



**HAL**  
open science

## Spatial analyses of ecological count data: a density map comparison approach

Claire Lavigne, Benoit Ricci, Pierre Franck, Rachid R. Senoussi

► **To cite this version:**

Claire Lavigne, Benoit Ricci, Pierre Franck, Rachid R. Senoussi. Spatial analyses of ecological count data: a density map comparison approach. *Basic and Applied Ecology*, 2010, 11 (8), pp.734-742. 10.1016/j.baae.2010.08.011 . hal-02668695

**HAL Id: hal-02668695**

**<https://hal.inrae.fr/hal-02668695v1>**

Submitted on 31 May 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1 **Spatial analyses of ecological count data: a density map comparison approach**

2  
3  
4 Claire Lavigne<sup>a\*</sup>, Benoît Ricci<sup>a</sup>, Pierre Franck<sup>a</sup> and Rachid Senoussi<sup>b</sup>

5 <sup>a</sup> INRA, UR 1115, Plantes et Systèmes de culture Horticoles, F-84000 Avignon

6 <sup>b</sup> INRA, UR546 Biostatistique et Processus Spatiaux, INRA, F-84000 Avignon

7  
8  
9 Running title: Map comparison for spatial analysis

10  
11  
12 Number of words: 4562 (see letter to the editor)

13  
14  
15  
16  
17 \* Corresponding author: Tel.: 33 (0)4 32 72 26 66; fax: 33 (0)4 32 72 24 32

18 E-mail address: [claire.lavigne@avignon.inra.fr](mailto:claire.lavigne@avignon.inra.fr)

19  
20  
21  
22

**1 Abstract**

2 Analysing spatial patterns of population distributions may help to infer the decisive underlying  
3 ecological processes. Here we propose a method adapted to the spatial analysis of count data.  
4 Named MAPCOMP (MAP COMParison), it is based on the calculation of a formal distance, the  
5 Hellinger distance, between the density map of counts and the density map of sampling effort.  
6 Statistical tests of spatial homogeneity are based on count permutations across sampling sites and  
7 on valuable properties of the Hellinger distance. We assessed the efficiency of MAPCOMP by  
8 simulating different types and locations of clusters of individuals and compared its performance  
9 to the classical red-blue SADIE method, used as a reference. The two methods were also  
10 compared with respect to counts of codling moth larvae in orchards. Thanks to its better  
11 theoretical properties than SADIE, MAPCOMP was efficient in detecting spatial inhomogeneity  
12 when clusters were located on square or elongated spatial domains and more or less close to the  
13 edges, even for small sample sizes. It also appeared not very sensitive to edge effects. Another  
14 advantage of MAPCOMP is a bandwidth parameter that allows assessing the spatial extent of  
15 heterogeneity, if any.

**17 Zusammenfassung**

18 Die Analyse der räumlichen Muster in den Verteilungen von Populationen kann dazu beitragen,  
19 entscheidende grundlegende ökologische Prozesse abzuleiten. Hier stellen wir eine Methode zur  
20 räumlichen Analyse von Zählraten mit dem Namen MAPCOMP (MAP COMParison) vor. Sie  
21 basiert auf der Berechnung der Hellinger-Distanz zwischen der Dichteverteilung der  
22 Beobachtungen und der Dichteverteilung der Untersuchungsintensität.

23 Die statistischen Tests auf räumliche Homogenität basieren auf Permutationen der  
24 Beobachtungszahlen über die Probestellen und auf nützlichen Eigenschaften der Hellinger-  
25 Distanz. Wir ermittelten die Effizienz von MAPCOMP, indem wir verschiedene Typen und  
26 Anordnungen der Cluster von Individuen simulierten und die Leistung von MAPCOMP mit der

## Postprint

Version définitive du manuscrit publié dans / Final version of the manuscript published in : Basic and Applied Ecology, 2010, In Press, DOI: 10.1016/j.baae.2010.08.011

1 der klassischen red-blue SADIE-Methode verglichen. Die beiden Methoden wurden auch auf  
2 Fänge von Apfelwicklerlarven in Obstplantagen angewendet.

3 Dank seiner, verglichen mit SADIE, besseren theoretischen Eigenschaften konnte MAPCOMP  
4 effektiv räumliche Inhomogenität selbst bei geringen Individuenzahlen aufdecken, wenn die  
5 Cluster auf quadratischen oder rechteckigen Gittern verteilt und mehr oder weniger nah am Rand  
6 positioniert waren. MAPCOMP erschien nicht sehr empfindlich gegenüber Randeffekten zu sein.  
7 Ein weiterer Vorteil von MAPCOMP ist ein Bandbreitenparameter, der es erlaubt, die räumlichen  
8 Ausdehnung der Heterogenität, so vorhanden, abzuschätzen.

9  
10 **Keywords:** *Cydia pomonella*, clustering, Hellinger distance, heterogeneity, MAPCOMP, Monte  
11 Carlo permutations, spatial pattern, spatial statistics.

12

## 1 Introduction

2 Analysing spatial patterns of population distributions at various scales, from local patches to  
3 landscapes may help to infer the underlying ecological processes (McIntire & Fajardo 2009).  
4 Spatial statistics offer a number of tools for point pattern analysis (reviewed in Dale, Dixon,  
5 Fortin, Legendre, Myers et al. 2002; Perry, Liebhold, Rosenberg, Dungan, Miriti et al. 2002). A  
6 persistent issue is the statistical detection and characterization of spatial heterogeneity, such as  
7 gradients or clustering (Dale et al. 2002; Perry et al. 2002). Numerous methods exist for  
8 presence/absence data (e.g. Diggle, Gomez-Rubio, Brown, Chetwynd & Gooding 2007). We  
9 shall focus here on more general count data. When the form of the count spatial distribution is  
10 known a priori from knowledge of ecological processes (e.g. marked Poisson processes),  
11 parametric methods can be used to characterize spatial patterns. We consider more frequent  
12 situations where the form of the data distribution is not known. Methods then may follow two  
13 approaches. The first one is based on the counts of pairs of sampling points that exhibit similar  
14 count values at a given distance (e.g. variograms). This approach makes it possible to test a  
15 hypothesis of global clustering and determine some of its characteristics but does not allow the  
16 explicit mapping of heterogeneity. The second approach is based on the local departure of the  
17 spatial density of observations from an expected density. This approach makes it possible to map  
18 particular spatial patterns (e.g. geographic clusters of disease cases, Gay, Barnouin & Senoussi  
19 2006). SADIE (Perry 1998; Perry, Winder, Holland & Alston 1999) is its most popular  
20 representative despite some drawbacks (edge effects: Xu & Madden 2005). We propose a new  
21 method (MAPCOMP) based on the second approach. We know of no other method that is  
22 adapted to counts, accounts for heterogeneous sampling effort, has well defined and stable  
23 statistical properties and allows mapping of heterogeneity. Note, furthermore, that MAPCOMP  
24 would also apply to continuous positive data on continuous or discrete spatial supports.

25 We first present a brief modelling framework. Using simulated data, we then compare our  
26 method to the red-blue analysis of SADIE, used here as a reference. Cluster detection methods

1 are known to be sensitive to the forms of clusters and the sampling area, as well as to edge effect  
2 (Yamada 2003). We thus focus on these issues. Finally, we analyze the spatial distribution of  
3 codling moth larvae in eight orchards using both methods and we compare these results.

