

---

## Application of Moran Eigenvector Maps (MEM) to irregular sampling designs

Brind'Amour Anik <sup>1,\*</sup>, Mahévas Stephanie <sup>1</sup>, Legendre Pierre <sup>2</sup>, Bellanger Lise <sup>3</sup>

<sup>1</sup> Unité Écologie et Modèles pour l'Halieutique, IFREMER, Rue de l'île d'Yeu, B.P. 21105, 44311 Nantes Cedex 03, France

<sup>2</sup> Département de Sciences Biologiques, Université de Montréal, C.P. 6128, succursale Centre-ville, Montréal, Québec, Canada H3C 3J7

<sup>3</sup> Laboratoire de Mathématiques Jean Leray UMR CNRS 6629, Université de Nantes, Nantes, France

\* Corresponding author : Anik Brind'Amour, email address : [Anik.Brindamour@ifremer.fr](mailto:Anik.Brindamour@ifremer.fr)

---

### Abstract :

Moran's eigenvectors maps (MEM) are attractive mathematical objects as they are fairly simple to calculate and can be used in most studies of spatially-explicit data. There is, however, an aspect of MEM analysis that still requires some investigation: the effect of irregular sampling on their modeling performance. This study investigates empirically the behavior of MEMs under different irregularity schemes. It is focusing on simulated scenarios representing sampling designs frequently encountered in ecology. We advocate that MEMs can be computed and correctly used with data coming from irregularly designed sampling surveys, given some precautions. We suggest that when the sampling sites are equally spaced but do not cover the entire study area, the MEMs can be computed directly on the coordinates of the sampling sites without any important loss of information. Whereas, when the phenomenon of interest is tackled using randomly stratified sampling designs, the MEMs should be computed on a reconstructed space of regular sampling sites followed by removal of the missing sites, before analysis. This solution of rebuilding a (regular) sampling space guarantees to capture the underlying process under study, improves the modeling results and relaxes the impact of the choice of the weighting matrix on the computation of MEMs.

**Keywords :** Autocorrelation, Irregular sampling, Sampling schemes, Spatial analyses, Statistical methods

## 37 1. INTRODUCTION

38 Ecosystems in general and marine communities in particular are complex systems  
39 composed of a large number of entities interacting with one another at various spatial and  
40 temporal scales. Characterization of these scales is an essential step to understand and predict  
41 the effects of changes in the processes governing these systems. It relies on mathematical and  
42 statistical methods that allow the quantitative description of the spatial and temporal  
43 complexity and are sufficiently robust to handle any type of sampling designs. Marine  
44 ecological surveys are often irregular (i.e. unevenly spaced) in space or time. Irregularity  
45 seems to be the rule rather than the exception. Sampling irregularity may have different  
46 causes and consequently display different patterns. In this study, we are dealing with two  
47 types of irregularity: i) a "random" irregularity encountered when the phenomenon of interest  
48 is tackled using randomly stratified sampling designs or the dataset contains missing sites or  
49 time points, and ii) a "constrained" irregularity, when the sampling sites are equally spaced  
50 but do not cover the entire study area because topography or other constraints prevent  
51 sampling in some sections (i.e. partial coverage).

52 Observed spatial distributions of species may arise from a plurality of endogenous and  
53 exogenous processes (e.g. species interactions, growth, population dynamics, physical  
54 forcing) occurring at multiple spatial scales (Vaclavik et al. 2012). This mixture of processes  
55 and scales clearly calls for mathematical tools capable of accounting for or modeling such  
56 patterns (Dray et al. 2012). Among the statistical methods, the Moran Eigenvector Maps  
57 (MEM, Dray et al. 2006) and its original form, the Principal Coordinates of Neighbor  
58 Matrices (PCNM, Borcard and Legendre 2002), are good candidates for analyzing such  
59 patterns. MEMs are derived from spectral graph theory and characterize a wide range of  
60 autocorrelation structures based on the survey design, i.e. the distances between the  $n$   
61 sampling sites or times (Dray et al. 2006). It is thus a spectral decomposition of the spatial (or

62 temporal) relationships among the sampling sites (or dates). This decomposition generates  
63  $(n-1)$  *eigenfunctions*, which are new orthogonal variables that can be used in statistical  
64 models as explanatory variables representing the spatial or temporal relationships among the  
65 study sites.

66 The MEMs and derived approaches have proved very helpful in studying the spatial and  
67 temporal distributions of ecological communities (Bellchambers et al. 2011; Brind'Amour et  
68 al. 2005; Mikulyuk et al. 2011). However, these studies have been conducted almost  
69 exclusively in a context of regular sampling. Although no technical reason prevents the  
70 spectral decomposition in a context of irregular sampling sites (i.e. unequally-spaced  
71 sampling sites), little is actually known regarding the behavior of the MEMs in such a  
72 context. When Borcard and Legendre (2002) first introduced the PCNM approach, they  
73 mentioned that irregular sampling schemes affected the amplitude, the phase and the periods  
74 of the sine waves generated by the PCNM. They suggested that PCNM developed with  
75 irregular sampling sites are suitable descriptors but likely consist of multiple spatial scales,  
76 making the interpretation of the spatial descriptors more difficult. Dray et al. (2006)  
77 elaborated a little more on the sensitivity of the connectivity matrix in the case of irregular  
78 distribution of sampling sites and illustrated the consequences of the sampling irregularity on  
79 the number of positive/negative eigenvalues and the spatial structures described by the  
80 associated eigenvectors. They came to the conclusion that sampling irregularity, defined  
81 through the spatial relationships among neighboring sites, may have substantial impact on the  
82 behavior and interpretation of the MEMs. Ecological surveys are designed to study processes  
83 overriding the spatial and temporal distribution of species. However, in some cases it is  
84 difficult or impossible to sample over the whole area where a specific ecological process  
85 occurs. For instance, recruitment of a species may take place in shallow coastal areas where  
86 the draft of the boat prevents access. In that case, the observation scale (study extent) at which

87 sampling is conducted is smaller than the ecological scale at which the process occurs. Such a  
88 mismatch between the observation and process scales may also have a significant effect on  
89 the interpretation of the MEM.

90 In cases of regular sampling with some points of the grid missing, it is common practice  
91 to develop the MEMs on the original matrix of sampling coordinates of the irregularly  
92 distributed sites (Fuentes-Rodriguez et al. 2013; Jombart et al. 2008, Sattler et al. 2010,  
93 Mikulyuk et al. 2011, Sharma et al. 2011, Vaclavik et al. 2012,). However, as suggested by  
94 Borcard and Legendre (2002), one can develop the MEMs on a transformed matrix of  
95 coordinates that has been filled with supplementary sampling sites to make it regular; by  
96 construction MEMs are orthogonal to one another. The added sites are then removed after the  
97 MEM have been computed (i.e. rows of the eigenvector matrix; Blanchet et al. 2013 and  
98 Borcard et al. 2004) . This procedure presents however the disadvantage of losing practical  
99 mathematical properties of the MEMs: the orthogonality among the MEMs and the  
100 maximization of spatial autocorrelation (Moran's I). The choice of filling or not the voids with  
101 supplementary sites, prior to the computation of the MEMs, and the number of supplementary  
102 sites needed to attain a sufficiently fine resolution without major loss of orthogonality, remain  
103 key questions in the development and interpretation of MEMs.

104 This study aims at empirically shedding light on some key questions about the  
105 development and interpretation of MEMs: Is the MEM approach relevant with irregular  
106 sampling designs? Do the MEMs developed with irregular sampling sites capture the spatial  
107 scales they are supposed to capture? Is there an irregularity threshold beyond which MEMs  
108 cannot be safely used? Is there a solution to counteract the problems caused by sampling  
109 irregularity? Should we compute MEMs using supplementary sites, or not? Using  
110 simulations, we empirically investigated the impact of various irregular sampling schemes  
111 (irregular distribution of sites and partial coverage of a study area) in the development and

112 interpretation of the MEMs. We focused on a selection of scenarios characterizing sampling  
113 designs frequently encountered in ecological studies.

