

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

2 **Running title:** A new statistical test for sPCA

3 **Word count:** 2982

4 Montano V^{1*} and Jombart T²

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

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

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

8 *Corresponding author: mirainoshojo@gmail.com

9 **Abstract**

- 10 • The spatial Principal Component Analysis (sPCA, Jombart 2008) is designed to
11 investigate non-random spatial distributions of genetic variation. Unfortunately, the
12 associated tests used for assessing the existence of spatial patterns (*global and*
13 *local test*; Jombart et al. 2008) lack statistical power and may fail to reveal existing
14 spatial patterns.
- 15 • Here, we present a non-parametric test for the significance of specific patterns
16 recovered by sPCA.
- 17 • We compared the performance of this new test to the original *global* and *local* tests
18 using datasets simulated under classical population genetic models. Results show
19 that our test outperforms the original *global* and *local* tests, exhibiting improved
20 statistical power while retaining similar, and reliable type I errors. Moreover, by
21 allowing to test various sets of axes, it can be used to guide the selection of
22 retained sPCA components. As such, it represents a valuable complement to the
23 original analysis, and should prove useful for the investigation of spatial genetic
24 patterns.

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

26 INTRODUCTION

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

34

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

47

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

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

50 eigenvalues as a test statistic, and uses a Monte-Carlo procedure to generate null

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

56 METHODS

57 *Test statistic*

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

$$66 \quad 1/(2n) X^T(L^T + L)X \quad (1)$$

67 We note λ the corresponding non-zero eigenvalues. We differentiate the r positive
68 eigenvalues λ^+ , corresponding to global structures, and the 's' negative eigenvalues λ^- ,
69 corresponding to local structures, so that $\lambda = \{\lambda^+, \lambda^-\}$. Without loss of generality, we
70 assume both sets of eigenvalues are ordered by decreasing absolute value, so that $\lambda_1^+ >$
71 $\lambda_2^+ > \dots > \lambda_r^+$ and $|\lambda_1^-| > |\lambda_2^-| > \dots > |\lambda_s^-|$.

72 Simply put, each eigenvalue quantifies the magnitude of the spatial genetic patterns in the
73 corresponding PC: larger absolute values indicate stronger global (respectively local)
74 structures. We note $V^+ = \{v_1^+, \dots, v_r^+\}$ and $V^- = \{v_1^-, \dots, v_s^-\}$ the sets of corresponding
75 PCs. The most natural choice of test statistic to assess whether a given PC contains
76 significant structure would seem to be the corresponding eigenvalue. This would, however,
77 not account the dependence on previous PCs: v_j^+ (respectively v_j^-) can only be significant if
78 all previous PCs $\{v_1^+, \dots, v_{j-1}^+\}$ are also significant. To account for this, we define the test
79 statistic for v_j^+ as:

$$80 \quad f_j^+ = \sum_{i=1, \dots, j} \lambda_i^+$$

81 and as:

$$82 \quad f_i^+ = \sum_{i=1, \dots, j} |\lambda_i|$$

83 for v_j^- .

84

85 ***Permutation procedure***

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

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

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

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

90 Fortunately, it can be approximated using a Monte-Carlo procedure, in which individual

91 genetic profiles are permuted randomly along the connection network, computing f_i^+ and f_i^-

92 for each permutation. Note that the original values of the test statistic are also included in

93 these distributions, as the initial spatial configuration is by definition a possible random

94 outcome. The p -values are then computed as the relative frequencies of permuted

95 statistics equal to or greater than the initial value of f_i^+ or f_i^- .

96

97 To guide the selection of global and local PCs to retain, this testing procedure can be used

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

99 tests, incremental Bonferroni correction is used to avoid the inflation of type I error, so that

100 the significance level for the i^{th} PC will be α / i , where α is the target type I error. The entire

101 testing procedure is implemented in the function `s_pca_randtest` in the package

102 *adeigenet* (Jombart 2008; Jombart and Ahmed 2011) for R (R Core Team 2016).