## 5 **Material and Method**

### 6 **Statistical methods**

7 Fig. 1 describes the MAPCOMP method. Symbols used in this study are summarized in Table 1.

#### 8 *General statistical framework*

9 The red-blue SADIE and the MAPCOMP methods are based on permutation tests. Permutation  
10 test methods consist in random permutations of indices  $s$  of an ordered data set  $\mathbf{X} = (X_s, s \in S)$ ,  
11 e.g. of counts observed at sites  $s$ . The index set  $S$  is structured via neighbourhood relations in a  
12 geographic domain  $D$ . To test a specific null hypothesis  $H_0$  (e.g. spatial homogeneity), one has to  
13 define an adequate test statistic,  $T(\mathbf{X})$ , and a specific subset,  $\Omega$ , of all data permutations. The  
14 choice of  $\Omega$  should reflect  $H_0$ . In the following, the set of all permutations of indices  $s$  was chosen  
15 as in SADIE, as we wished to keep the same  $H_0$  hypothesis. By its very definition,  $T(\mathbf{X})$  is  
16 designed to statistically behave differently under  $H_0$  and the alternative hypothesis  $H_1$ . The test  
17 relies on the comparison between the single statistic value calculated using the observed data  $T_{obs}$   
18 and the statistic values calculated using a large number of independently sampled data  
19 permutations.

#### 21 *Test statistic $T(\mathbf{X})$*

22 The four SADIE statistics are based on the minimal sum of distances that individuals have to  
23 move from site to site to reach a distribution such that either (1) the number of individuals  
24 (possibly fractional) is exactly the same over all sampling sites (statistics  $I_A, \bar{v}_i, \bar{v}_j$ ; associated  
25 probabilities  $P_A, P\bar{v}_i, P\bar{v}_j$ ), or (2) all individuals are situated on a single sampling site (statistic  
26  $J_A$ , probability  $Q_A$ ) (Perry 1998, Perry et al. 1999).

1 The suggested statistical test measures a formal (not spatial) distance between the density map of  
 2 counts and the density map of sampling effort. It increases with increasing inhomogeneity of  
 3 count distribution over sample sites. We thus needed (1) an estimation of density maps for  
 4 sample sites and observations and (2) a statistical test to assess the closeness of spatial densities.

5 *Comparing spatial densities*

6 We relied on the Hellinger distance between two probability densities  $p(s)$  and  $q(s)$  over a domain  
 7  $D$  (Gibbs & Su 2002). It is defined as follows:

$$8 \quad H(p, q) = \left( \frac{1}{2} \int_D \left( \sqrt{p(s)} - \sqrt{q(s)} \right)^2 ds \right)^{1/2} \quad s = (x, y) \in D \quad \text{eq. 1}$$

9 *Estimation of density maps*

10 Let  $S = \{s_i = (x_i, y_i) \quad i=1, \dots, \#S\}$  be the set of the geographic coordinates of  $\#S$  sites (Fig. 1).

11 Considering the data  $\mathbf{X}$  as a realization of a random measure whose theoretical normalized  
 12 intensity is denoted  $p(s)$ , we could estimate  $p$  at any location  $s$  via a local smoothing using a  
 13 probability density function, named kernel,  $K$ , that provides for every location  $s$  the weights of all  
 14 observations (Scott 1992):

$$15 \quad \hat{p}_h(s) = \sum_{s_i \in S} K_h(s, s_i) X_{s_i} / \sum_{s_i \in S} X_{s_i} \quad \text{eq. 2}$$

16 where, the function  $K_h(s, s_i)$  indexed by the positive parameter  $h$  (bandwidth parameter) is usually  
 17 the  $h$ -scaled and renormalized kernel  $K$  on the plane as follows:

$$18 \quad K_h(s, s_i) = K((s-s_i)/h)/h^2. \quad \text{eq. 3}$$

19 However, to account for the fact that bordering points intrinsically suffer from a lack of  
 20 neighbouring observation sites, we introduced the following edge correction:

$$21 \quad K_h(s, s_i) = K\left(\frac{s-s_i}{h}\right) / \int_D K\left(\frac{s-s_i}{h}\right) ds,$$

22 where  $D$  is the study domain.

23  
 24 We chose the following square supported kernel: for  $\mathbf{s}=(x, y)$ ,

1 
$$K(s) = C_1 \exp\left(-\left(\frac{1}{1-x^2} + \frac{1}{1-y^2}\right)\right) \text{ if } -1 \leq x, y \leq 1 \text{ and } 0 \text{ elsewhere} \quad \text{eq. 4}$$

2 and  $C_1$  is such that  $\int_{-1}^{+1} \int_{-1}^{+1} K(s) ds = 1$ , i.e.  $C_1 \simeq 20.3$ .

3 The estimated density  $\hat{p}_h(s)$  is thus a local average of counts of neighbouring sites within a  
4  $2h \times 2h$  square (Appendix A: Fig. 1).

5 The normalized sampling intensity  $q(s)$  is calculated similarly replacing  $\mathbf{X}$  by a vector of  
6 sampling effort. In the case of a uniform distribution of sampling sites, one could simply take  
7  $\hat{q}_h(g_{kl}) = 1/\#S$ .

8  
9 *The bandwidth h: a focusing tool*

10 In functional estimation, the bandwidth  $h$  is usually optimized to balance the global bias and  
11 variance of  $\hat{p}_h(\cdot)$  yielding either  $h_{opt} = C(\#S)^{-1/6}$  where  $C$  is a constant that can be chosen a priori  
12 (depending on domain area) or by a cross-validation method (Hardle, 1989). Here we used  $h$   
13 differently, giving it the role of a scale parameter to investigate the range at which sites sharing  
14 similar values clustered together. The tested  $h$  values should be such that the minimum value is  
15 larger than the distance between neighbouring sampling points (otherwise no smoothing occurs)  
16 and in the range of values at which dependencies among observations are expected. Exploring a  
17 range of values for  $h$  may reveal specific patterns at different spatial scales.

18  
19 *Computation procedure*

20 Statistics and maps were computed using R 2.1.1 (R Development Core Team 2005). The  
21 estimated densities were approximated over a regular grid  $\mathbf{G}$  of mesh  $\delta$  with nodes  $(g_{kl})$  covering  
22 the study domain  $D$  (e.g.  $\delta$  can be chosen as half the minimal distance between sampling sites).

23 The density of counts,  $p$ , was given values  $\hat{p}_h(g_{kl})$  whereas  $q$  estimating the sampling density



1 was given values  $\hat{q}_h(g_{kl})$  defined the same way as  $\hat{p}_h(g_{kl})$  except that all  $X_{s_i}$  were set to one.

2 The test statistic measuring the distance to homogeneity was thus:

$$3 \quad T_h(X) = \frac{\delta}{\sqrt{2}} \left( \sum_{g_{kl} \in G \cap D} \left( \sqrt{\hat{p}_h(g_{kl})} - \sqrt{\hat{q}_h(g_{kl})} \right)^2 \right)^{1/2} \quad \text{eq. 5}$$

## 5 Data sets

### 6 Simulated data

7 Sampling domains  $D$  were represented by square or rectangular grids with approximately the  
 8 same size (*i.e.* either 200 x 200 or 284 x 142). Sampling sites were chosen on the domains setting  
 9 the first site at coordinates (10, 10). Then coordinates of the sampling sites were incremented by  
 10 value 20 ((10, 30), (30, 10), (30, 30), etc.). With this procedure, the total of sampling sites  
 11 amounted to 100 for square domains and 98 for rectangular domains.

12 We then considered spatial distributions of either  $N=20$  or 100 sampled individuals,  
 13 corresponding to species with moderate and high population abundances. To study the effect of  
 14 the domain shape and the count distribution over this domain, individuals were randomly  
 15 distributed over sampling sites following five configurations: (1) uniform distribution in both the  
 16 square and the rectangular domain; (2 & 3) two Gaussian clusters located along the central  
 17 horizontal axis of the square domain either close to each other at sites (100, 75) and (100,125) or  
 18 farther away at sites (100, 50) and (100, 150); (4) two Gaussian clusters on one side of the square  
 19 domain, *i.e.* at sites (100, 25) and (100, 75); and (5) the same two clusters as in (3) but in the  
 20 middle of the rectangular domain, *i.e.* at sites (70,90) and (70,190) (Appendix A: Fig.2).