## 114 2. METHODS AND DATA

### 115 *2.1 Simulation approach*

116 Investigation of the MEM behavior with regard to sampling irregularity was done by  
117 testing the ability of the MEMs to correctly detect spatial structures commonly observed in  
118 ecology under three different scenarios of sampling irregularity. The simulations were  
119 designed to mimic sampling strategies frequently encountered in ecological studies. They  
120 were based on three components that may affect the computation and interpretation of the  
121 MEMs: i) irregularity of the sampling sites produced by sub-sampling the sampling zone (i.e.  
122 random sampling), ii) irregularity of the sampling sites generated by partial coverage of the  
123 sampling zone (i.e. blocks of unsampled sites), and iii) the process-observation mismatch, i.e.  
124 whether the observations match or not the scale of the ecological process under study. We  
125 focused on three combinations of these components that we called scenarios (see below).  
126 Each scenario is thus answering a specific question regarding irregularity in sampling  
127 strategies. All simulations were replicated 100 times by modifying the spatial coordinates of  
128 the samples. Precisely, we randomly sampled over the complete grid the number of cells  
129 corresponding to the subsampling thresholds.

130 In each scenario, we compared two approaches for computing MEMs (Fig. 1):

- 131 • The complete-grid approach ( $MEM_{comp}$ ): compute the MEM from the coordinates of  
132 all sites in the full grid, including the unsampled points, then remove from the MEM  
133 matrix the rows that correspond to the unsampled sites.
- 134 • The reduced-grid approach ( $MEM_{red}$ ): use only the geographic coordinates of the  
135 sampled sites to construct the MEM data matrix. That approach is commonly used in  
136 the literature.

### 137 *2.1.1 Scenario 1: Random sampling design (S1)*

138 The first scenario tested the ability of the MEM to correctly capture the spatial structure  
139 in the case of a random sampling strategy. Random sampling was assessed at eight different  
140 thresholds (25, 50, 60, 70, 80, 90, 100%), representing the percentage of the studied area  
141 covered by the sampled sites. Although the lower thresholds (25%) may seem quite small,  
142 they are representative of what ecologists use in randomly stratified scientific surveys  
143 (Brind'Amour et al. 2014). In this scenario, irregularity was random and created unequal  
144 distances between neighboring sampled sites (Fig. 2). The process under study matched the  
145 observation scale and occurred at "global" scale, that is, over the entire zone under study. In  
146 that scenario, the distances between neighboring sites increased in irregularity as the threshold  
147 decreased.

### 148 *2.1.2 Scenario 2: Blocks of missing data (S2)*

149 The second scenario was developed to test the ability of MEMs to correctly capture the  
150 spatial structure when the survey includes blocks of missing observations, such as  
151 inaccessible areas (Fig. 2; Jones et al. 2008). In this scenario, the sampled sites covered 50%  
152 of the whole area. They were regularly spaced, but as the area was greatly reduced, the  
153 distances between sites were irregular. The process under study occurred at "global" scale and  
154 did not match the observation scale.

### 155 *2.1.3 Scenario 3: Random sampling design and blocks of missing data on a global structure* 156 *(S3)*

157 This scenario tested the capacity of MEMs to capture the spatial structure when the  
158 survey was randomly designed and included blocks of missing observations (as in the second  
159 scenario). In that scenario, the sampled area covered 50% of the total area and we tested the  
160 effect of different sampling thresholds on the MEMs (Fig. 2). The process occurred at  
161 "global" scale and did not match the observation scale.

## 162 2.2 MEM computation

163 There are various ways of computing MEMs (see Dray et al. 2006 and Legendre and  
164 Legendre 2012 for details). In our simulations, we computed both MEMs and db-MEMs  
165 following the steps using the packages *spacemake R* (Dray 2013) and *spdep* (Bivand 2011) in  
166 R (R Core Team 2014). It is worth mentioning that the package *adespatial* recently developed  
167 by Dray et al. (2016) can now compute the MEMs and db-MEMs. Differences between the  
168 two types of MEMs are summarized in Appendix A1. The two types were calculated on a  
169 matrix of  $20 \times 20 = 400$  sites for the three scenarios. They were computed using a distance  
170 matrix transformed into a similarity matrix (Legendre and Legendre 2012, p.861) weighed by  
171 a connectivity matrix (see details in Appendix A1).

## 172 2.3 Predefined spatial structures

173 Different ways can be used to simulate spatial structures. For instance, one can use the  
174 MEMs themselves to generate response values or simulate independent geostatistical  
175 distributions. In here, we simulated the spatial structures using the MEM themselves and  
176 using empirical variograms with various ranges to modify the degree of spatial  
177 autocorrelation. The use of the MEMs themselves was done as an "experiment" to verify if we  
178 could correctly capture the predefined MEM as a spatial structure. With that approach we  
179 were expecting to capture perfectly the modeled spatial structure with the  $MEM_{comp}$  given that  
180 the same MEM served as the response and it was also included in the set of explanatory  
181 variables. The predefined spatial structures were created from MEMs computed on a 2D  
182 regular grid of 20 by 20 cells ( $n = 400$  cells). For the three scenarios, we selected four MEMs  
183 from the entire set produced (MEM #1 called MEM01, 10, 150, 350; Fig. 3). In the three  
184 scenarios, the selected MEMs were analyzed as separate response variables (i.e. single  
185 variable), corresponding to a gradient of spatial structures varying from coarse to very fine  
186 spatial scales. Coarse spatial scale was characterized by large positive eigenvalues, medium

187 spatial scale by intermediate positive eigenvalues, fine spatial scale by small negative  
 188 eigenvalues, and very fine spatial scale by large negative eigenvalues. For each  
 189 studied  $j^{th}$  MEM, called  $Y^j$ , representing a selected predefined spatial structure, a random  
 190 noise was added using values sampled from a normal distribution with mean of 0 and  $\sigma$  of  
 191 0.05,  $N(0,\sigma)$ .

192 The second type of predefined spatial structures were developed using geostatistical  
 193 distributions. The four predefined spatial structures modeled corresponded to a gradient of  
 194 spatial scales varying from coarse to very fine spatial scales (Fig. 3). Empirical variograms  
 195 were developed using spatially correlated random fields computed on a 2D regular grid of 20  
 196 by 20 cells ( $n = 400$  cells) followed by unconditional Gaussian simulations (Pebesma 2004).  
 197 The variogram ( $\gamma$ ) was modeled using a spherical model (Cressie 1993, for more details):

$$\begin{aligned}
 \gamma(h) &= \frac{c}{2} \left( \frac{3h}{a} - \frac{h^3}{a^3} \right) & h \leq a \\
 \gamma(h) &= c & h > a
 \end{aligned} \tag{1}$$

201 where  $\gamma(h)$  is the variogram value at distance  $h$ ,  $a$  is the range, i.e. the distance,  $h$ , beyond  
 202 which the autocorrelation is presumably zero, and  $c$  is the sill, i.e. the value at which the  
 203 variogram levels off. For the simulations we arbitrarily fixed  $c$  at a value equal to 5 and varied  
 204 the range from a coarse ( $a = 20$  cell length), to medium ( $a = 5$ ), to fine ( $a = 3$ ), and very fine  
 205 spatial autocorrelation structure ( $a = 1$ ). To simulate the four spatial structures we fixed the  
 206 three beta coefficients to zero (i.e. no linear trend) and as mentioned above we varied the  
 207 range. The simulations were repeated 100 times per spatial structure using R package *gstat*  
 208 (Pebesma 2004).

#### 209 *2.4 MEM evaluation*

210 The predefined spatial structures were used as response variables  $Y$  in linear regression  
 211 models to evaluate the MEM behavior. We compared the two following:

$$212 \quad Y_i^{(j)} = m + (\text{MEM}_{comp})_i \times \beta_{comp}^{(j)} + \varepsilon_i^{(j)} \text{ and,} \quad (2)$$

$$213 \quad Y_i^{(j)} = m + (\text{MEM}_{red})_i \times \beta_{red}^{(j)} + \varepsilon_i^{(j)} \quad (3)$$

214 where  $\varepsilon_i^{(j)} \sim N(0; \sigma)$  and  $Y_i^{(j)}$  is fixed to be either one of the four MEMS or variogram  
 215 simulations, in the three scenarios.  $n$  is the number of observations and depends on the  
 216 number of missing observations.  $\text{MEM}_{comp}$  (resp.  $\text{MEM}_{red}$ ) is the matrix of predictor  
 217 variables,  $(MEM)_i$  is a line vector containing the  $i^{\text{th}}$  line of matrix MEM,  $\beta$  is the column  
 218 vector of parameters (regression coefficients) to be estimated, and  $m$  is the intercept.

