



**Repositorio Institucional de la Universidad Autónoma de Madrid**

<https://repositorio.uam.es>

Esta es la **versión de autor** del artículo publicado en:

This is an **author produced version** of a paper published in:

Journal of Multivariate Analysis 146 (2016): 237-247

DOI: <http://dx.doi.org/10.1016/j.jmva.2015.09.017>

**Copyright:** © 2016 Elsevier

El acceso a la versión del editor puede requerir la suscripción del recurso

Access to the published version may require subscription

# Shape classification based on interpoint distance distributions

José R. Berrendero<sup>a</sup>, Antonio Cuevas<sup>a</sup>, Beatriz Pateiro-López<sup>b,\*</sup>

<sup>a</sup>*Departamento de Matemáticas, Universidad Autónoma de Madrid, Spain*

<sup>b</sup>*Departamento de Estadística e Investigación Operativa, Universidad de Santiago de  
Compostela, Spain*

---

## Abstract

According to Kendall (1989), in shape theory... *The idea is to filter out effects resulting from translations, changes of scale and rotations and to declare that shape is “what is left”.* While this statement applies in principle to classical shape theory based on landmarks, the basic idea remains also when other approaches are used. For example, we might consider, for every shape, a suitable associated function which, to a large extent, could be used to characterize the shape. This finally leads to identify the shapes with the elements of a quotient space of sets in such a way that all the sets in the same equivalence class share the same identifying function. In this paper, we explore the use of the interpoint distance distribution (i.e. the distribution of the distance between two independent uniform points) for this purpose. This idea has been previously proposed by other authors [e.g., Osada et al. (2002), Bonetti and Pagano (2005)]. We aim at providing some additional mathematical support for the use of interpoint distances in this context. In particular, we show the Lipschitz continuity of the transformation taking every shape to its corresponding interpoint distance distribution. Also, we obtain a partial identifiability result showing that, under some geometrical restrictions, shapes with different planar area must have different interpoint distance distributions. Finally, we address practical aspects including a real data example on shape classification in marine biology.

*Keywords:* Functional data, Identifiability, Interpoint distance, Shape analysis, Volume function.

---

## 1. Introduction

We are concerned here with the problem of classifying *shapes*, where, in informal terms, a shape is the family of all plane figures that can be obtained from a basic template figure (e.g., a square) by applying isometry transformations (rigid movements + symmetries) together with changes of scale. Also, we would like to include all the “deformed versions” (within some limits) of these basic elements, subject again to isometry transformations and/or scale changes. So, to mention just a very simple example,

---

\*Corresponding address: Rúa Lope Gómez de Marzoa s/n. 15782 Santiago de Compostela. Spain

*Email address:* `beatriz.pateiro@usc.es` (Beatriz Pateiro-López)

18 one could think that we want to automatically discriminate between two  
19 capital letters, say “B” and “D”, manually drawn with a thick line marker,  
20 whatever their size or their orientation.

21 In marine biology, one might be interested on classifying fish species us-  
22 ing shape analysis techniques. In some cases the basis for the recognition  
23 method is the fish image itself; see Storbeck and Daan (2001). Other re-  
24 searches have used the so-called *otholits*, small pieces present in the inner  
25 ear of the fish, which can be considered as “microfossils” whose shapes are  
26 useful in species recognition, among other applications; see Lombarte et al.  
27 (2006). In Section 5 we will use this otolith example as an illustration for  
28 the methodology we propose.

29 Whatever the practical problem at hand, we need to define, in precise  
30 mathematical terms, what we mean for “shapes” in our setting. Then we  
31 will be ready to use the statistical methods for classification, either super-  
32 vised (discrimination) or unsupervised (clustering) from the available data  
33 set of shapes. In the example of Section 5 we will focus on clustering but  
34 discrimination methods could be considered as well.