103 **Simulation study**

104 To assess the performance of our test, we simulated genetic data under three migration
105 models: island (IS) and stepping stone (SS), using the software GenomePop 2.7 (Carvajal-
106 Rodríguez 2008), and isolation by distance (IBD), using *IBDSimV2.0* (Leblois 2009). We
107 simulated the IS and SS models with 4 populations, each with 25 individuals, and a single
108 population under IBD with 100 individuals. 200 unlinked SNPs diploid loci were simulated.
109 Populations evolved under constant effective population size $\theta = 20$, and interchanged
110 migrants at three different symmetric and homogeneous rates (0.005, 0.01, and 0.1). We
111 performed 100 independent runs for each of the three migration rates, for a total of 300
112 simulated dataset per migration model.

113

114 To quantify rates of errors type I for the *spca_randtest*, *global* and *local tests*, we extracted
115 100 random coordinates from 10 square 2D grids, using the function *spsample* from the
116 *spdep* package (Bivand et al. 2013). In order to evaluate the rate of false negatives for
117 global patterns, we manually generated 10 sets of 100 pairs of coordinates simulating
118 gradients and/or patches from 2D grids. To test for the rate of false negatives for local
119 patterns, we perform a principal component analysis on 10 random datasets simulated
120 under the SS model with 0.005 migration rate. We used the coordinates of the individuals
121 on the first principal component and set the second coordinate to zero for all individuals
122 (1D). With the coordinates so produced, we used the function *chooseCN* in *adeigenet* to
123 obtain 10 neighbouring graphs where the most genetically distinct individuals (falling in the
124 upper quartile of the pairwise genetic distances) are considered as neighbors, while the
125 others are non-neighbors.

126 We tested 100 simulations each for all the 30 sets of geographic coordinates (random,
127 positive and negative), for each of the three migration rates (0.005, 0.01 and 0.1), for each

128 of the three migration models (IS, SS, IBD; total of 9,000 tests per migration model). We
129 repeated all tests using a subset of 40 SNPs per individual, for a total of 18,000 tests in the
130 absence of spatial structures, and 36,000 tests in the presence of global or local
131 structures.

132 **RESULTS**

133 ***Statistical power of the *spca_randtest****

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

148

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

154 **CONCLUSIONS**

155 We introduced a new statistical test associated to the sPCA to evaluate the statistical
156 significance of global and local spatial patterns. Using simulated data, we show that this
157 new approach outperforms previously implemented tests, having greater statistical power
158 (lower type II errors) whilst retaining consistent type I errors. Our simulations also suggest
159 that demographic settings and migratory models can substantially impact the ability to
160 detect spatial patterns. The impact of specific factors such as the effective population size
161 or the number of individuals sampled per population remain to be investigated.

162 **Acknowledgements**

163 The authors declare no conflict of interest

164 **Author contributions**

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

166 **Literature**

- 167 1. Bivand RS, Pebesma E, Gómez-Rubio V (2013) *Applied Spatial Data Analysis with*
168 *R*. Springer, New York, 378pp.
- 169 2. Balkenhol N, Gugerli F, Cushman SA, Waits LP, Coulon A, Arntzen JW, Holderegger
170 R, Wagner HH (2009) Identifying future research needs in landscape genetics:
171 where to from here? *Landscape Ecology*, **24**, 455.
- 172 3. Carvajal-Rodríguez A (2008) GENOMEPOP: A program to simulate genomes in
173 populations. *BMC Bioinformatics*, **9**, 223.
- 174 4. Cushman SA, Landguth EL (2010) Spurious correlations and inference in landscape
175 genetics. *Molecular Ecology*, **19**, 3592–3602.
- 176 5. Hotelling H (1933). Analysis of a complex of statistical variables into principal
177 components. *Journal of educational psychology* **24**, 417.
- 178 6. Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic
179 markers. *Bioinformatics*, **24**, 1403–1405.
- 180 7. Jombart T, Devillard S, Dufour AB, Pontier D (2008) Revealing cryptic spatial
181 patterns in genetic variability by a new multivariate method. *Heredity*, **101**, 92–103.
- 182 8. Jombart T, Pontier D, Dufour AB (2009). Genetic markers in the playground of
183 multivariate analysis. *Heredity* **102**, 330–341.
- 184 9. Jombart T, Devillard S, Balloux F (2010). Discriminant analysis of principal
185 components: a new method for the analysis of genetically structured populations.
186 *BMC genetics* **11**, 94.
- 187 10. Jombart T, Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-
188 wide SNP data. *Bioinformatics* **27**, 3070–3071.
- 189 11. Moran PAP (1950) Notes on Continuous Stochastic Phenomena. *Biometrika*, **37**,
190 17–23.

