sPCA, would clarify
201 the extent of the spatial information that can be obtained with such methods without
202 comparing explicit evolutionary hypotheses. In fact, the sPCA and the associated
203 *sPCA_randtest* cannot distinguish between explicit migration models. However, the
204 possibility to detect which eigenvalues contain the spatial information provides the user
205 with further information to interpret the biological meaning of the spatial structure, by
206 focusing on few meaningful dimensions.

207

208 Our data application seems to confirm that the *sPCA_randtest* is more effective than *global*

209 or *local* tests. We chose indeed a previously published dataset of human populations
210 which span a subcontinental area of Africa and had been originally detected to be a highly
211 structured dataset with a geographic cline of population differentiation (Montano et al
212 2013). On the basis of the original results, we would have expected a spatial global
213 structure to be present in the data and thus detected with an sPCA. While the *global* test
214 failed to provide statistical significance, the *spca_randtest* did obtain significant results and
215 pointed to the three first most positive components to be also significant after Bonferroni
216 correction. In agreement with the original interpretation of the genetic structure within the
217 samples, spatial component 1 (SP1) shows a clear differentiation of populations in the
218 Gabon-Congo region, while SP2 detects differentiation of Central Nigerian and North
219 Cameroonian populations, on one hand, and extreme Western populations of Senegal, on
220 the other hand (Figure 3). The colored combination of the first and second most positive
221 component (Figure 3) also correctly detects a more fragmented differentiation across
222 Central forested areas (Cameroon, Gabon and Congo) compared to more homogeneous
223 Central-Western populations, which was the main result of the original publication based
224 on very different approaches (Montano et al 2013). We limited the analysis to these two
225 component as the third did not add much information to the previous.

226

227 Our simulation approach coupled with a real data application well illustrates the
228 informativeness of our new test to retrieve significant spatial patterns, being these global
229 or local structures and highlights the usefulness of selecting a specific number of
230 significant components to interpret the biological meaning of the results.

231 **Declarations**

232 **Data Accessibility**

233 https://github.com/thibautjombart/adegetnet/blob/master/R/spca_randtest.R

234 **Acknowledgements**

235 The authors declare no conflict of interest

236 **Author contributions**

237 Test development: VM and TJ. Data analysis: VM. Wrote the manuscript: VM and TJ.

238 **Literature**

- 239 1. Bivand RS, Pebesma E, Gómez-Rubio V (2013) *Applied Spatial Data Analysis with*
240 *R*. Springer, New York, 378pp.
- 241 2. Balkenhol N, Gugerli F, Cushman SA, Waits LP, Coulon A, Arntzen JW, Holderegger
242 R, Wagner HH (2009) Identifying future research needs in landscape genetics:
243 where to from here? *Landscape Ecology*, **24**, 455.
- 244 3. Carvajal-Rodríguez A (2008) GENOMEPOP: A program to simulate genomes in
245 populations. *BMC Bioinformatics*, **9**, 223.
- 246 4. Cushman SA, Landguth EL (2010) Spurious correlations and inference in landscape
247 genetics. *Molecular Ecology*, **19**, 3592–3602.
- 248 5. Hotelling H (1933). Analysis of a complex of statistical variables into principal
249 components. *Journal of educational psychology* **24**, 417.
- 250 6. Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic
251 markers. *Bioinformatics*, **24**, 1403–1405.
- 252 7. Jombart T, Devillard S, Dufour AB, Pontier D (2008) Revealing cryptic spatial
253 patterns in genetic variability by a new multivariate method. *Heredity*, **101**, 92–103.
- 254 8. Jombart T, Pontier D, Dufour AB (2009). Genetic markers in the playground of
255 multivariate analysis. *Heredity* **102**, 330–341.
- 256 9. Jombart T, Devillard S, Balloux F (2010). Discriminant analysis of principal
257 components: a new method for the analysis of genetically structured populations.
258 *BMC genetics* **11**, 94.
- 259 10. Jombart T, Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-
260 wide SNP data. *Bioinformatics* **27**, 3070–3071.
- 261 11. Montano V, Marcari V, Pavanello M, Anyaele O, Comas D, Destro-Bisol G, Batini C
262 (2013) The influence of habitats on female mobility in Central and Western Africa