21 For Gaussian clusters, the value of standard deviation  $\sigma$ , hereafter denoted cluster radius, was  
 22 taken in set {12, 25, 35, 70}. For  $\sigma=12$  all sampled individuals were in sampling sites  
 23 neighbouring the centre of the cluster, hereafter denoted cluster focus, whereas for  $\sigma=70$  some  
 24 individuals could be distributed along the edge of domains (Fig. 2).

25 In total, we thus considered 32 clustered patterns (4 configurations x 2  $N$  values x 4  $\sigma$  values).

1 We then simulated ten replicate distributions for each pattern, resulting in a total of 320  
2 distributions.

### 3 4 *Codling moth sampling*

5 We sampled codling moth diapausing larvae in eight apple orchards located in south-eastern  
6 France (WGS84 : 43°46'27"N to 43°51'23"N and 4°51'12"E to 4°57'34"E). The mean area of  
7 orchards was  $0.40 \pm 0.07$  ha with values varying from 0.13 to 0.76 ha. Larvae were caught on 10-  
8 cm wide corrugated cardboard strip traps wrapped around tree trunks in July 2008 and collected  
9 the following October. Approximately 30 traps were distributed on a regular grid over each  
10 orchard.

### 11 12 **Testing homogeneity of spatial patterns**

13 We tested the homogeneity of spatial distributions of individuals over sampling points with the  
14 help of both the SADIE method (free software SADIEShell 1.22) and of MAPCOMP using R  
15 2.1.1 (R Development Core Team 2005). When using SADIE on simulated data, we used the  
16 maximum allowed number of Monte-Carlo simulations for significance test ( $k5psimul=153$ , *i.e.*  
17 5967 randomisations). We recorded the *P*-value of each of the three permutation tests associated  
18 with the three statistics  $I_A$ ,  $\bar{v}_i$ ,  $\bar{v}_j$  (*i.e.*  $P_{I_A}$ ,  $P_{\bar{v}_i}$ ,  $P_{\bar{v}_j}$ ) for each of the 320 spatial distributions, as  
19 Perry (1998) does not recommend considering the fourth statistic,  $J_A$ , in case of more than one  
20 cluster.

21 To use the MAPCOMP method, grids with a 2x2 mesh size over domains were chosen to  
22 calculate density maps. Then, for each bandwidth  $h$  taken in  $\{12, 25, 50, 90\}$ , the density maps of  
23 both sampling sites and count distributions were calculated. The smallest value  $h=12$  corresponds  
24 to a smoothing over only the nearest neighbours and the largest value,  $h=90$ , to a smoothing over  
25 about half the sampling sites. For each  $h$  value the Hellinger distance was calculated between the  
26 two density maps and the homogeneity hypothesis was tested using 10,000 permutations of

1 simulated counts over sampling sites. We recorded the lowest  $P$ -value among those obtained with  
2 the four  $h$  values, i.e. the result corresponding to the most appropriate  $h$  value,  $h_{Sig}$ . These  $P$ -  
3 values were used for comparison with SADIE.

4 We also tested the homogeneity of the spatial pattern in codling moth data using both SADIE and  
5 MAPCOMP. For the latter method, a 2 m x 2 m approximating grid was used and five  $h$  values  
6 ranging from a value encompassing only the nearest neighbouring traps to a value encompassing  
7 about  $1/4^{\text{th}}$  of the number of traps per orchard were chosen. The  $P$ -values for all  $h$  values in each  
8 orchard were then recorded.

9

## 10 **Results**

### 11 **Power comparison of the two tests**

12 As expected, the two methods provided non-significant results when individuals were randomly  
13 uniformly distributed over sites, whatever the shape of the domain or the sample size (not  
14 shown), meaning that neither method detected false positives. As results are generally very close  
15 in terms of  $P_A$ ,  $P\bar{v}_i$ ,  $P\bar{v}_j$  in SADIE, only results obtained with  $P_A$  are used for comparison.

16

#### 17 *Impact of the sample size on detection of spatial heterogeneity*

18 On a square plot, both methods detected the two central clusters efficiently, whether they were  
19 separated by 50 (configuration 2) or 100 (configuration 3) distance units when samples were  
20 large ( $N=100$ ) except the clusters with the smallest radius ( $\sigma=12$ ) that were not detected with  
21 SADIE (Figs 3A, 3C). For both methods, detection was also marginally worse in configuration 3  
22 than in configuration 2 when  $\sigma=70$ .

23 Decreasing the number of observed individuals to  $N=20$  resulted in less detection of clustering by  
24 both methods (Figs 3A, 3C). However, because MAPCOMP provided significant values in 64/80  
25 simulated patterns and only poorly detected clusters of the largest radius ( $\sigma=70$ ), its detection  
26 capacity can be deemed good. In contrast, SADIE provided significant results only in 7/80

1 patterns that corresponded to clusters with radius  $\sigma=25$  and  $\sigma=35$ .

2 Obviously, not all  $h$  values were equally efficient with the MAPCOMP method whatever  $N$ .  $h_{Sig}$ ,  
3 *i.e.* the  $h$  value providing more significant results for the ten replicates, tended to increase with  
4 increasing radius width  $\sigma$  (not shown).

#### 6 *Impact of the domain shape on detection of spatial heterogeneity*

7 The elongated shape of the sampling domain (configuration 5) slightly modified the behaviour of  
8 MAPCOMP (Figures 3C vs 3D): results were not modified for  $N=100$ , but for  $N=20$ ,  
9 MAPCOMP somewhat less efficiently detected clusters with either small or large radii (from  
10 30/40 detections in square domains to 26/40). In contrast, the detection ability of SADIE largely  
11 decreased for  $N=100$  (from 26/40 to 4/40) and remained close to 0 for  $N=20$ .

#### 13 *Impact of cluster distance to border*

14 Setting the foci of the clusters closer to a domain edge (configuration 4) had a marginal effect on  
15 inhomogeneity detection with MAPCOMP (Figure 3A vs 3B). Thanks to the border effect  
16 correction, all clustered patterns were still detected with certainty for  $N=100$ . Detection  
17 probability only increased for  $\sigma=70$  and  $N=20$  as compared with the reference situation of centred  
18 clusters (configuration 2). On the contrary, the SADIE method was very sensitive to border  
19 effects: in this case, the method detected inhomogeneity for  $\sigma=12$  and  $N=100$ , and the probability  
20 of detecting a cluster rose from 6/40 to 37/40 for  $N=20$ .

#### 22 **Codling moth data**

23 Clustering was detected by MAPCOMP in 3 out of the 8 tested orchards (Table 2). These  
24 orchards were also pinpointed using SADIE, testing for regularity (for two orchards) and  
25 crowding (for the third one, orchard #G). We used this latter test as data suggested the presence  
26 of a single cluster. The other five orchards showed no significant departure from homogeneity

1 using both methods. Interestingly, in the three heterogeneous orchards, the density of codling  
2 moth larvae was higher at the edges (Figure 4).

## 4 **Discussion**

5 We have presented a new method, MAPCOMP, for detecting the inhomogeneity of spatial  
6 patterns in count data. We compared it to SADIE, a reference method for ecologists to detect the  
7 spatial heterogeneity in count data (*e.g.* Thomas, Parkinson, Griffiths, Garcia & Marshall 2001;  
8 Schellhorn, Bellati, Paull & Maratos 2008). Our results confirm that SADIE is efficient in many  
9 standard situations (*e.g.* regularly shaped domains, large sample size). However, MAPCOMP is  
10 as efficient as SADIE for large sample sizes and wide enough cluster patterns. Moreover, it better  
11 detected clustered patterns for small sample sizes and clusters with a small radius. MAPCOMP  
12 also appeared less sensitive than SADIE to the shape of the sampled domain and to edge effects.  
13 These differences could be explained by the different nature of the test statistic used. In SADIE,  
14 the test statistic measures the total distance that individuals have to move to reach a regular  
15 distribution of individuals. Consequently, the SADIE method has the drawback that the value of  
16 the statistic associated with observations strongly depends on large distances between sampling  
17 points. For instance, its value for a cluster at the short side of an elongated domain will be much  
18 larger than the value assessed for this same cluster located in the central part of this same domain.  
19 The associated observed *P*-value will thus be smaller since the distribution of the distance under  
20  $H_0$  is the same in both cases. This sensitivity of SADIE to localization of clusters has already  
21 been pointed out (Xu & Madden 2005). This explains the difficulty of SADIE to detect small  
22 clusters in elongated areas. However, it must not be forgotten that, in case of a single cluster, the  
23 SADIE method proposes another test for detection of clustering based on the minimum distance  
24 that individuals would have to move to be all clustered on a single sampling site. This test should  
25 be very efficient at detecting clusters of small radii. Conversely, the better efficacy of  
26 MAPCOMP for detecting heterogeneity in small samples is probably due to the use of a

1 normalized version of smoothing technique for densities. Its better efficacy with respect to edge  
 2 effects and shape of the sampled area is certainly due to the insensitivity of the Hellinger distance  
 3 to size and shape of the observation domain.