- 191 12. Novembre J, Stephens M (2008) Interpreting principal component analyses of
192 spatial population genetic variation. *Nature Genetics*, **40**, 646–649.
- 193 13. Pearson K (1901) On lines and planes of closest fit to systems of points in space.
194 *Philosophical Magazine Series 6*, **2**, 559–572.
- 195 14. Peres-Neto PR, Jackson DA, Somers KM (2005) How many principal components?
196 stopping rules for determining the number of non-trivial axes revisited.
197 *Computational Statistics & Data Analysis*, **49**, 974 – 997.
- 198 15. R Core Team (2017) R: A language and environment for statistical computing. *R*
199 *Foundation for Statistical Computing*, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)
200 [project.org/](https://www.R-project.org/).
- 201 16. Schiffers K , Travis JMJ (2014) ALADYN - a spatially explicit, allelic model for
202 simulating adaptive dynamics. *Ecography*, **37**, 1288–1291.

203 **Legends**

204

205 **Figure 1.** Graphical representation of the results reported in Tables 1 and 2. Only test with
206 threshold 0.05 are plotted.

207

208 **Figure 2.** Distributions of significant eigenvalues detected in the presence of global (blue
209 bars) and local (green bars) spatial patterns after hierarchical Bonferroni correction, for
210 100 significantly positive and 100 significantly negative patterns. Black bars correspond to
211 eigenvalues which are significant without Bonferroni correction. Bars' height indicates the
212 frequency of observing a significant eigenvalue in a certain position (from most positive to
213 most negative) over the 100 tested patterns.

214 **Table 1.** Significant results for *global test* (g test), *local tests* (l test), and *spca_randtest* (r test +/-) for random, global and local patterns
 215 using 200 loci per individual. IS, SS, IBD indicate the migration models (see Methods); different migration rates are coded by number: 1 =
 216 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
 217 thresholds .05 and .01.

200 SNPs	Models	<i>p</i> -value*	Random Patterns				Global Patterns				Local Patterns			
			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 *0.009* 0.011 0.012 **0.939** **0.999** *0.005* *0.000* *0.009* 0.023 **0.225** **0.684**