- 263 inferred from human mitochondrial variation. *BMC evolutionary biology* **13**, 1:24.
- 264 12. Moran PAP (1950) Notes on Continuous Stochastic Phenomena. *Biometrika*, **37**,
- 265 17–23.
- 266 13. Novembre J, Stephens M (2008) Interpreting principal component analyses of
- 267 spatial population genetic variation. *Nature Genetics*, **40**, 646–649.
- 268 14. Pearson K (1901) On lines and planes of closest fit to systems of points in space.
- 269 *Philosophical Magazine Series 6*, **2**, 559–572.
- 270 15. Peres-Neto PR, Jackson DA, Somers KM (2005) How many principal components?
- 271 stopping rules for determining the number of non-trivial axes revisited.
- 272 *Computational Statistics & Data Analysis*, **49**, 974 – 997.
- 273 16. R Core Team (2017) R: A language and environment for statistical computing. *R*
- 274 *Foundation for Statistical Computing*, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)
- 275 [project.org/](https://www.R-project.org/).
- 276 17. Schiffers K, Travis JMJ (2014) ALADYN - a spatially explicit, allelic model for
- 277 simulating adaptive dynamics. *Ecography*, **37**, 1288–1291.

278 **Legends**

279

280 **Figure 1.** Flow chart illustrating the steps of the *sPCA_randtest*. The first step on the top
281 panel assess the statistical significance of global either local patterns. If at least one of the
282 two is significant, the second step of the test exploits the eigenvalue distribution recorded
283 over the permutations to obtain an empirical p -value for each eigenvalue, starting from the
284 most positive (or most negative). As the first eigenvalue is significant in comparison with a
285 chosen threshold, the following is tested and compared to a more stringent threshold
286 (Bonferroni correction) until a non-significant eigenvalue is found and the routine stops.

287

288 **Figure 2.** Graphical representation of island and stepping stone migration models (IS and
289 SS) in the panel above. Black rows represent the presence and direction of migration rates
290 among populations (purple circles). The panel below represents two examples of
291 simulated global patterns, where a set of 100 pairs of coordinates are picked from a set of
292 1000 random pairs of coordinates built in 2D squares at different scales (in the example
293 here reported the scales are 1:1e4 and 1:1e5, respectively). Every 25 pairs of coordinates
294 are assigned to a different simulated population, distinguished by red, blue, black and
295 yellow colors, in order to obtain spatially segregated populations. These simulated spatial
296 distributions are used to calculate the matrix L of spatial connection (see Figure S1).

297

298 **Figure 3.** Plot of the first and second most positive observed eigenvalues of the mtDNA
299 dataset here analysed. The background map represents the countries from where the
300 populations included into the original study were sampled (from West to East: Senegal,
301 Guinea-Bissau, Guinea, Sierra Leone, Liberia, Ivory Coast, Ghana, Togo, Benin, Nigeria,
302 Cameroon, Equatorial Guinea, Gabon, Congo). sPC1 and sPC2 are represented

303 independently using a square size proportional to the value of each population along the
304 first and second component, respectively. Whites squares show negative values and black
305 squares the positive values, with size being proportional to the absolute value of the
306 coordinate. sPC1-sPC2 is a summarized representation of the values along the first and
307 second component assumed by each population, using a color gradient.

308

309 **Figure S1.** Distributions of significant eigenvalues detected in the presence of global (blue
310 bars) and local (green bars) spatial patterns after hierarchical Bonferroni correction, for
311 100 significantly positive and 100 significantly negative patterns. Black bars correspond to
312 eigenvalues which are significant without Bonferroni correction. Bars' height indicates the
313 frequency of observing a significant eigenvalue in a certain position (from most positive to
314 most negative) over the 100 tested patterns.