35 The classical theory of shape analysis is largely based on the use of  
36 “landmarks” (i.e., finite vectors of coordinates characterizing the shapes). It  
37 was developed, to a large extent, by D. Kendall who expressively referred to  
38 shape analysis studies in the following terms: *The idea is to filter out effects*  
39 *resulting from translations, changes of scale and rotations and to declare that*  
40 *shape is “what is left”*; see Kendall (1989). A general perspective of this  
41 theory can be found in Kendall (1989), Kendall et al. (1999) or Kendall and  
42 Le (2010).

43 We should mention however that other, less general, notions of shapes  
44 have been proposed. As Kent (1995) points out, “... *statistical models for*  
45 *shapes may be based on underlying models for the landmarks themselves, or*  
46 *they may be constructed directly within shape space. In some special cases*  
47 *specialized models may be constructed*”. Our approach here could be un-  
48 derstood as one of these specialized models: roughly speaking, we propose  
49 to identify a shape with the corresponding *interpoint distance distribution*,  
50 that is, the distribution of the distance (normalized to 1) between two ran-  
51 domly chosen points in the figure.

52

### 53 *Related literature*

54 In fact, the idea of using the interpoint distance distribution to identify  
55 the shapes has been previously proposed by other authors, with different  
56 applications in mind. For example, the very much cited paper by Osada et  
57 al. (2002) explores the practical aspects of using the interpoint distance in  
58 the problem of discriminating shapes in image analysis. As these authors  
59 point out, “*The primary motivation for this approach is to reduce the shape*  
60 *matching problem to the comparison of probability distributions, which is*  
61 *simpler than traditional shape matching methods that require pose registra-*  
62 *tion, feature correspondence, or model fitting. We find that the dissimi-*

63 *larities between sampled distributions of simple shape functions (e.g., the*  
64 *distance between two random points on a surface) provide a robust method*  
65 *for discriminating between classes of objects (e.g., cars versus airplanes) in*  
66 *a moderately sized database, despite the presence of arbitrary translations,*  
67 *rotations, scales, mirrors, tessellations, simplifications, and model degenera-*  
68 *cies”. See also Bonetti and Pagano (2005) for a different use of interpoint*  
69 *distance distributions in the context of medical research.*

70 In Kent (1994) interpoint distances (between landmarks) are used, via  
71 multi-dimensional scaling, in shape analysis. Our approach here is some-  
72 what different as it avoids the use of landmarks at the expense of some loss  
73 in generality.

74 Let us finally mention that the use of interpoint distance distributions  
75 entails the precise definition of a corresponding, suitable “space of shapes”;  
76 see Section 2 below, where the whole approach makes sense. Other related  
77 shape spaces can be found in the literature, in particular those based on  
78 “deformable templates”: see Grenander (1976), Amit et al. (1991), Hobolt  
79 and Vedel-Jensen (2000), Hobolt et al. (2003).

80

81 *The purpose and contents of this paper*

82 On the theoretical side, we will provide some support for the use of in-  
83 terpoint distance distributions to characterize shapes: first, we relate, in  
84 Theorem 1 below, the distance between interpoint distance distributions  
85 with a natural, geometrically motivated, distance between shapes defined  
86 in Section 2. Second, we consider the problem of providing a sufficient  
87 condition on the sets in the Euclidean plane in order to ensure that two dif-  
88 ferent sets fulfilling this condition must necessarily have different interpoint  
89 distance distributions. Theorem 2 provides a quite general identifiability  
90 criterion, which is in fact the most general result of this type we are aware  
91 of. In the Supplementary Material section we also briefly consider the con-  
92 nection between the interpoint distance distribution and the covariogram  
93 (sometimes called “set covariance”), another popular function which has  
94 been used sometimes to characterize sets and shapes; see Cabo and Badde-  
95 ley (1995, 2003).

96 Finally, in Section 5 our methodology based on interpoint distance distri-  
97 butions is used in a problem of fishes otoliths classification, via hierarchical  
98 clustering.