219 The ability of MEM analysis to correctly capture the predefined spatial structures was  
 220 evaluated using five criteria: the number of significant  $\text{MEM}_{comp}$  and  $\text{MEM}_{red}$  in each  
 221 scenario, the adjusted  $R^2$  of the fitted models, and the collinearity (estimated by Pearson  
 222 correlation coefficients) between the members of the subset of  $\text{MEM}_{comp}$  considered as  
 223 predictors. The first criterion (i.e. number significant MEM) was obtained from multiple  
 224 regression analyses after forward selection between the response variable  $Y$  and the  
 225 explanatory MEM (spatial descriptors). Given that  $(n-1)$  MEMs are generated, the choice of a  
 226 method to correctly select significant MEMs with regard to the overestimation of the variance  
 227 explained is an important issue that has been discussed (Blanchet et al. 2008). In here, we  
 228 used a forward selection analysis based on a permutation procedure (*forward.sel* function)  
 229 developed by S. Dray in the R package *packfor*. It follows the recommendations of Blanchet  
 230 et al. (2008) and Munoz et al. (2009) to split the regression model in different parts (or  
 231 submodels) to circumvent the problem of over parameterization. In here we divided MEMs  
 232 into four submodels corresponding to a gradient of spatial structures varying from coarse to  
 233 very fine spatial scales. Coarse spatial scale was mainly characterized by large positive  
 234 eigenvalues, medium spatial scale by intermediate positive eigenvalues, fine spatial scale by  
 235 small negative eigenvalues, and very fine spatial scale by large negative eigenvalues. The  
 236 significance of regressions was tested as suggested in Blanchet et al. (2008), by applying a

237 forward selection on each submodel with a double stopping rule (i.e.  $\alpha$  threshold and a  
238 maximum threshold for the global model). This procedure controls for type I error inflation.  
239 The predicted values were estimated using linear regression models by fitting the significant  
240 MEM (previously identified) to the response variable  $Y$  (i.e. predefined spatial structure).

### 241 3. RESULTS

#### 242 3.1 Scenarios

243 The simulations were conducted using MEMs and db-MEMs. As no difference was  
244 found between the two types of MEMs, only the results with MEM, generalization of db-  
245 MEM, are presented here for two extreme cases of the predefined spatial structures (Fig. 3):  
246 the coarser and finer spatial scales using the MEMs (MEM01 and MEM350) and the  
247 variogram simulations (ranges = 20 and 3). The Appendix B contains results for all the spatial  
248 scales (including those presented here) calculated for the two types of spatial structures  
249 (MEMs and variograms) and the three scenarios.

##### 250 3.1.1 Scenario 1: Random sampling design (S1)

251 *MEM spatial structures*– Comparison of the  $MEM_{comp}$  and the  $MEM_{red}$  suggests that the first  
252 approach globally outperformed the second (Fig. B1, Appendix B). The  $MEM_{comp}$  always  
253 captured the predefined spatial structure and showed adjusted  $R^2$  always above those of  
254  $MEM_{red}$  (Fig. 4a and 4d). That result stands for the sub-sampling thresholds that could be  
255 tested (i.e. > 50%) independently of the spatial scale. On the other hand, when the MEMs are  
256 computed directly from the reduced matrix of coordinates (the reduced approach), it takes  
257 between 2 to 7 MEMs to capture the predefined spatial structure (regardless of the sub-  
258 sampling threshold). That approach of computing MEMs succeeds in modeling the coarse  
259 predefined spatial structure with adjusted  $R^2$  comparable to those obtained for the  $MEM_{comp}$   
260 but fails in capturing fine spatial structures with adjusted  $R^2$  varying between 0 and 0.5.

261 With the  $MEM_{red}$  the property of orthogonality is preserved and the MEMs are thus  
262 uncorrelated to one another. This property also holds for all sub-sampling thresholds. This is  
263 not the case for the  $MEM_{comp}$  where orthogonality is lost when the missing sampling sites are  
264 removed. Nevertheless, the correlation coefficients among the  $MEM_{comp}$  are very low as they  
265 never reach values higher than 0.12 (Fig. B1, Appendix B).

266 *Variogram simulations*— Results of the modeling of the variogram simulations using the  
267  $MEM_{comp}$  and the  $MEM_{red}$  are in line with the previous results using the MEM as spatial  
268 structures (Fig. B4, Appendix B): the  $MEM_{comp}$  slightly outperformed the  $MEM_{red}$ . The  
269  $MEM_{comp}$  always captured the predefined spatial structure and showed adjusted  $R^2$  almost  
270 always above those of  $MEM_{red}$  (Fig. 5a). However, as shown previously, at low sampling  
271 thresholds (i.e.  $< 50\%$ ) when the number of explanatory variables equals or exceeds the  
272 number of sites, the selection procedure stops and no  $MEM_{comp}$  is included in the models (Fig.  
273 5d), thereby lowering the adjusted  $R^2$  to 0. On the other hand, when the MEMs are computed  
274 directly from the reduced matrix of coordinates (the *reduced* approach), it takes between 10 to  
275 40 of the generated MEMs to capture the predefined spatial structure. That approach of  
276 computing MEMs succeeds in capturing the predefined spatial structure below 50% but  
277 displays higher variability at low sampling thresholds. For the two approaches, the effect of  
278 sub-sampling on the global fit of the predefined spatial structure grows worst as the scale of  
279 the spatial structure decreases. For instance, when the structures are characterized by coarse  
280 spatial scales, the adjusted  $R^2$  stabilizes at  $\sim 0.85$  for all the sub-sampling thresholds above  
281 50% of the sampled area (Fig. 5a). When spatial structures are defined at fine spatial scales,  
282  $MEM_{comp}$  shows on average a better fit of  $\sim 10\%$ . The two approaches fail in capturing the  
283 spatial structures at low sampling thresholds.

### 3.1.2 Scenario 2: Blocks of missing data (S2)

MEM spatial structures. – When the spatial structures occur at coarse (MEM01 and MEM10) spatial scales, the MEM<sub>red</sub> and MEM<sub>comp</sub> give similar results (Fig. B2, Appendix B). The two approaches slightly differ when the spatial structures are at fine scales (MEM ≤ 150). In these cases, it takes on average 5 to 15 MEM<sub>red</sub> to detect the predefined spatial structure but it never succeeds in modeling the spatial structure as efficiently as with MEM<sub>comp</sub> (Fig. 4b). Indeed, when the MEM<sub>comp</sub> are used, they systematically captured the modeled MEM and showed adjusted R<sup>2</sup> values 5 to 15% higher than when using the MEM<sub>red</sub>. In these cases, the collinearity induced by removing the supplementary sites is always well below 0.1% (Fig. B2, Appendix B).

Variogram simulations. – When the predefined spatial structures are developed using variogram simulations, the results are very similar to those presented above, i.e. the MEM<sub>red</sub> and MEM<sub>comp</sub> give similar results notwithstanding the simulated spatial scales (Fig. 5b and Fig. B5, Appendix B). They both show decreasing adjusted R<sup>2</sup> and increasing uncertainty at medium and fine spatial structures.

### 3.1.3 Scenario 3: Random sampling design and blocks of missing data on a global structure (S3)

MEM spatial structures. – Results for that scenario are similar to those obtained in S1: the MEM<sub>comp</sub> clearly outperformed the MEM<sub>red</sub> (Fig. 4c and 4f). In contrast to S1, the results in S3 indicate a stronger impact of increasing irregularities (i.e. increasing sub-sampling) on the capacity of the MEM to correctly detect the predefined spatial structure (Fig. B3, Appendix B). The MEM<sub>comp</sub> systematically selected the predefined spatial structure (Fig. 4f) and reached better fits than the MEM<sub>red</sub>, regardless of the level of sub-sampling and the nature of the predefined spatial structure (i.e. coarse to very fine structures; Fig. B3 in Appendix B). The effect of irregularity is most obvious using the MEM<sub>red</sub> and notably when modeling fine

309 spatial structures. In that case, the adjusted  $R^2$  drops more rapidly than it did in S1 and never  
310 reaches values above 0.3 (Fig. 4c). The counterpart of using  $MEM_{comp}$  is emphasized by the  
311 collinearity, which sometimes reaches values equal to 0.1 (Fig. B3, Appendix B).