315 **Table 1.** Significant results for *global test* (g test), *local tests* (l test), and *spca_randtest* (r test +/-) for random, global and local patterns
316 using 200 loci per individual. IS, SS, IBD indicate the migration models (see Methods); different migration rates are coded by number: 1 =
317 0.005, 2 = 0.01 and 3 = 0.1. Results show the proportion of significant tests over 1,000 replicates, based on 1,000 permutations with
318 thresholds .05 and .01.

200 SNPs		Random Patterns				Global Patterns				Local Patterns			
Models	Significance level	g test	r test (+)	l test	r test (-)	g test	r test (+)	l test	r test (-)	g test	r test (+)	l test	r test (-)
IS-1	.05	0.054	0.059	0.041	0.047	0.947	0.985	0.029	0.001	0.047	0.071	0.061	0.284
	.01	0.011	0.007	0.009	0.010	0.822	0.948	0.005	0.001	0.008	0.010	0.015	0.113
IS-2	.05	0.040	0.041	0.058	0.056	0.227	0.564	0.044	0.018	0.056	0.059	0.050	0.123
	.01	0.007	0.009	0.009	0.013	0.067	0.302	0.005	0.002	0.011	0.007	0.012	0.026
IS-3	.05	0.051	0.040	0.053	0.041	0.055	0.049	0.045	0.047	0.049	0.047	0.044	0.059
	.01	0.010	0.014	0.013	0.008	0.010	0.013	0.007	0.013	0.002	0.014	0.008	0.019
SS-1	.05	0.053	0.058	0.053	0.050	0.986	0.996	0.022	0.000	0.063	0.064	0.124	0.582
	.01	0.007	0.011	0.010	0.010	0.960	0.988	0.002	0.000	0.017	0.010	0.041	0.398
SS-2	.05	0.044	0.058	0.058	0.063	0.798	0.909	0.047	0.004	0.034	0.044	0.059	0.316
	.01	0.011	0.011	0.013	0.016	0.676	0.771	0.010	0.000	0.004	0.005	0.014	0.147
SS-3	.05	0.047	0.046	0.057	0.049	0.054	0.128	0.040	0.042	0.044	0.054	0.049	0.071
	.01	0.014	0.007	0.011	0.013	0.014	0.036	0.006	0.010	0.003	0.009	0.006	0.009
IBD-1	.05	0.044	0.050	0.053	0.048	0.962	0.999	0.021	0.000	0.025	0.087	0.438	0.809
	.01	0.008	0.012	0.009	0.010	0.926	0.997	0.003	0.000	0.009	0.023	0.192	0.694
IBD-2	.05	0.052	0.045	0.061	0.038	0.967	0.998	0.023	0.000	0.046	0.076	0.451	0.794
	.01	0.009	0.008	0.011	0.009	0.932	0.997	0.004	0.000	0.009	0.018	0.208	0.672
IBD-3	.05	0.052	0.046	0.053	0.050	0.977	0.999	0.015	0.000	0.050	0.083	0.441	0.824

.01	0.013	<i>0.009</i>	0.011	0.012	0.939	0.999	<i>0.005</i>	<i>0.000</i>	<i>0.009</i>	0.023	0.225	0.684
-----	-------	--------------	-------	-------	--------------	--------------	--------------	--------------	--------------	-------	--------------	--------------

319 **p-values* are in italic when non significant and in bold when the fraction of true positive is above 20%

320 **Table 2.** Results for the same simulations reported in Table 1 using a subset of 40 loci per individual.