99

## 100 **2. The space of shapes**

101 In what follows we will mainly focus on the case of shapes in the plane  
102  $\mathbb{R}^2$  (the most important, by far, in practical applications). However, some of  
103 the ideas we will develop can be also adapted to more general, multivariate  
104 cases. Our starting point will be the family  $\mathcal{C}$  of compact non-empty sets in  
105  $\mathbb{R}^2$  with diameter 1; this means that  $\text{diam}(C) = \max\{\|x - y\|, x, y \in C\} =$   
106 1, for all  $C \in \mathcal{C}$ , where  $\|\cdot\|$  stands for the Euclidean norm. We may think

107 that the family  $\mathcal{C}$  is the result of transforming the set of all possible plane  
 108 images by a uniform change of scale (where “uniform” means that the same  
 109 transformation scale is applied in both coordinates) in such a way that all  
 110 of them have a common diameter. We will define our space of shapes as the  
 111 quotient space obtained from a natural equivalence relation in  $\mathcal{C}$ . However,  
 112 the family  $\mathcal{C}$  is too large to work with (in particular, to define a meaningful,  
 113 tractable distance between shapes). So we will need to restrict ourselves to  
 114 a smaller subset  $\mathcal{C}_1 \subset \mathcal{C}$  which, still, will include most “black-and-white”  
 115 images arising in practical applications.

116 To be more specific, given two positive constants  $a$  and  $m_1$ , we define  $\mathcal{C}_1$   
 117 as the class of sets  $C \in \mathcal{C}$  fulfilling the following conditions:

- 118 (i)  $\mu(C) \geq a$ , where  $\mu$  denotes the Lebesgue measure in  $\mathbb{R}^2$ .
- 119 (ii) All the sets in  $\mathcal{C}_1$  are regular, that is, every  $C \in \mathcal{C}_1$  fulfills  $C = \overline{\text{int}(C)}$ .
- 120 (iii)  $\mu(B(\partial C, \epsilon)) < m_1 \epsilon$ ,  $\forall \epsilon \in (0, 1]$ .

121 Here  $\partial A$  denotes the topological boundary of the set  $A$ ,  $B(A, \epsilon)$  stands  
 122 for the “parallel set”  $B(A, \epsilon) = \{x : d(x, A) \leq \epsilon\}$  and  $d(x, A) = \inf\{\|x -$   
 123  $y\|, y \in A\}$  (when  $A = \{x\}$  we will use the standard notation  $B(x, \epsilon)$  instead  
 124 of  $B(\{x\}, \epsilon)$ ).

We assume that the space  $\mathcal{C}_1$  is endowed with the metric,

$$d_{HH}(C, D) = d_H(C, D) + d_H(\partial C, \partial D),$$

125 where  $d_H$  stands for the ordinary Hausdorff metric between compact sets.

126 Let us now define on  $\mathcal{C}_1$  the *isometry* equivalence relation: we will say  
 127 that  $C, D \in \mathcal{C}_1$  are *isometric* (and denote it by  $C \sim D$ ) when there exists a  
 128 isometry (i.e., a map  $i : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  satisfying  $\|i(x) - i(y)\| = \|x - y\|$ ) such  
 129 that  $i(C) = D$ . The family of all sets in  $\mathcal{C}_1$  equivalent to a set  $C$  will be  
 130 represented by  $[C]$ .

131 Finally, denote by  $\mathcal{S}$  the family of equivalence classes and define in  $\mathcal{S}$   
 132 the *quotient metric*,  $\tilde{d}_{HH}$ , using the standard definition method [see, for  
 133 example, Burago et al. (2001, p. 62)],

$$\tilde{d}_{HH}([C], [D]) = \inf\left\{\sum_{i=1}^n d_{HH}(P_i, Q_i) : [P_1] = [C], [Q_n] = [D], n \in \mathbb{N}\right\}, \quad (1)$$

134 where the infimum is taken on all finite sequences such that  $[Q_i] = [P_{i+1}]$  for  
 135  $i = 1, \dots, n - 1$ . In principle, the general method (1) to translate a metric  
 136 to the quotient space defines only a semi-metric, but we will see below that  
 137 in this case it provides a true metric; in fact, we will also see in Proposition  
 138 1 that (1) can be expressed in a much simpler way in our case.