312 *Variogram simulations.*—Results for that scenario are very similar to those obtained in S1: the  
313  $MEM_{comp}$  slightly outperformed the  $MEM_{red}$  (Fig. 5c and 5f). In contrast to S1, the results in  
314 S3 indicate a stronger difference between the two approaches on the capacity of the MEM to  
315 correctly detect the spatial structure (Fig. B6, Appendix B). The  $MEM_{comp}$  selected the  
316 simulated spatial structure (Fig. 5c) and reached better fits between 10 to 40% in comparison  
317 to the  $MEM_{red}$ , at levels of sub-sampling above 50% (Fig. B6 in Appendix B). As in S1, at  
318 low sampling thresholds (i.e. < 50%), very few  $MEM_{comp}$  are included in the models whereas  
319 5 to 10  $MEM_{red}$  are selected to reach Adjusted  $R^2$  varying between 0.70 (coarse spatial scales)  
320 and 0.1 (fine spatial scales).

#### 321 4. DISCUSSION

322 The number of studies using the MEM approach and its derivatives has more than  
323 doubled in recent years (Bellchambers et al. 2011; Blanchet et al. 2013; Fuentes-Rodriguez et  
324 al. 2013; Jombart et al. 2008; Mikulyuk et al. 2011; Sattler et al. 2010; Sharma et al. 2011;  
325 Sharma et al. 2012; Vaclavik et al. 2012); the original papers describing the method received  
326 hundreds of citations listed on Web of Science and Google Scholars. Most of these  
327 applications used irregular sampling designs. However, very few of them have actually  
328 discussed the effect of irregularity on the development and interpretation of the MEM  
329 (Blanchet et al. 2013 and Borcard et al. 2004). Our study aimed at investigating empirically  
330 the capacity of MEM analysis to correctly identify predefined spatial structures at various  
331 spatial scales, under different scenarios of irregularity. This was done to help ecologists use  
332 the full potential of the MEM approach in ecological modeling. We suggest to develop the  
333 MEMs on a regular sampling grid, followed by removal of the missing sites. We also warn

334 against sampling irregularity when the sampling sites cover a low proportion of the studied  
335 area and when one wishes to model ecological processes occurring at very fine spatial scales.

#### 336 *4.1 MEM<sub>red</sub> or the common way of computing the MEMs*

337 Our study tested the performance of two ways of computing MEMs (MEM<sub>comp</sub> and MEM<sub>red</sub>)  
338 to correctly captured different predefined spatial structures. This was done with the objective  
339 of comparing the commonest way of computing the MEM (MEM<sub>red</sub>) with another less  
340 common approach (MEM<sub>comp</sub>). We used two types of predefined spatial structures, one  
341 based on the MEM themselves, that can seen as tautological (or dependent) with the  
342 MEM<sub>comp</sub>, and another one using independent spatial structures. We considered the primer  
343 predefined spatial structure as a "controlled situation" where we were expecting to capture  
344 perfectly the pattern using the MEM<sub>comp</sub>. In that sense, the MEM<sub>comp</sub> responded as expected  
345 and gave almost a perfect fit notwithstanding the scenario and the spatial scales. It was more a  
346 less a test for the MEM<sub>red</sub> as most of the studies using the MEMs are developing the MEMs  
347 directly on the sample sites without filling the voids (e.g. Mikulyuk et al. 2011; Sattler et al.  
348 2010). For that type of MEMs construction, results were generally considered good or  
349 lukewarm at broader and finer spatial scales respectively.

350 The independent predefined spatial structures (i.e. variograms) were used as comparison  
351 between the two types of MEMs. When the selection procedure allowed the MEM<sub>comp</sub> to be  
352 computed (threshold > 50%), they MEMs showed between 10 to 15% better fit than the  
353 MEM<sub>red</sub>. At low thresholds, i.e. when the proportion of sites sampled is low given the surface  
354 of the studied zone (< 50% or less) and broad spatial structures are expected, we suggest that  
355 the MEM<sub>red</sub> can be safely used.

#### 356 *4.2 Irregularity: effect of random design vs blocks of missing data*

357 In this study we showed that removing blocks of sampling sites (e.g. scenario *S2*) was  
358 less harmful to the conclusions than randomly removing the same number of sites (50%

359 threshold in  $S1$ ) in a regularly-spaced design. This can be explained by the fact that the  
360 proportion of regular distances among the sites in the first case is kept relatively high in  
361 comparison to the second case where any distance can be eliminated. Recent studies using the  
362 MEMs (Bellchambers et al. 2011; Fuentes-Rodriguez et al. 2013; Mikulyuk et al. 2011;  
363 Sattler et al. 2010; Sharma et al. 2011; Vaclavik et al. 2012) fell in our  $S2$  and  $S3$  scenarios  
364 with varying sub-sampling thresholds (all below 40%). They developed the MEMs directly on  
365 the coordinates of the sampling sites without adding supplementary sites. Given the results of  
366 our simulations, these studies might have missed some spatial scales of variability and  
367 presumably underestimated the importance of the predictors in terms of their contributions to  
368 the overall goodness of fit of their models. While this shortcoming most likely did not affect  
369 the spatial patterns they observed, it might have had some influence on the relative  
370 contributions of the spatial components they estimated in their variance partitioning analyses  
371 (Fuentes-Rodriguez et al. 2013; Mikulyuk et al. 2011; Sattler et al. 2010).

372         Munoz (2009) developed and tested a smoothing model to select significant distance-  
373 based eigenvector maps (DBEM, a particular case of MEM), on a regular and irregular  
374 sampling designs. He found no differences between the two designs and concluded that the  
375 DBEM approach was highly suitable for analyzing ecological surveys. Munoz results cannot  
376 be compared directly with ours, as his smoothing model does not keep, by definition,  
377 individual elements (i.e. MEMs) but rather combines them in sub-models using smoothing  
378 windows. Notwithstanding this difference, our results showed that at a low sub-sampling  
379 threshold (>90% of the sites kept) the models developed using the  $MEM_{red}$  produced similar  
380 but not as good results as those developed with the  $MEM_{comp}$  approach. In a sense, this is in  
381 agreement with what Munoz observed in his work as his regular and irregular sampling  
382 schemes were composed of the same number of sites (2500 points) and only differed in their

383 spatial positions. Whether the conclusions of Munoz (2009) would still hold under different  
384 sub-sampling thresholds and other irregular schemes remains an open question.

#### 385 *4.3 Rebuilding regularity: an efficient solution*

386 When they first introduced the PCNM method, Borcard and Legendre (2002) suggested  
387 to outwit the problem of missing observations by adding geographic coordinates in the dataset  
388 prior to MEM computation. The solution used in here slightly differs from theirs, as we are  
389 filling the voids as they suggested, but we are adding supplementary sites to mimic a regular  
390 sampling scheme. As advocated by Dray et al. (2006, p.487), the choice of a spatial weighting  
391 matrix  $W$  is crucial in the computation of MEMs and in the case of regular sampling, the  
392 structures defined by the eigenvectors (i.e. MEMs) are less sensitive to the choice of  $W$ .  
393 Therefore, recreating a regular sampling matrix offers the advantage of allowing the  
394 computation of the MEMs using any neighboring relationships in  $W$ , in addition to keeping a  
395 fine spatial resolution among the sites. In here we rebuilt a "rectangular cuboid" grid by using  
396 the maximum and minimum values on the X and Y axes. This way of recreating a complete  
397 and regular sampling grid may not always be the optimal technique, particularly when the  
398 sampling zone has the shape of a "rectangular parallelepiped". In that case, our technique may  
399 artificially expand the sampling zone and thus the number of MEM. We suggest that special  
400 care should be taken when developing the complete sampling grid (i.e. cell size and shape of  
401 the total extent). On the other hand, if one decide to use the MEMs that are computed directly  
402 on the matrix with missing observations, the choice of the weighting matrix should be  
403 optimized (Dray et al. 2006).