4 Further, MAPCOMP, as SADIE, not only detects spatial heterogeneity but also provides maps of  
 5 lower and higher than expected individual densities, i.e. patches and gaps, based on the difference  
 6  $\hat{p}(g_{kl}) - \hat{q}(g_{kl})$ . These maps help to gain insight into the factors underlying the perceived  
 7 inhomogeneity.

8 Finally, from a completely different perspective, SADIE provides an association statistical test  
 9 (Perry & Dixon, 2002) that allows the comparison of the spatial distributions of a single species  
 10 at different dates (*e.g.* Blackshaw & Vernon 2006) or of two species at the same date (*e.g.* Pearce  
 11 & Zalucki 2006). As MAPCOMP is based on a formal distance between spatial densities, it could  
 12 readily be used in these contexts by replacing the sampling density with the density of the second  
 13 species or previous observations of the species under study (for a first application, see Debras,  
 14 Senoussi, Rieux, Buisson & Dutoit 2008).

15  
 16 From a more general point of view, MAPCOMP has further desirable properties for the statistical  
 17 analysis of count data because of its use of the Hellinger distance: (1)  $H(p,q)$  is a true metric, that  
 18 is non-negative, null only if  $p=q$ , symmetric in  $p(s)$  and  $q(s)$  (so that there is no need to define a  
 19 reference population) and satisfies the triangular inequality  $H(p_1,p_2) \leq H(p_1,q) + H(p_2,q)$ ,  
 20 allowing for instance for the comparison (and partial ordering) between different species  
 21 distributions. (2)  $H$  has bounded values in  $[0,1]$ , it thus avoids formally large and asymmetric  
 22 values and can be considered as well scaled for value comparisons. (3)  $H$  behaves smoothly with  
 23 respect to data aggregation. For example, if several distinct study regions  $D_j$  with respective  
 24 densities  $p_j$  and  $q_j$  were aggregated together and given weights  $\alpha_j \geq 0, \sum_j \alpha_j = 1$  (*e.g.*  
 25  $\alpha_j = \text{area}(D_j) / \sum_k \text{area}(D_k)$ ), then the squared Hellinger distance of the resultant densities  
 26 would simply equal the mean of squared Hellinger distances over the areas, *i.e.*

1  $H^2\left(\sum_j \alpha_j p_j, \sum_j \alpha_j q_j\right) = \sum_j \alpha_j H^2(p_j, p_j)$ . (4)  $H$  has a robust and non-erratic behaviour when

2 adding new observations in the given domain (temporal sampling or species grouping) (Beran,  
3 1977): adding new data within a given domain can only lead to a distance decrease that will  
4 converge to the 'true'  $H$  distance between densities as the number of data points increases. For  
5 example, the grouping of distinct subpopulations  $p_j$  which are given weights  $\alpha_j \geq 0, \sum_j \alpha_j = 1$   
6 (e.g.  $\alpha_j =$  proportion of biomass of population  $j$ ) observed within a given domain gives

7  $H^2\left(\sum_j \alpha_j p_j, q\right) \leq \sum_j \alpha_j H^2(p_j, q)$ . Another advantage of MAPCOMP is the fact that continuous

8 variables can be used instead of count data without the model undergoing any modifications.  
9 These variables could be either abiotic factors (e.g. temperature, altitude, etc.) or community or  
10 population parameters (e.g. species diversity, biomass, growth, etc.). Finally, the last but not the  
11 least asset rests on the possible use of the bandwidth parameter as a scale analysis tool,  
12 specifically to test the spatial extent of heterogeneity, if any.

13 A very large number of statistical methods are available to ecologists for the detection of spatial  
14 inhomogeneity. Using the Dale et al. (2002) classification for these methods indicates that  
15 MAPCOMP is not redundant with any of those taken into consideration in their study. Being  
16 based on continuous density functions and varying bandwidths, MAPCOMP bears relationships  
17 to the wavelets methods, but is not based on many fitting parameter estimation. It also shares  
18 similarities with the block and quadrat variance methods as it attempts to detect larger than  
19 expected variance in adjacent data using windows of varying sizes, but it is not based on a direct  
20 measure of local variances.

21 While the comparison of methods on simulated data tends to favour the MAPCOMP method, it  
22 should be noted that MAPCOMP and SADIE behaved rather similarly in the case of our field  
23 study dataset. However, the cluster in orchard  $G$  was detected by SADIE only when using the  
24 crowding test which is not the most recommended one (Perry et al. 1999). Interestingly, the

1 detection of spatial heterogeneity did not depend on sample size or orchard area. From the  
2 ecological point of view, the heterogeneity of codling moth distribution in all three cases was due  
3 to higher than expected numbers of larvae on orchard edges not bordered by a hedgerow or  
4 another orchard (Fig. 4).

5 To conclude, both methods proved to be efficient at detecting spatial heterogeneity for large  
6 sample sizes with regular domains. Moreover, MAPCOMP efficiently detects inhomogeneity on  
7 small sample sizes and less regular domains. This could be particularly useful in the case of  
8 conservation biology of rare species as well as in the case of agricultural pests where population  
9 densities are expected to be low and population habitats may be geometrically intricate.

10

### 11 **Acknowledgments**

12 We are grateful to Jean-François Toubon for contacts with farmers, to members of EPI for  
13 counting codling moth larvae, to farmers for access to their orchards and to regional experimental  
14 station La Pugère. B.R. was funded by INRA and PACA district. This research was partly funded  
15 by the ECOGER project “Ecco des vergers”. Tanis Plant and D.G. Biron edited a draft of this  
16 paper in English.

17

### 18 **Appendix A: Supplementary Material**

19 The online version of this article contains additional supplementary data. Please visit XXXX.

20



1 **References**

- 2
- 3 Beran, R. (1977). Minimum Hellinger distance estimates for parametric models. *The annals of*
- 4 *Statistics*, 5, 445-463.
- 5 Blackshaw, R.P., & Vernon, R.S. (2006). Spatiotemporal stability of two beetle populations in
- 6 non-farmed habitats in an agricultural landscape. *Journal of Applied Ecology*, 43, 680-689.
- 7 Dale, M.R.T., Dixon, P., Fortin, M.J., Legendre, P., Myers, D.E., & Rosenberg, M.S. (2002).
- 8 Conceptual and mathematical relationships among methods for spatial analysis. *Ecography*,
- 9 25, 558-577.
- 10 Debras, J.F., Senoussi, R., Rieux, R., Buisson, E., & Dutoit, T. (2008). Spatial distribution of an
- 11 arthropod community in a pear orchard (southern France) - Identification of a hedge effect.
- 12 *Agriculture Ecosystems & Environment*, 127, 166-176.
- 13 Diggle, P.J., Gomez-Rubio, V., Brown, P.E., Chetwynd, A.G., & Gooding, S. (2007). Second-
- 14 order analysis of inhomogeneous spatial point processes using case-control data. *Biometrics*,
- 15 63, 550-557.
- 16 Gay, E., Barnouin, J., & Senoussi, R. (2006). Spatial and Temporal Patterns of Herd Somatic Cell
- 17 Score in France. *Journal Dairy Science*, 89, 2487-2498.
- 18 Gibbs, A.L., & Su, F.E. (2002). On choosing and bounding probability metrics. *International*
- 19 *Statistical Review*, 70, 419-435.
- 20 Härdle, W. (1989) *Applied non parametric regression*, Cambridge University Press.
- 21 McIntire, E.J.B., & Fajardo, A. (2009). Beyond description: the active and effective way to infer
- 22 processes from spatial patterns. *Ecology*, 90, 46-56.
- 23 Pearce, S., & Zalucki, M.P. (2006). Do predators aggregate in response to pest density in
- 24 agroecosystems? Assessing within-field spatial patterns. *Journal of Applied Ecology*, 43, 128-
- 25 140.
- 26 Perry, J.N. (1998). Measures of spatial pattern for counts. *Ecology*, 79, 1008-1017.
- 27 Perry, J.N., & Dixon, P.M. (2002). A new method to measure spatial association for ecological