40 SNPs		Random Patterns				Global Patterns				Local Patterns			
Models	Significance level	g test	r test (+)	l test	r test (-)	g test	r test (+)	l test	r test (-)	g test	r test (+)	l test	r test (-)
IS-1	.05	0.052	0.061	0.046	0.050	0.591	0.807	0.033	0.004	0.036	0.000	0.055	0.077
	.01	0.016	0.013	0.010	0.007	0.393	0.592	0.005	0.000	0.004	0.000	0.015	0.022
IS-2	.05	0.053	0.047	0.038	0.042	0.103	0.226	0.046	0.020	0.073	0.000	0.057	0.038
	.01	0.011	0.009	0.006	0.006	0.022	0.072	0.011	0.005	0.012	0.000	0.010	0.006
IS-3	.05	0.047	0.050	0.050	0.045	0.048	0.060	0.044	0.042	0.036	0.000	0.053	0.026
	.01	0.009	0.011	0.008	0.007	0.009	0.011	0.011	0.011	0.002	0.000	0.013	0.001
SS-1	.05	0.052	0.054	0.039	0.049	0.898	0.949	0.017	0.000	0.050	0.001	0.067	0.169
	.01	0.009	0.012	0.005	0.011	0.826	0.865	0.006	0.000	0.007	0.000	0.021	0.052
SS-2	.05	0.046	0.045	0.050	0.046	0.528	0.588	0.044	0.009	0.052	0.000	0.048	0.081
	.01	0.013	0.010	0.010	0.015	0.377	0.370	0.016	0.000	0.005	0.000	0.011	0.014
SS-3	.05	0.068	0.040	0.050	0.048	0.066	0.055	0.053	0.033	0.026	0.000	0.047	0.023
	.01	0.014	0.005	0.013	0.012	0.012	0.009	0.005	0.006	0.006	0.000	0.008	0.000
IBD-1	.05	0.049	0.053	0.052	0.057	0.822	0.883	0.027	0.002	0.034	0.055	0.124	0.480
	.01	0.005	0.008	0.013	0.013	0.755	0.742	0.004	0.000	0.005	0.008	0.032	0.278
IBD-2	.05	0.043	0.054	0.060	0.049	0.835	0.880	0.028	0.001	0.043	0.051	0.111	0.458
	.01	0.011	0.007	0.015	0.009	0.755	0.732	0.005	0.000	0.008	0.015	0.026	0.259
IBD-3	.05	0.043	0.042	0.051	0.050	0.844	0.899	0.026	0.002	0.048	0.058	0.115	0.465
	.01	0.012	0.013	0.012	0.010	0.763	0.756	0.007	0.000	0.009	0.010	0.023	0.263

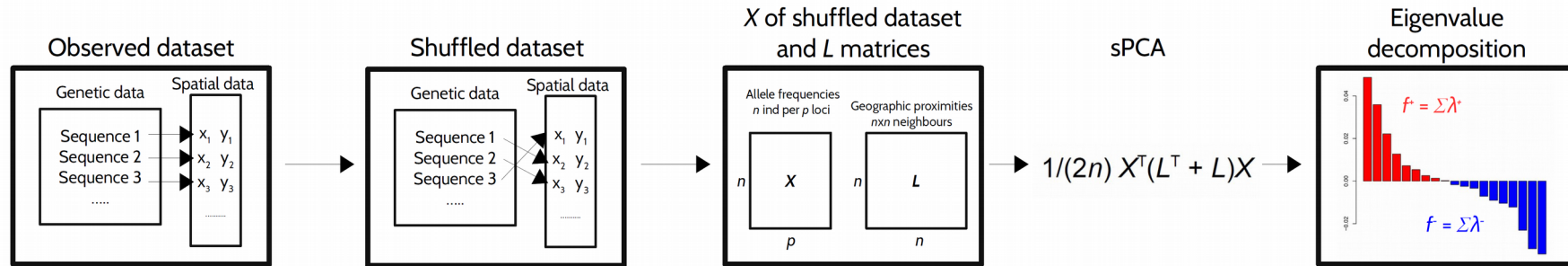
321 **p-values* are in italic when non significant and in bold when the fraction of true positive is above 20%

322 **Table 3.** Results of the *spca_randtest* with 1e4 permutations on the human mtDNA dataset (Montano et al, 2013). The simulated
 323 distribution of the f_i^+ and f_i^- statistics are compared to the f_i^+ and f_i^- statistics observed for the original dataset. A significant global pattern
 324 (or significant f_i^+ observed statistics) is found with the *spca_randtest* (p-value < 0.01). Thus, each eigenvalue is compared with its
 325 simulated distribution and assigned to be significant if its observed p -value is lower than the corrected Bonferroni p -value, with starting
 326 threshold of 0.05. Significant observed p -values as compared with Bonferroni corrected p -values are highlighted in bold.