404 Building the MEMs on a reconstructed matrix of regular sampling has two drawbacks. First, it  
405 introduces correlations among the MEM, thus losing, to a certain degree, the orthogonality  
406 property of the MEM. This was already pointed out by Borcard and Legendre (2002) and  
407 Borcard et al. (2004, p. 1828). In here we confirm that statement: with irregular sampling

408 surveys, one has to accept the compromise of losing the appealing property of orthogonality  
409 in the modeling process. Second, the selection procedure (if correctly applied) stops when the  
410 number of variables equals or exceeds the number of sites. This situation can be circumvented  
411 by maintaining the number of variables lower than the number of sampling sites by, for  
412 instance, dividing the MEMs selection in different submodels and correcting accordingly for  
413 type 1 error. In This situation cannot happen with the  $MEM_{red}$ , because their number will  
414 always be less than the total number of sites and they in our study, they may captured a spatial  
415 structure notwithstanding the sub-sampling level. In cases where sampling irregularity is very  
416 high and induces strong correlations between  $MEM_{comp}$ , one may use  $MEM_{red}$ . However, in  
417 such cases the MEM approach may not be the most appropriate modeling technique, although  
418 alternatives are scarce (*e.g.* Empirical Orthogonal Functions, Kutzbach 1967).

419 In this study, we explored the efficiency of rebuilding a regular sampling grid prior to  
420 calculation of the MEMs by testing the ability of the “reconstructed” MEMs to correctly  
421 capture the different predefined spatial structures. Application of such a solution indicated  
422 that the predefined spatial structure was identified and showed very good adjustment for all  
423 the scenarios, using an appropriate statistical selection procedure. That solution also  
424 succeeded well in modeling the various spatial structures tested, from coarse to medium  
425 spatial scales. Rebuilding a regular sampling grid has an interesting advantage of assessing  
426 the inter-annual comparison of spatial structures using randomly stratified sampling designs.  
427 Given that the reconstructed grid is common to all the sampling years, the spatial analyses can  
428 thus be done on the same basis thereby allowing direct comparison of the spatial scales  
429 among years.

430 *4.4 Strong effect of irregularity with fine spatial structures.*—Whether very small and  
431 negative eigenvalues should be included in a modeling process remains an open debate  
432 (Munoz 2009 and Mahevas et al. unpublished) and is beyond the scope of this study.

433 Nevertheless, our simulations and other authors suggest that it would be probably safer to  
434 discard them in highly irregular sampling surveys (Blanchet et al. 2011). Indeed, the MEMred  
435 were less efficient in capturing the finer spatial scales, as showed their lower  $R^2$ . The failure  
436 of the two MEMs to capture the fine spatial scales using the variogram simulations (i.e. range  
437 equal to 1) can be explained by the choice of the grid size that we used (distance between two  
438 sites equals 1). This underlines that when the scale of the pattern is smaller than the sampling  
439 design, we are not able to detect a signal.

## 440 5. CONCLUSION

441 Our simulations dealt with relatively simple spatial structures (i.e. simulations with  
442 various ranges or predefined MEMs), while in most ecological studies, the spatial distribution  
443 of species is more complex and varies over a wide range of spatial scales. We tested two ways  
444 of computing the MEMs and conclude that both approaches can be used. Nevertheless some  
445 precautions must be taken to prevent their misuse. When the MEMs are computed on the  
446 “*complete*” or reconstructed space of the sampling sites prior to analysis, we suggest that the  
447 sets of eigenvectors with positive and negative eigenvalues can be used safely together in  
448 further analyses, given that a relevant selection procedure of significant variables is used. In  
449 our simulations no difference was found between the computation of MEMs or db-MEMs. In  
450 all scenarios, developing MEM over the complete area and subsequently reducing them to fit  
451 the sampling design created correlations among the MEMs (i.e. non-orthogonal eigenvectors),  
452 however, the values of the correlations were low (maximum of 0.10 in absolute value) and did  
453 not preclude the use of the  $MEM_{comp}$  as spatial descriptors in (partial) regression or canonical  
454 analyses. The importance of the correlations among the MEMs in highly complex ecosystems  
455 remains to be tested. For that particular aspect we call upon mathematicians to study the  
456 properties of the  $MEM_{comp}$  in a reduced sampling design and particularly the loss of  
457 equivalence between eigenvalues and Moran's  $I$ . We also showed that MEMs can be

458 computed directly on the coordinates when blocks of sites are missing, without any significant  
459 loss of information, and correctly interpreted if the process under study matches the scale of  
460 observation, which is generally the case. However, when the MEMs are computed directly  
461 using the spatial coordinates, special care should be taken in defining a relevant connectivity  
462 matrix and thus choosing appropriate neighboring relationships. The developments in here  
463 were applied in a spatial context, although similar conclusions could likely be drawn for  
464 temporal analyses.

#### 465 *Acknowledgments*

466 This work was carried out under the project COSELMAR funded by the Regional  
467 Council of the Pays de la Loire. The authors would like to thank the scientists and crews who  
468 participated in the NURSE surveys in the Bay of Vilaine nursery grounds between 2008 and  
469 2010. The authors acknowledge Stéphane Dray for his review and his comments that greatly  
470 improved the manuscript.

#### 471 REFERENCES

- 473 Bellchambers, L.M., Meeuwig, J.J., Evans, S.N., Legendre, P., 2011. Modelling habitat  
474 associations of the common spider conch in the Cocos (Keeling) Islands. *Mar. Ecol.*  
475 *Prog. Ser.* 432, 83-90
- 476 Blanchet, F.G., Bergeron, J.A.C., Spence, J.R., He, F., 2013. Landscape effects of  
477 disturbance, habitat heterogeneity and spatial autocorrelation for a ground beetle  
478 (*Carabidae*) assemblage in mature boreal forest. *Ecography* 36, 636-647
- 479 Blanchet F. G., Legendre, P., Borcard D., 2008. Forward selection of explanatory variables.  
480 *Ecology* 89, 2623-2632
- 481 Borcard, D., Legendre, P., 2002. All-scale spatial analysis of ecological data by means of  
482 principal coordinates of neighbour matrices. *Ecol. Model.* 153, 51-68
- 483 Borcard, D., Legendre, P., Avois-Jacquet, C., Tuomisto, H., 2004. Dissecting the spatial  
484 structure of ecological data at multiple scales. *Ecology* 85, 1826-1832