139 The elements of the quotient metric space  $\mathcal{S}$  will be called *shapes*. So  
 140 the shapes are in fact classes of equivalence  $[C]$  for  $C \in \mathcal{C}_1$ .

141

142 *Some motivation*

143 Regarding the intuitive meaning of the assumptions imposed on  $\mathcal{C}_1$ , let  
 144 us note that they do not entail any serious restriction for the practical

145 classification problems of pattern recognition. To explain the meaning of  
 146 these assumptions let us identify our shapes with figures drawn with a sign  
 147 painting marker:

148 Assumption (i) just states that, after re-scaling, our shapes must have  
 149 a minimum “thickness”, expressed in a minimum “drawing area”  $a$ .

150 Condition (ii) is usual in geometric probability models. Under this as-  
 151 sumption, the set  $C$  cannot consist of a closed “central core”  $C_1$  plus some  
 152 “superfluous” parts  $P$  (such as rays or isolated points) with  $\mu(P) = 0$ .

153 Condition (iii) rules out involved drawings, with a very large boundary.  
 154 To see this, let us briefly recall the notion of (*boundary*) *Minkowski content*,  
 155 which is perhaps the simplest way (among several others, see e.g. Mattila  
 156 (1995)) to define the “boundary measure” of a set  $C \subset \mathbb{R}^d$ . Of course, for  
 157 the two-dimensional case, “boundary measure” is synonymous with “length  
 158 perimeter”. In precise terms, the  $(d-1)$ -dimensional (boundary) Minkowski  
 159 content of  $C$  is defined by the limit

$$L_0(C) = \lim_{\epsilon \rightarrow 0} \frac{\mu(B(\partial C, \epsilon))}{2\epsilon}, \quad (2)$$

160 A closely related notion is the *one-sided (outer) Minkowski content*, defined  
 161 by

$$L_0^+(C) = \lim_{\epsilon \rightarrow 0} \frac{\mu(B(C, \epsilon) \setminus C)}{\epsilon}, \quad (3)$$

162 See Ambrosio et al. (2008) for a comprehensive study of this notion, includ-  
 163 ing conditions under which  $L_0(C) = L_0^+(C)$ . For statistical aspects related  
 164 to the Minkowski content we refer to Cuevas et al. (2007) and Berrendero  
 165 et al. (2014). Note that under condition (iii),  $L_0(C) \leq m_1$  for all  $C \in \mathcal{C}_1$ .

166 *A simpler, alternative expression for the distance between shapes.*

167 While (1) gives the “canonical” expression for the distance in a quotient  
 168 metric space, the effective calculation of this metric looks rather trouble-  
 169 some. The following proposition provides a simpler, more natural expression  
 170 for (1) and shows that  $\tilde{d}_{HH}$  is in fact a metric instead of just a semi-metric:  
 171 this means that  $\tilde{d}_{HH}([C], [D]) = 0$  implies  $[C] = [D]$ .

173 **Proposition 1.** *The semi-metric (1) can be expressed as*

$$\tilde{d}_{HH}([C], [D]) = \inf\{d_{HH}(C', D') : C' \in [C], D' \in [D]\}. \quad (4)$$

174 *Moreover, this expression defines in fact a true metric.*

175 *Proof.* This result follows from Th. 2.1 in Cagliari et al. (2014). In part  
 176 (i) of this theorem it is proved that a expression of type (4) holds for the  
 177 semi-distance (1) in the quotient space whenever the equivalence classes of  
 178 this space are the orbits of the action of a group of isometries. This is the  
 179 case here.