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

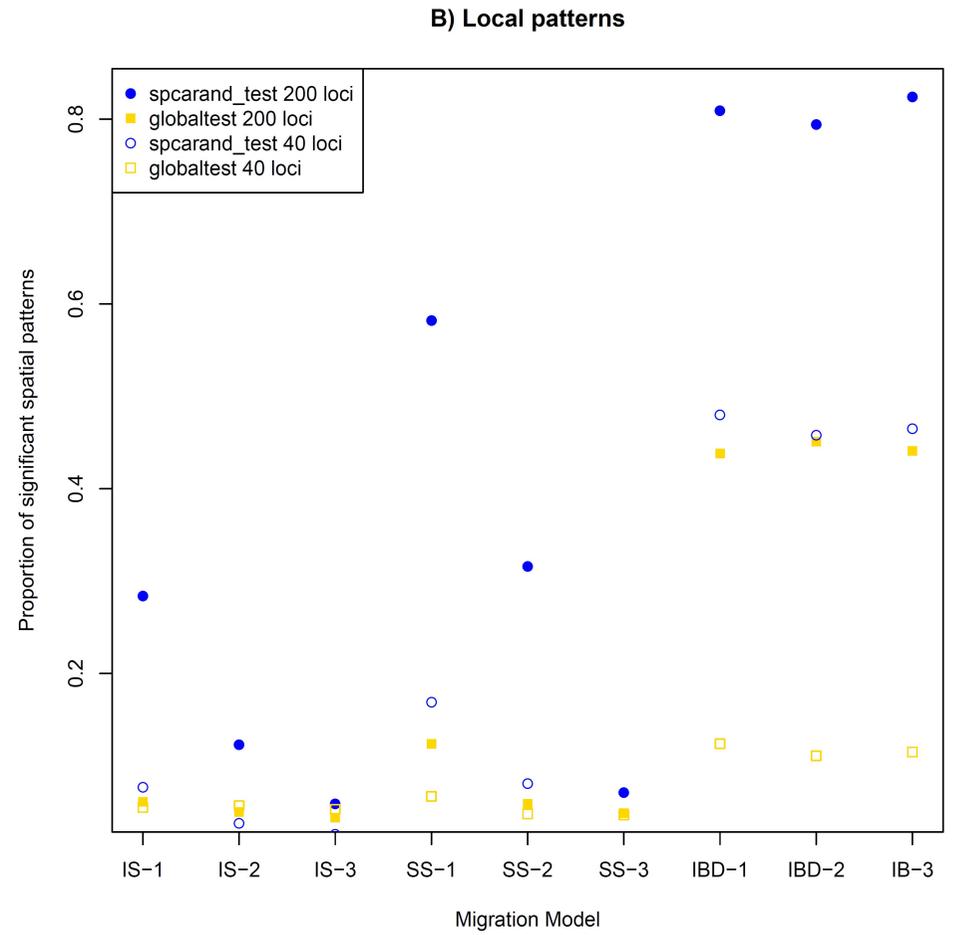
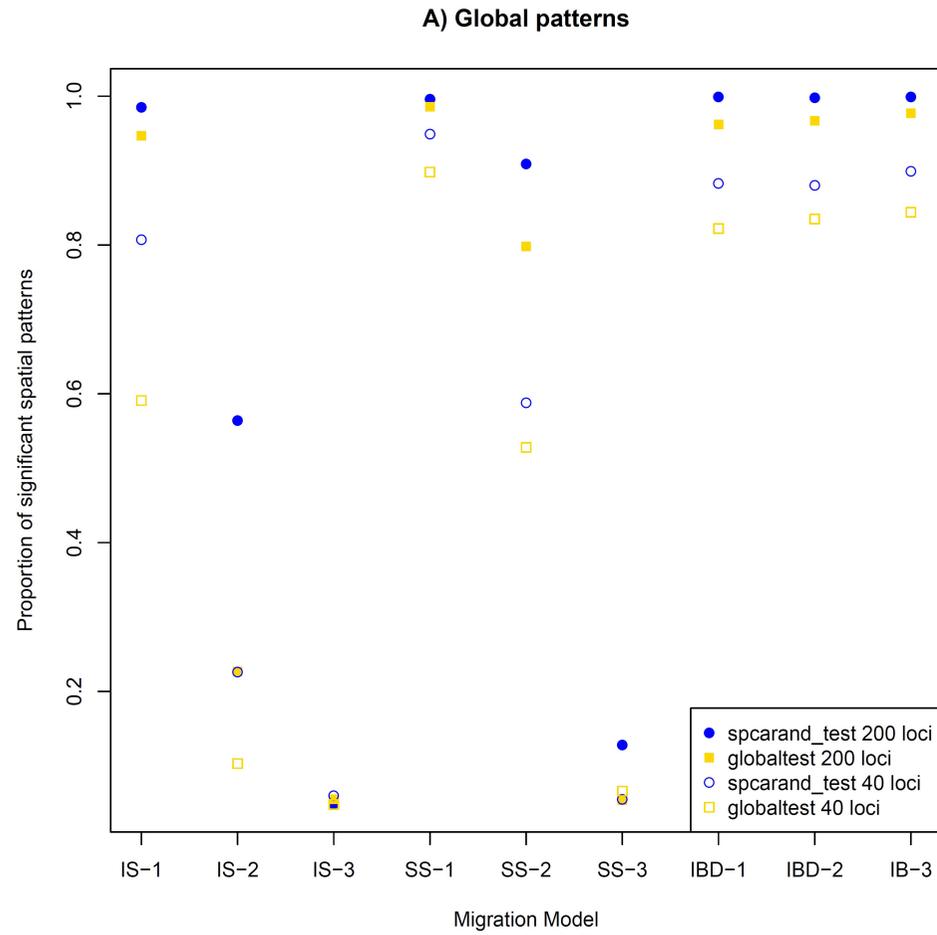
219 **Table 2.** Results for the same simulations reported in Table 1 using a subset of 40 loci per individual.

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

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

221 **Figure 1**

222



223 **Figure 2**

224

225