- 485 Brind'Amour, A., Boisclair, D., Legendre, P., Borcard, D., 2005. Multiscale spatial  
486 distribution of a littoral fish community in relation to environmental variables.  
487 *Limnol. Oceanogr.* 50, 465-479
- 488 Brind'Amour, A., P., L., J., M., S., V., Fovau, A., Le Bris, H., 2014. Morphospecies and  
489 taxonomic sufficiency of benthic megafauna in scientific bottom trawl surveys. *Cont.*  
490 *Shelf Res.* 72, 1-9
- 491 Dray, S., Legendre, P., Peres-Neto, P.R., 2006. Spatial modelling: a comprehensive  
492 framework for principal coordinate analysis of neighbour matrices. *Ecol. Model.* 196,  
493 483-493
- 494 Dray, S., Péliissier, R., Couteron, P., Fortin, M.-J., Legendre, P., Peres-Neto, P.R., Bellier, E.,  
495 Bivand, R., Blanchet, F.G., De Cáceres, M., Dufour, A.-B., Heegaard, E., Jombart, T.,  
496 Munoz, F., Oksanen, J., Thioulouse, J., Wagner, H.H., 2012. Community ecology in  
497 the age of multivariate multiscale spatial analysis. *Ecol. Monogr.* 82, 257-275
- 498 Dray, S., Blanchet, G., Borcard, D., Guenard, G., Jombart, T., Larocque, G., Legendre, P.,  
499 Madi, N., Wagner, H.H., 2016. *adespatial: Multivariate Multiscale Spatial Analysis*. R  
500 package version 0.0-7. <http://CRAN.R-project.org/package=adespatial>
- 501 Fuentes-Rodriguez, F., Juan, M., Gallego, I., Lusi, M., Fenoy, E., Leon, D., Penalver, P.,  
502 Toja, J., Casas, J.J., 2013. Diversity in Mediterranean farm ponds: trade-offs and  
503 synergies between irrigation modernisation and biodiversity conservation. *Fresh. Biol.*  
504 58, 63-78
- 505 Jombart, T., Devillard, S., Dufour, A.-B., Pontier, D., 2008. Revealing cryptic spatial patterns  
506 in genetic variability by a new multivariate method. *Heredity* 101, 92-103
- 507 Jones, M.M., Tuomisto, H., Borcard, D., Legendre, P., Clark, D.B., Olivas, P.C., 2008.  
508 Explaining variation in tropical plant community composition: influence of  
509 environmental and spatial data quality. *Oecologia* 155, 593-604
- 510 Kutzbach, J.E., 1967. Empirical eigenvectors of sea-level pressure, surface temperature and  
511 precipitation complexes over North America. *Journal of Applied Meteorology* 6, 791-  
512 802
- 513 Mahevas, S., Bellanger, L., Trenkel, V.M., 2008. Cluster analysis of linear model coefficients  
514 under contiguity constraints for identifying spatial and temporal fishing effort patterns.  
515 *Fish. Res.* 93, 29-38
- 516 Mikulyuk, A., Sharma, S., Van Egeren, S., Erdmann, E., Nault, M.E., Hauxwell, J., 2011. The  
517 relative role of environmental, spatial, and land-use patterns in explaining aquatic  
518 macrophyte community composition. *Can. J. Fish. Aquat. Sci.* 68, 1778-1789
- 519 Munoz, F., 2009. Distance-based eigenvector maps (DBEM) to analyse metapopulation  
520 structure with irregular sampling. *Ecol. Model.* 220, 2683-2689
- 521 Pebesma, E.J., 2004. Multivariable geostatistics in S: the gstat package. *Comp. & Geosc.* 30,  
522 683-691

- 523 R Development Core Team, 2014. R: A language and environment for statistical computing.  
524 R Foundation for Statistical Computing, Vienna, Austria
- 525 Rochet, M.-J., Trenkel, V., Bellail, R., Coppin, F., Le Pape, O., Mahé, J.-C., Morin, J.,  
526 Poulard, J.-C., Schlaich, I., Souplet, A., Vérin, Y., Bertrand, J., 2005. Combining  
527 indicator trends to assess ongoing changes in exploited fish communities: diagnostic  
528 of communities off the coasts of France. *ICES J. Mar. Sci.* 62, 1647-1664
- 529 Sattler, T., Borcard, D., Arlettaz, R., Bontadina, F., Legendre, P., Obrist, M.K., Moretti, M.,  
530 2010. Spider, bee, and bird communities in cities are shaped by environmental control  
531 and high stochasticity. *Ecol.* 91, 3343-3353
- 532 Sharma, S., Legendre, P., Boisclair, D., Gauthier, S., 2012. Effects of spatial scale and choice  
533 of statistical model (linear versus tree-based) on determining species-habitat  
534 relationships. *Can. J. Fish. Aquat. Sci.* 69, 2095-2111
- 535 Sharma, S., Legendre, P., De Caceres, M., Boisclair, D., 2011. The role of environmental and  
536 spatial processes in structuring native and non-native fish communities across  
537 thousands of lakes. *Ecography* 34, 762-771
- 538 Vaclavik, T., Kupfer, J.A., Meentemeyer, R.K., 2012. Accounting for multi-scale spatial  
539 autocorrelation improves performance of invasive species distribution modelling  
540 (ISDM). *J. Biogeogr.* 39, 42-55

541

542 Figure captions

543 Fig. 1. Flowchart of the methodology used to develop the *reduced* and the *complete* MEM  
544 sets in this study. Texts in bold indicate the final dimensions of the MEM matrix. The full  
545 grid in this example includes 72 sites whereas in the real study, we used 200 sites (see  
546 Methods for details).

547

548 Fig. 2. (a) Schematic 2D illustration of the four scenarios tested in this study. X and Y axes  
549 represent longitude and latitude respectively. The gray scale corresponds to a predefined  
550 spatial structure used as an example of a response variable  $Y$  in this study. Detailed  
551 description is found in text with the corresponding scenario.

552

553 Fig. 3. Illustration of the two types of predefined spatial structures used in the study,  
554 corresponding to a gradient of spatial scales varying from coarse to very fine spatial scales.  
555 The spatial structures (a) matched four selected eigenvectors or MEMs, or (b) were simulated  
556 using geostatistical distributions with varying ranges (see Methods for details).

557

558 Fig. 4. MEM simulations. Comparison of the two approaches to compute MEMs in the three  
559 scenarios. By column, (a-d) display the adjusted  $R^2$  and (e-h) the number of  $MEM_{comp}$  or  
560 percentage of  $MEM_{red}$  required to model the various predefined spatial structures at different  
561 sampling thresholds (abscissa for S1 and S3: 10, 25, 50, 60, 70, 80, 90, 100%). Results are  
562 presented for the two most contrasted spatial structures (coarser spatial scale: filled circles,  
563 black; finer spatial scale: empty squares, grey). Solid lines correspond to  $MEM_{comp}$  while the  
564 dashed lines correspond to  $MEM_{red}$ . Intervals correspond to plus and minus the standard  
565 deviation estimated from the 100 replicated simulations. Results for all simulations are  
566 presented in Appendix B (Fig. B1 to B3).

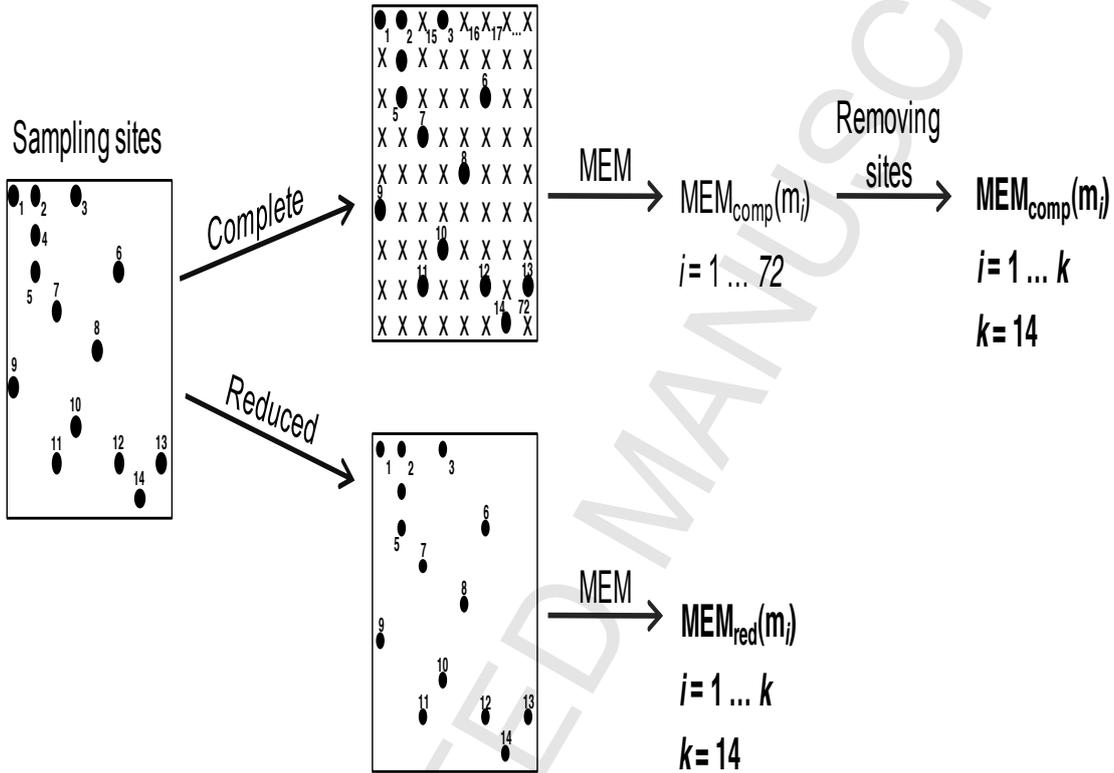
567

568 Fig. 5. Variogram simulations. Comparison of the two approaches to compute MEMs in the  
569 three scenarios. By column, (a-c) display the adjusted  $R^2$  and (d-f) the number of MEM  
570 required to model the various predefined spatial structures at different sampling thresholds  
571 (abscissa: 25, 50, 60, 70, 80, 90, 100%). Results are presented for the two most contrasted  
572 spatial structures (coarser spatial scale: filled circles, black; finer spatial scale: empty squares,  
573 grey). Solid lines correspond to  $MEM_{comp}$  while the dashed lines correspond to  $MEM_{red}$ .  
574 Intervals correspond to plus and minus the standard deviation estimated from the 100  
575 replicated simulations. Results for all simulations are presented in Appendix B (Fig. B4 to  
576 B6).