- 1 count data. *Ecoscience*, 9, 133-141.
- 2 Perry, J.N., Liebhold, A.M., Rosenberg, M.S., Dungan, J., Miriti, M., Jakomulska, A., & Citron-  
3 Pousty, S. (2002). Illustrations and guidelines for selecting statistical methods for quantifying  
4 spatial pattern in ecological data. *Ecography*, 25, 578-600.
- 5 Perry, J.N., Winder, L., Holland, J.M., & Alston, R.D. (1999). Red-blue plots for detecting  
6 clusters in count data. *Ecology Letters*, 2, 106-113.
- 7 Schellhorn, N.A., Bellatib, J., Paullb, C.A., & Maratos, L. (2008). Parasitoid and moth movement  
8 from refuge to crop. *Basic and Applied Ecology*, 9, 691-700.
- 9 Scott, D.W. (1992). *Multivariate Density Estimation. Theory, Practice and Visualization*. New  
10 York: Wiley.
- 11 Thomas, C.F.G., Parkinson, L., Griffiths, G.J.K., Garcia, A.F., & Marshall, E.J.P. (2001).  
12 Aggregation and temporal stability of carabid beetle distributions in field and hedgerow  
13 habitats. *Journal of Applied Ecology*, 38, 100-116.
- 14 Xu, X.M., & Madden, L.V. (2003). Considerations for the use of SADIE statistics to quantify  
15 spatial patterns. *Ecography*, 26, 821-830.
- 16 Yamada, I., & Rogerson, P. (2003). An Empirical Comparison of Edge Effect Correction  
17 Methods Applied to K -function Analysis. *Geographical Analysis*, 35, 97-109.
- 18

1 **Table 1:** Symbols used in this study.  
 2  
 3

Symbol	Meaning
$\mathbf{X} = (X_s, s \in S)$	A set of ordered observations (Data)
$X_s$	Observation at site $s$
$S$	Set of observation sites
$\#S$	Number of sites in $S$
$s = (x,y)$	Observation site $s$ identified by its coordinates $(x,y)$
$D$	Geographical domain over which observations are performed
$\Omega$	A subset of permutations of sites used for testing the null hypothesis $H_0$ among all possible permutations
$T(\mathbf{X})$	The test statistic calculated from the ordered data set $\mathbf{X}$
$T_{obs}$	The test statistic calculated from the actual observations
$I_A, \bar{v}_i, \bar{v}_j, J_A$	The four test statistics from SADIE
$P_A, P\bar{v}_i, P\bar{v}_j, Q_A$	The four P-values associated to the four above statistics in SADIE
$p$	Probability density function of observations
$q$	Probability density function of sampling effort
$H(p,q)$	Hellinger distance between $p$ and $q$ (integral over domain $D$ )
$K$	The basic kernel smoothing function ( any probability density function)
$h$	Bandwidth parameter for smoothing
$K_h$	The $h$ -scaled and renormalized smoothing kernel $K$ ( $K_l=K$ )
$\mathbf{G}$	Grid over $D$ used to compute the integral $H$
$\delta$	Mesh size of $\mathbf{G}$
$g_{kl}$	Node of coordinates $(k,l)$ on $\mathbf{G}$
$\hat{p}_h(g_{kl})$	Estimation of $p$ using $K_h$ at node $g_{kl}$
$\hat{q}_h(g_{kl})$	Estimation of $q$ using $K_h$ at node $g_{kl}$
$T_h(\mathbf{X})$	Computed test statistic in MAPCOMP for bandwidth $h$ and data set $\mathbf{X}$
$N$	Number of sampled individuals in simulated data sets
$\sigma$	Radius of simulated clusters

4  
 5  
 6

Manuscrit d'auteur / Author manuscript

1 **Table 2:** Characteristics of sampled orchards and tests of codling moth spatial homogeneity.  
 2  $h_i$ :  $i^{\text{th}}$  bandwidth parameter used for test and  $P$ : P-value using MAPCOMP;  $P_A, P\bar{v}_i, P\bar{v}_j, Q_A$ : P-values for tests in SADIE. Orchards exhibiting  
 3 significant heterogeneity are in bold.

Orchard Id.	Number of traps	Area (ha)	Codling moths/trap mean $\pm$ se	MAPCOMP					SADIE								
				$h_1$ P	$h_2$ P	$h_3$ P	$h_4$ P	$h_5$ P	$P_A$	$P\bar{v}_i$	$P\bar{v}_j$	$Q_A$					
A	32	0.32	0.59 $\pm$ 0.20	<i>11</i>	<i>13</i>	<i>15</i>	<i>18</i>	<i>20</i>	0.94	0.95	0.96	0.96	0.95	0.58	0.61	0.54	0.77
B	33	0.76	12.48 $\pm$ 1.64	<i>19</i>	<i>21</i>	<i>23</i>	<i>25</i>	<i>28</i>	0.32	0.28	0.27	0.25	0.20	0.21	0.21	0.39	0.36
C	34	0.13	4.30 $\pm$ 0.58	<i>7</i>	<i>9</i>	<i>11</i>	<i>13</i>	<i>15</i>	0.76	0.62	0.61	0.53	0.46	0.65	0.83	0.83	0.93
D	32	0.19	1.06 $\pm$ 0.18	<i>8</i>	<i>11</i>	<i>14</i>	<i>17</i>	<i>20</i>	0.29	0.14	0.18	0.37	0.57	0.53	0.50	0.56	0.11
E	31	0.45	9.35 $\pm$ 1.15	<i>16</i>	<i>20</i>	<i>24</i>	<i>28</i>	<i>32</i>	0.95	0.93	0.92	0.92	0.92	0.97	0.945	0.905	0.43
F	<b>30</b>	<b>0.33</b>	<b>5.56<math>\pm</math>1.11</b>	<i>12</i>	<i>15</i>	<i>18</i>	<i>21</i>	<i>23</i>	<b>8.10<sup>-4</sup></b>	<b>4.10<sup>-4</sup></b>	<b>4.10<sup>-4</sup></b>	<b>&lt;10<sup>-4</sup></b>	<b>&lt;10<sup>-4</sup></b>	<b>2.10<sup>-4</sup></b>	<b>0.003</b>	<b>5.10<sup>-4</sup></b>	0.29
G	<b>33</b>	<b>0.70</b>	<b>2.06<math>\pm</math>0.42</b>	<i>19</i>	<i>21</i>	<i>24</i>	<i>26</i>	<i>28</i>	<b>0.012</b>	<b>0.008</b>	<b>0.009</b>	<b>0.008</b>	<b>0.009</b>	0.105	0.159	0.115	<b>0.021</b>
H	<b>30</b>	<b>0.27</b>	<b>0.87<math>\pm</math>0.25</b>	<i>10</i>	<i>12</i>	<i>13</i>	<i>14</i>	<i>15</i>	<b>0.015</b>	<b>0.009</b>	<b>0.007</b>	<b>0.006</b>	<b>0.008</b>	<b>0.016</b>	<b>0.019</b>	<b>0.048</b>	<b>0.005</b>

1 **Figure legends:**

2 **Fig. 1:** Schematic representation of the MAPCOMP method to analyze spatial patterns of  
3 population distributions for a given species. The density map of observed counts  $p_{obs}^h$  is derived  
4 from the smoothing of the observed counts using kernel  $K_h$  on grid  $G$ . The density map of  
5 sampling  $q_h$  is derived from locations of sampling sites using the same kernel and grid. The test  
6 statistic  $T_{obs}(X_s)$  based on the Hellinger distance between these two maps is compared to each of  
7 those obtained between the same sampling density map  $q_h$  and each map calculated from  
8 permuted counts  $T(X_{permut})$ .