Spatial patterns	Eigenvalue	Observed p-value	Bonferroni p-value	
Global pattern	0.0058	3.4e-2	0.0105	0.05
Local pattern	0.8826	8.5e-3	0.0137	0.025
		4.1e3	0.0136	0.016
		1.6e-3	0.506	0.

Flow chart of *sPCA_randtest*

Step 1. Detecting global or local spatial patterns

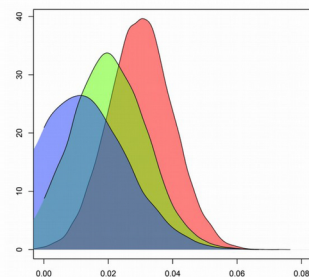


Permutation process is repeated x times to produce empirical distributions of f^* and f statistics

Step 2. Assessing statistical significance of single eigenvalues conditional on step 1

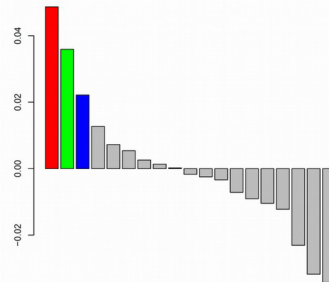
The value of single eigenvalues generated at each permutation is recorded separately, from most positive to most negative

Example of empirical eigenvalue distributions

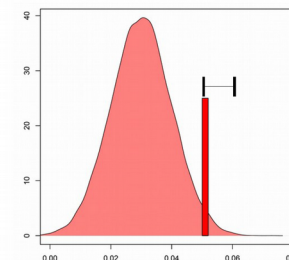


First, second and third most positive eigenvalue distributions

sPCA eigenvalues of observed dataset



Probability of observed first positive eigenvalue in respect with the simulated distribution of first positive eigenvalues



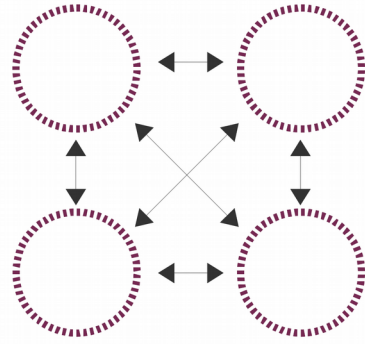
Fraction of simulated eigenvalues that are \geq observed first eigenvalue is used to calculate an empirical p-value