1

2 Figure 1

3



4

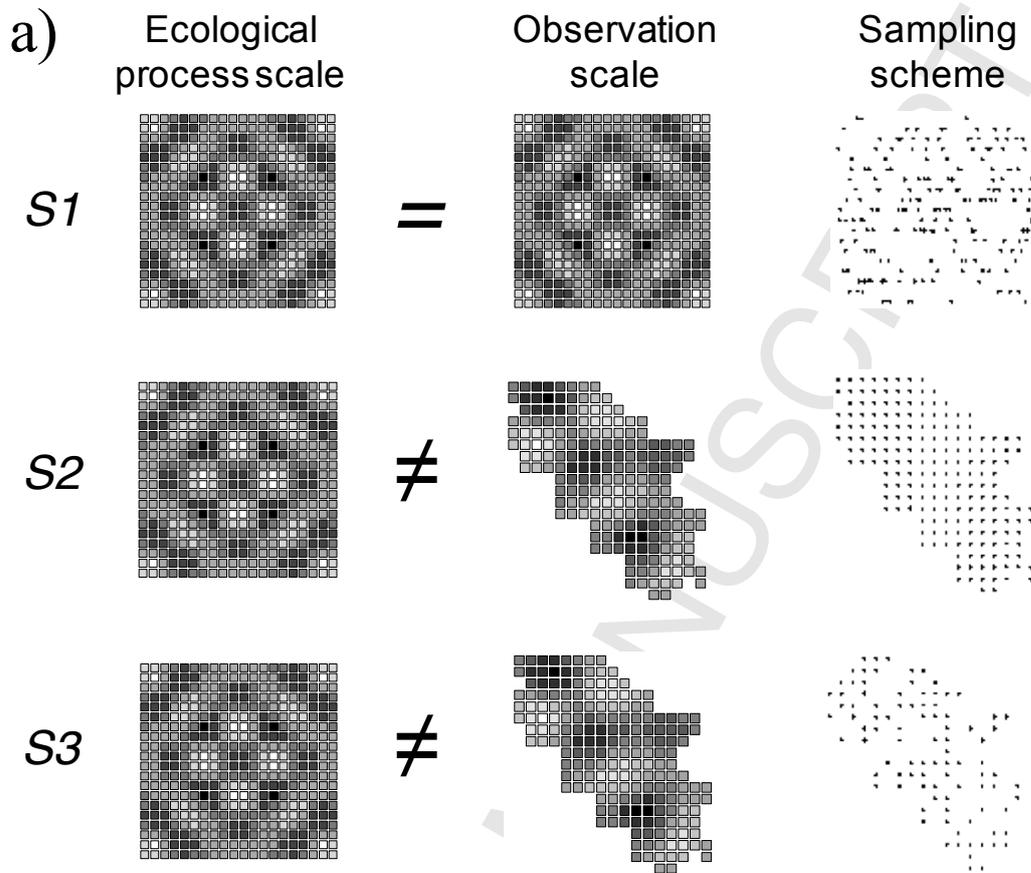
5

6

7

8

1 Figure 2

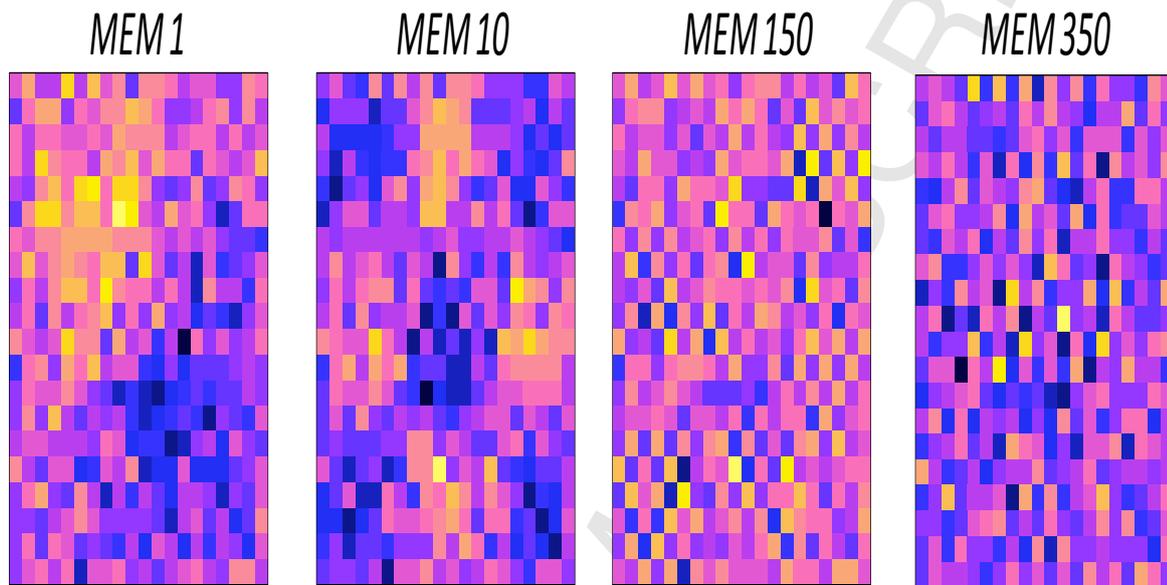


2

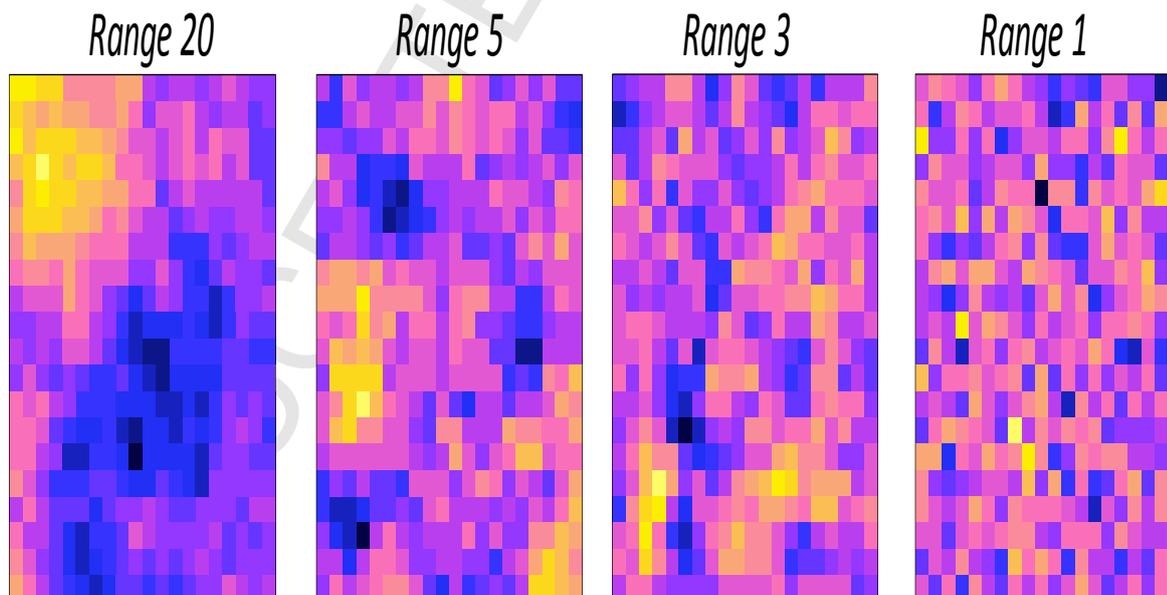
3

1 Figure 3

(a) Eigenvectors (MEM)