9  
10 **Fig. 2:** Examples of simulated cluster: A)  $N=100$ ,  $\sigma=25$ , configuration 2, B)  $N=100$ ,  $\sigma=25$ ,  
11 configuration 3, C)  $N=100$ ,  $\sigma=25$ , configuration 4, D)  $N=100$ ,  $\sigma=25$ , configuration 2 (See  
12 Appendix A: Fig. 2 for other examples).

13  
14 **Fig. 3:** Detection probability of clustered patterns in simulated data, i.e. number of cases (out of  
15 10) where the  $P$ -value is below 0.05 as a function of  $\sigma$ , the radius of the cluster. White bars:  
16 MAPCOMP, grey bars: SADIE.

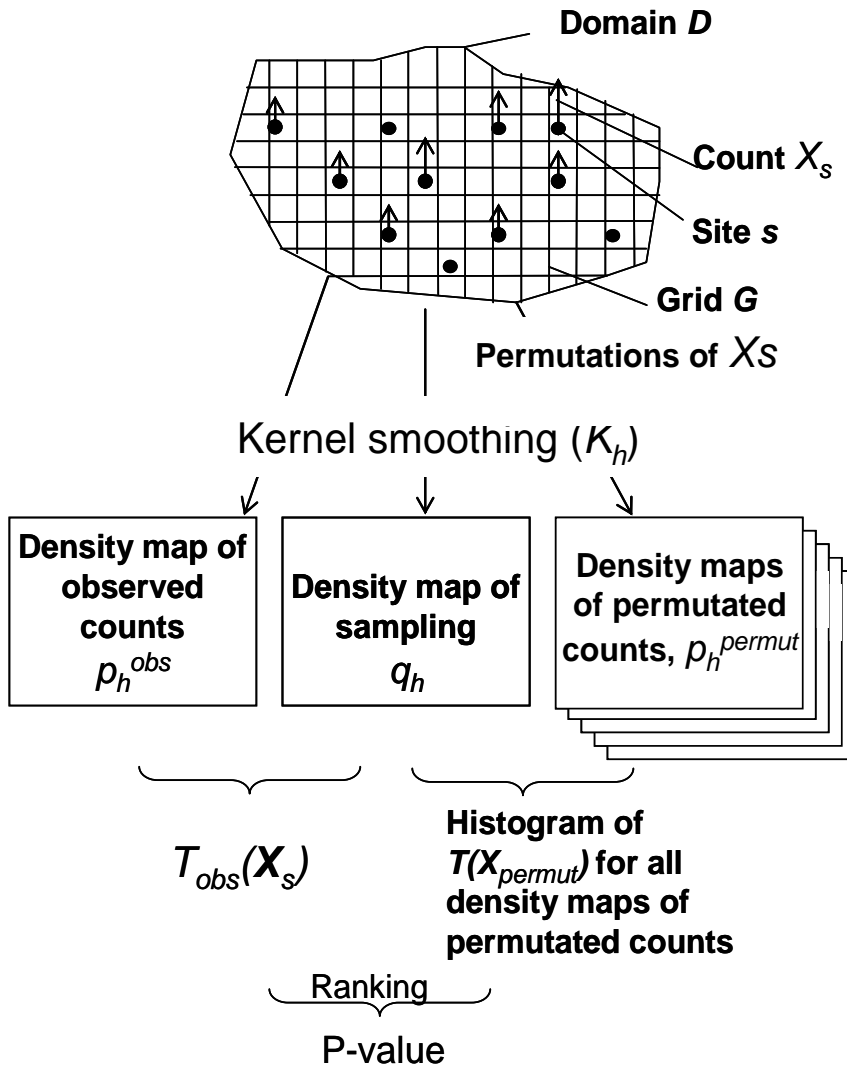
17  
18 **Fig. 4:** Density maps of the three orchards on which we detected clustering of codling moths  
19 larvae.  $h$ : bandwidth used for MAPCOMP, dots= traps, thick lines=hedgerows or neighbouring  
20 orchard or forest edge.

21

22

1  
2

Figure 1

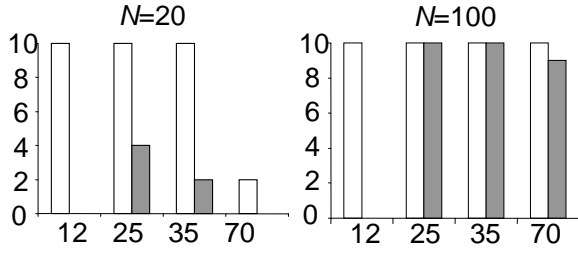


3  
4

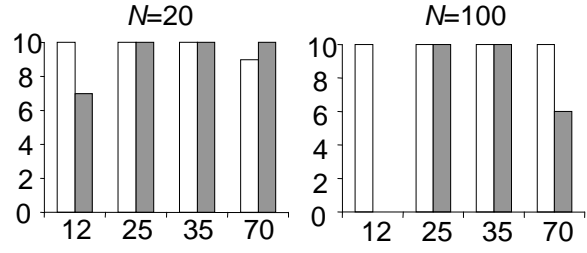


1 **Figure 3**

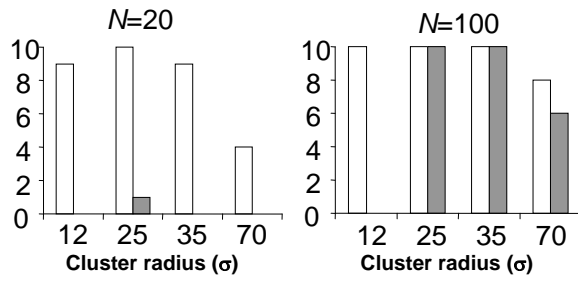
**A) Configuration 2: square  $D$ , close clusters**



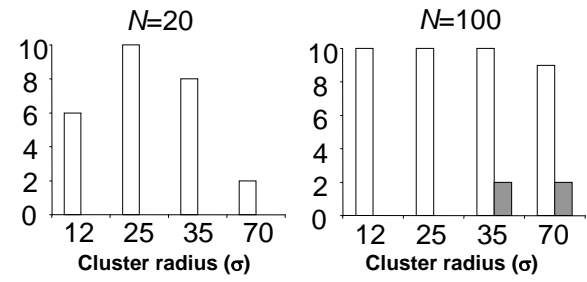
**B) Configuration 4: square  $D$ , clusters close to edge**