329 **Figure 2**

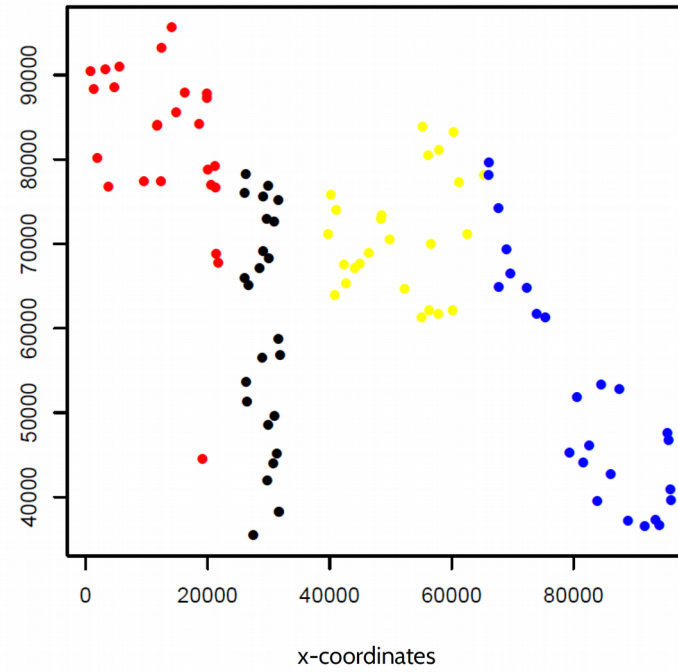
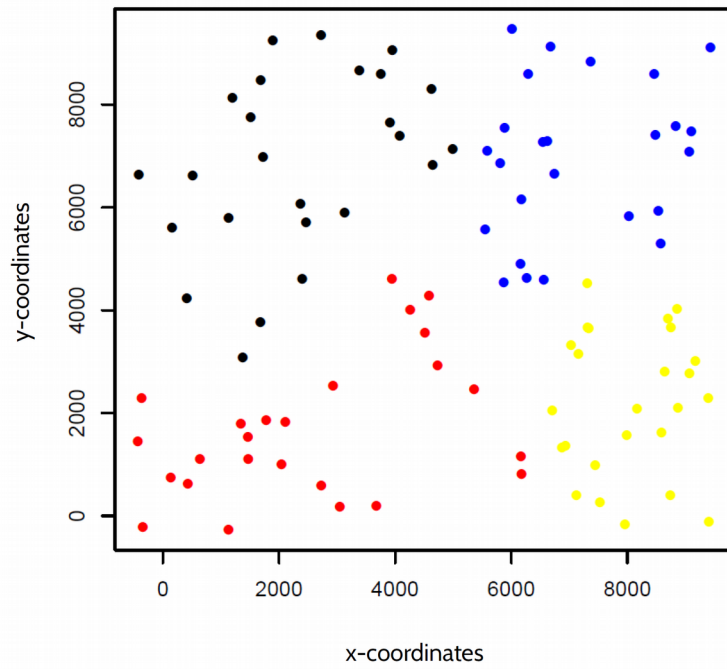
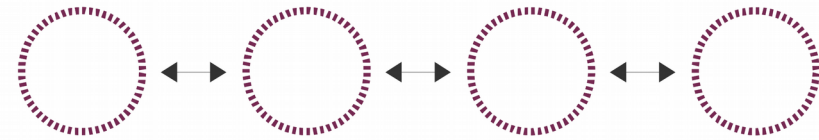
330

331

Island model



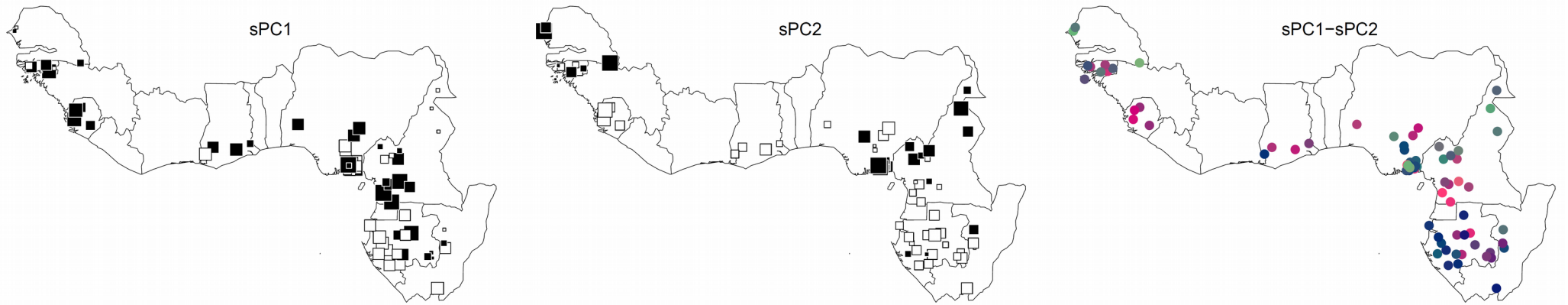
Stepping stone model



332 **Figure 3**

333

334



335 **Figure S1.**

336