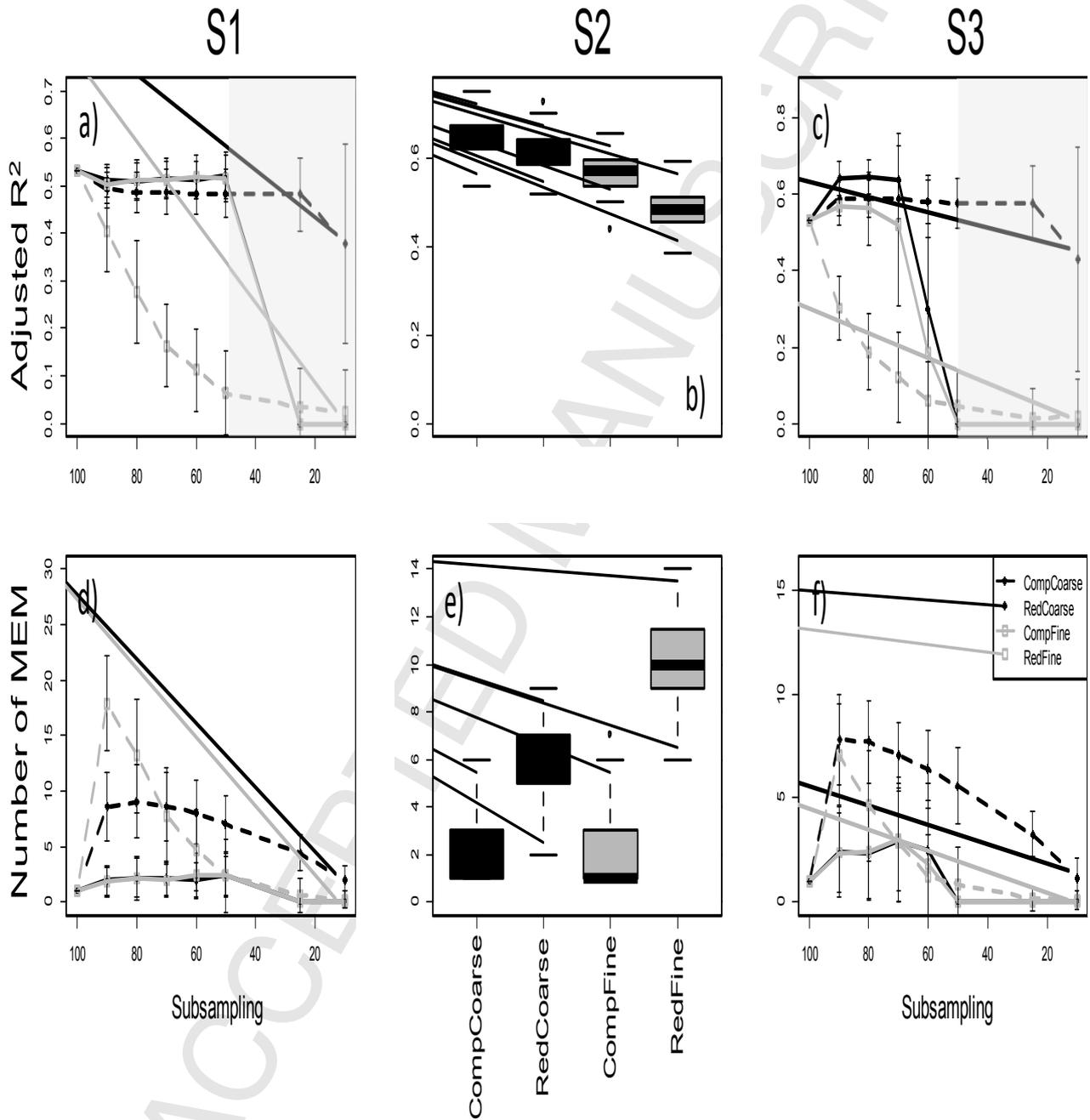
(b) Geostatistical distributions



2

1 Figure 4

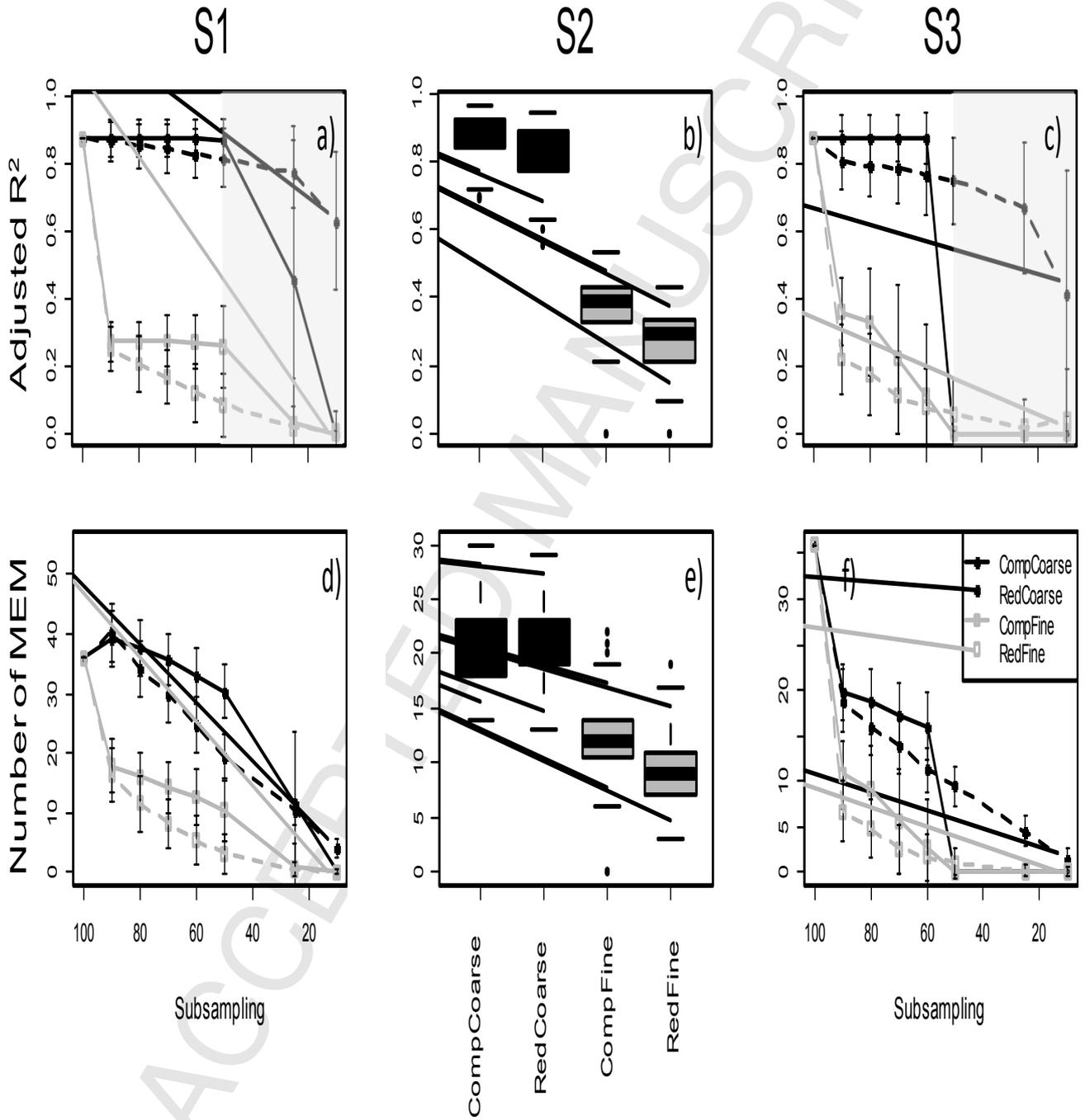
2



3

1 Figure 5

2



3