**C) Configuration 3: square  $D$ , far clusters**

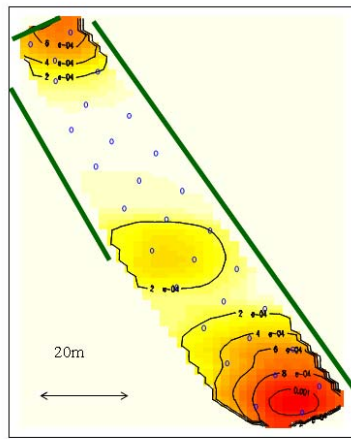


**D) Configuration 5: rectangular  $D$ , far clusters**

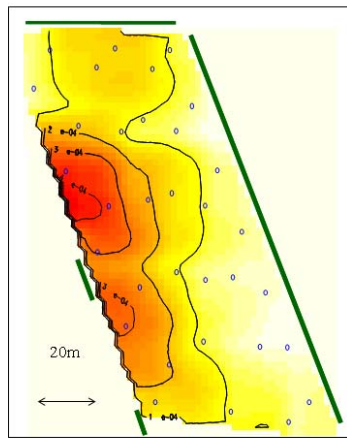




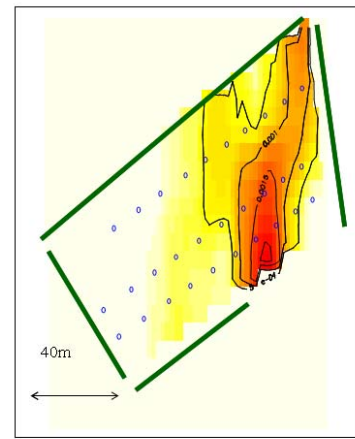
1 **Figure 4**



9 Orchard F,  $h=18$



10 Orchard G,  $h=24$

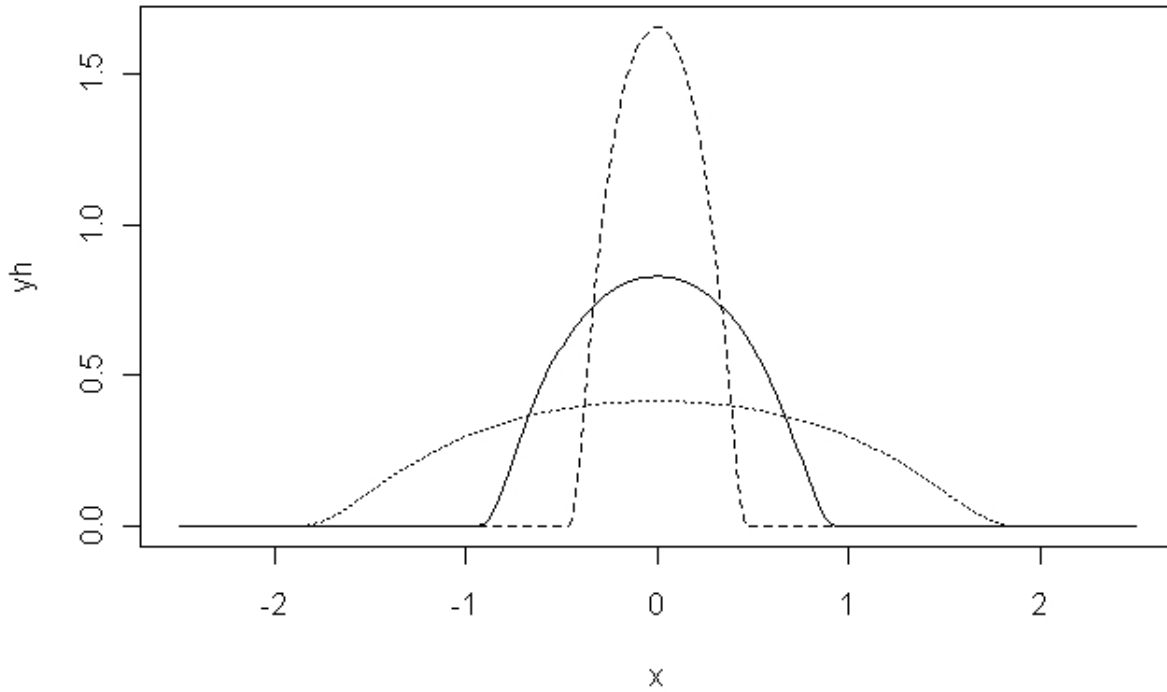


13 Orchard H,  $h=14$

1 **Appendix A:**

2 **Figure 1:** 2D representation of the kernel used for smoothing: —  $h=1$ , .....  $h=2$ , - -  $h=0.5$ .

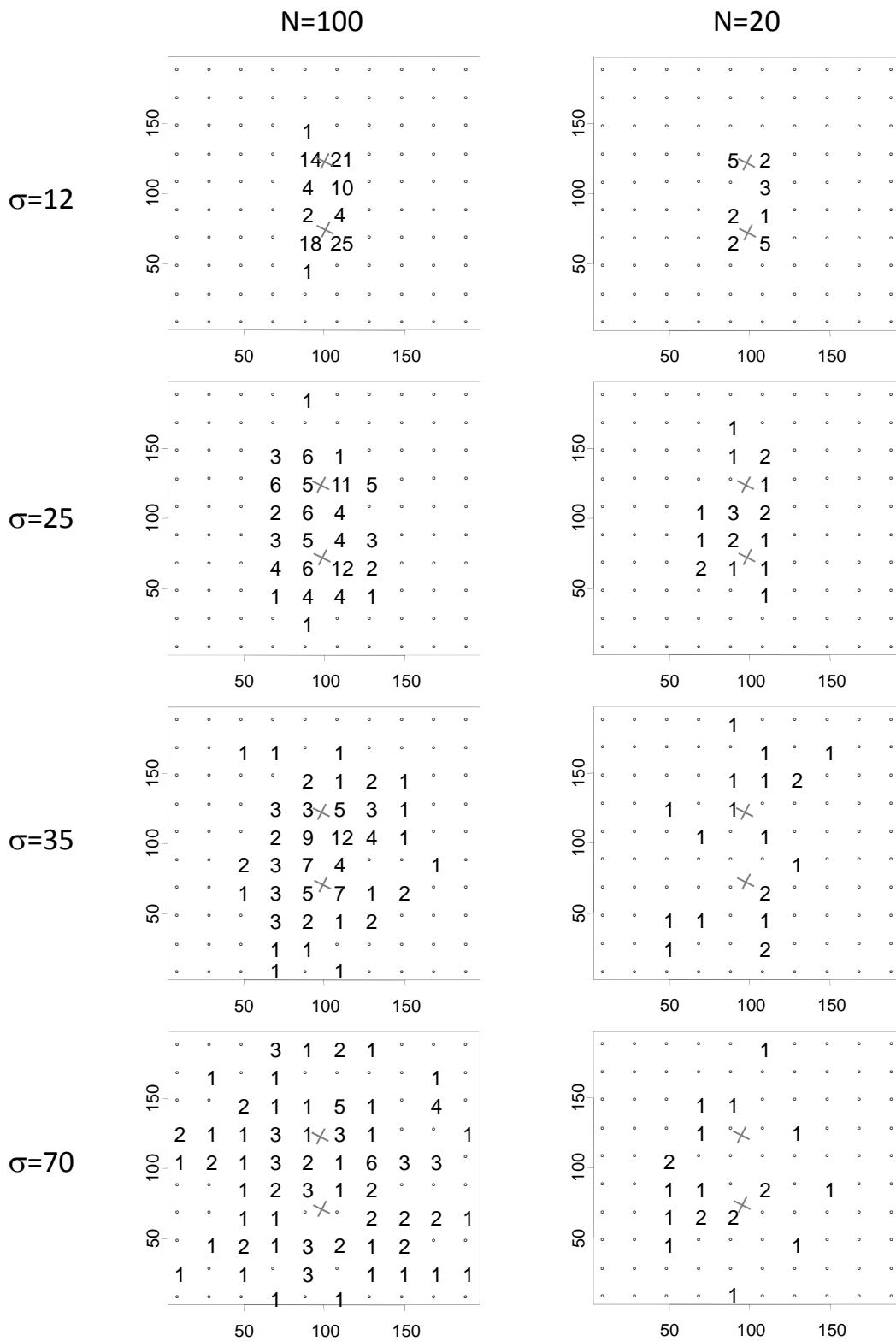
3



4  
5

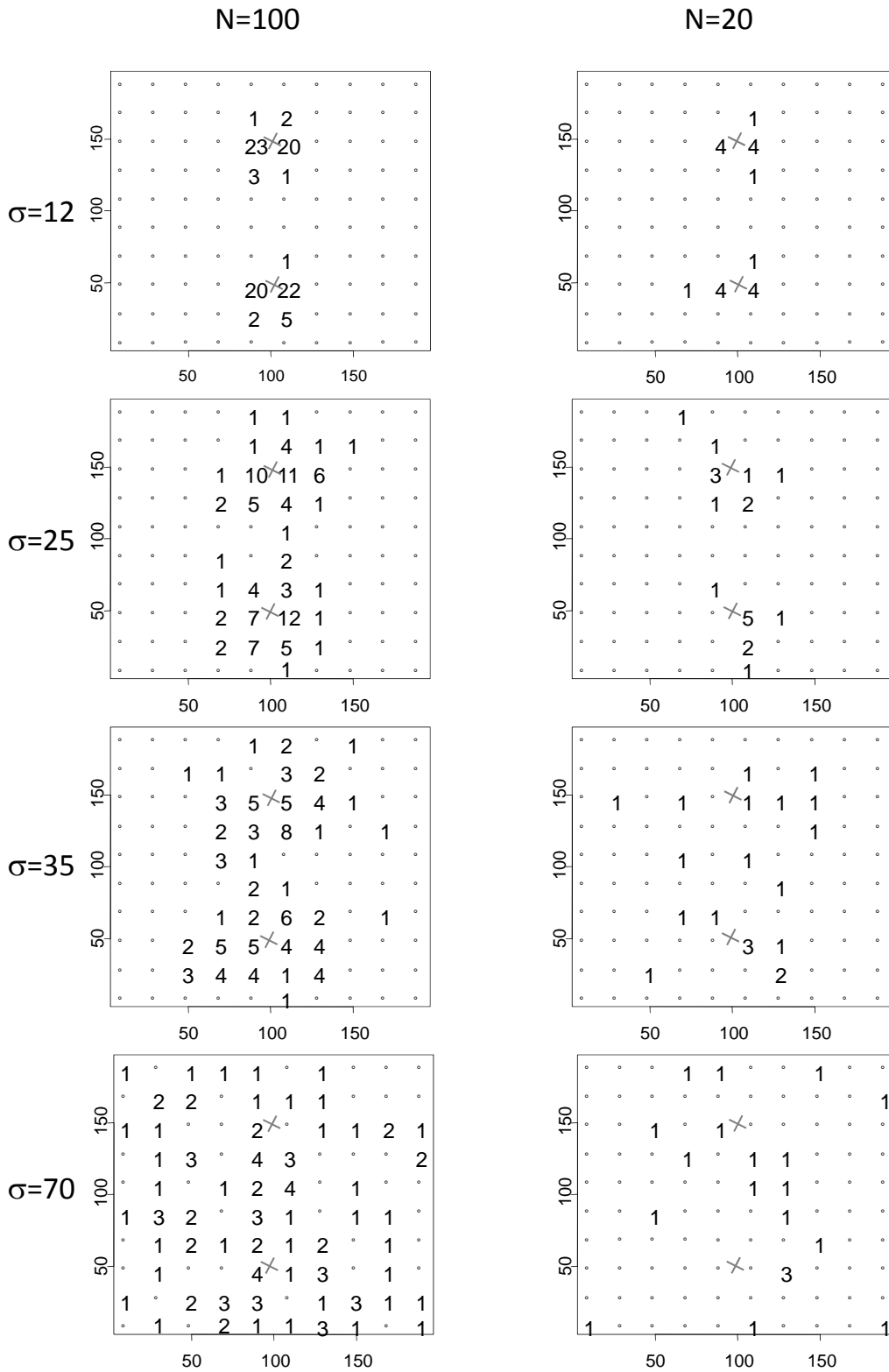
1 **Figure 2:** One example of each of the 32 different simulated clustered patterns.  $N$ =number of  
 2 sampled individuals,  $\sigma$ =radius of simulated clusters.

Configuration 2: Square D, close central clusters



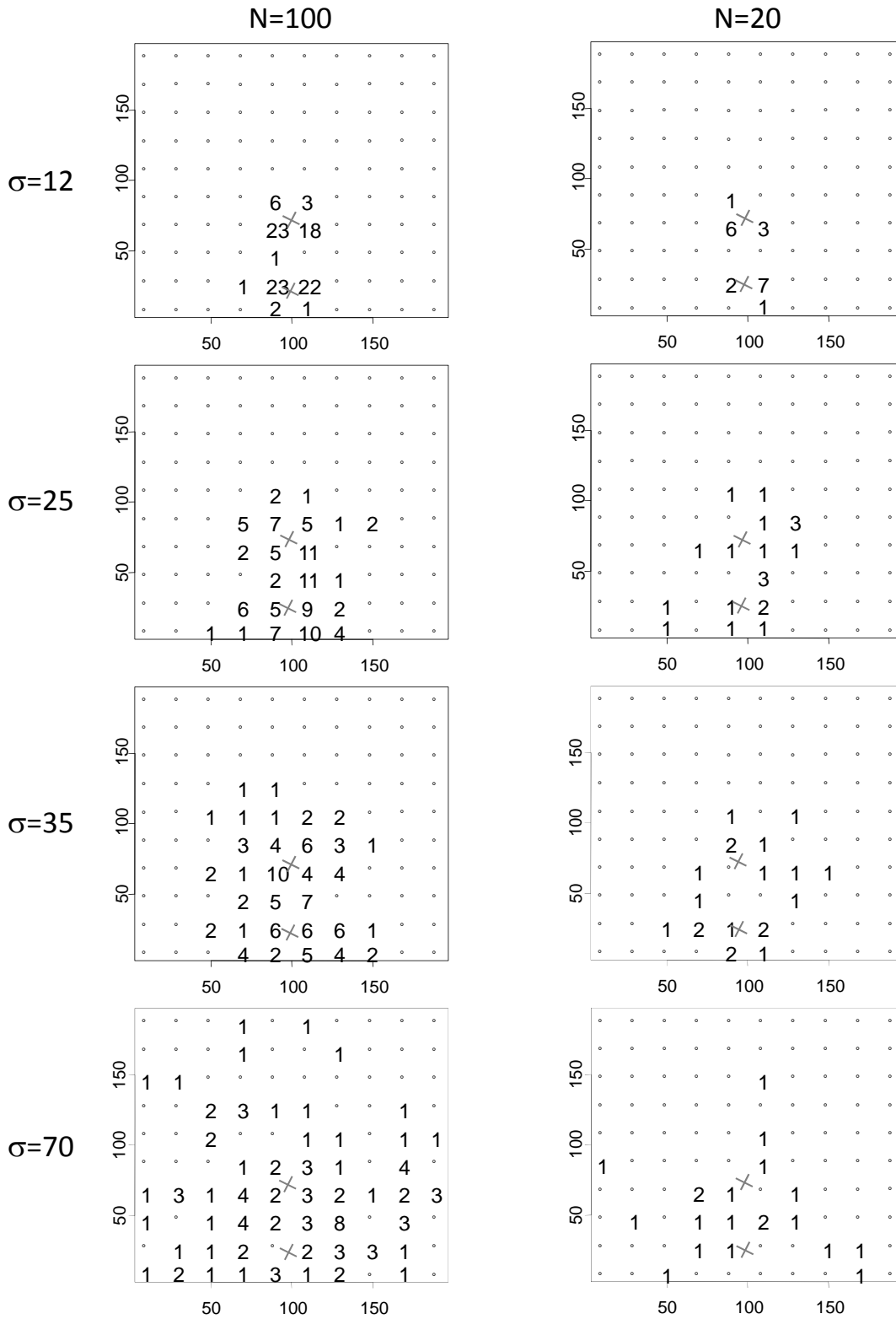
3

Configuration 3: Square  $D$ , far central clusters



1  
2

Configuration 4: Square D, clusters close to edge



Configuration 5: Rectangular  $D$ , central clusters