180 The fact that expression (1), or (4), defines a true metric is a consequence  
 181 of conclusion (iv) in the aforementioned theorem where the authors prove

182 that (4) is a metric if and only if the orbits of the action are closed sets. To  
183 see that  $[C]$  is a closed set let us consider a convergent sequence  $\{C_n\}$  of  
184 elements  $C_n \in [C]$  with  $n \geq 1$ ; denote by  $C_0$  the limit, i.e.,  $d_{HH}(C_n, C_0) \rightarrow 0$ .  
185 By definition of  $[C]$ , any  $C_n$  can be obtained as  $C_n = t_n(C)$ , where  $t_n$  is an  
186 isometry. Since  $\|t_n(x) - t_n(y)\| = \|x - y\|$ , it turns out that the sequence  
187  $\{t_n\}$  is equicontinuous; moreover, for each  $x \in \mathbb{R}^2$  the sequence  $\{t_n(x)\}$  is  
188 bounded; this is clearly true when  $x \in C$ , since the sequence  $C_n = t_n(C)$  is  
189  $d_H$ -convergent. Then, for a general  $x \in \mathbb{R}^2$ ,  $\{t_n(x)\}$  is also bounded (since,  
190 given  $x_0 \in C$ ,  $\|t_n(x) - t_n(x_0)\| = \|x - x_0\|$ ). So, from Ascoli-Arzelà Theorem  
191 [e.g., Folland (1999, p. 137)] we can ensure that there exists a subsequence  
192 of  $\{t_n\}$ , denoted again  $\{t_n\}$ , such that  $t_n \rightarrow t$ , uniformly on compacts, for  
193 some transformation  $t$ , which must be necessarily an isometry. We thus  
194 have  $d_H(t_n(C), t(C)) \rightarrow 0$ , but, since  $t_n(C) = C_n$  and  $d_H(C_n, C_0) \rightarrow 0$ , we  
195 get  $C_0 = t(C)$ . Finally to see  $C_0 \in [C]$  we only have to prove that  $C_0$  fulfills  
196 conditions (i), (ii) and (iii) stated above in the definition of the class  $\mathcal{C}_1$ . But  
197 this is a trivial consequence of the *Classification Theorem for Isometries on the*  
198 *Plane* [see, for example, Martin (1982, p. 65)] which states that each non-  
199 identity isometry on the plane is either a translation, a rotation, a reflection,  
200 or a glide-reflection (i.e., the composition of a reflection and a translation  
201 in the direction of the reflection axis). This shows that the plane isometries  
202 are “measure preserving” (i.e.,  $\mu(A) = \mu(t(A))$ ) and “boundary preserving”  
203 (i.e.,  $\partial t(C) = t(\partial C)$ ) and therefore, (i)-(iii) hold also for  $t(C) = C_0$ . We  
204 conclude that  $[C]$  is closed.  $\square$

### 205 3. The interpoint distance distribution

206 As mentioned in the introduction, our approach is based on eventually  
207 identifying a shape  $[C]$  with a density function, supported on  $[0, 1]$ . This is  
208 the density function of the distribution of the random variable defined as  
209 the distance between two points randomly chosen on  $C$ .

210 To be more precise, for each  $C \in \mathcal{C}_1$ , define the random variable

$$Y_C = \|X_1 - X_2\|, \quad (5)$$

211 where  $X_1, X_2$  are iid random variables uniformly distributed on  $C$ . It is  
212 readily seen that  $Y_C$  is absolutely continuous with respect to the Lebesgue  
213 measure  $\mu$ . Let us denote by  $f_C$  the density function of  $Y_C$ .

214 Theorem 1 below provides a partial mathematical motivation for the  
215 identification  $[C] \simeq f_C$  by showing that the transformation  $[C] \mapsto f_C$  is  
216 continuous (in fact it is Lipschitz), so that if two shapes are close enough  
217 then the corresponding interpoint distance densities must be also close to-  
218 gether. The problem of analyzing to what extent  $f_C$  is helpful in order to  
219 identify  $C$  will be discussed in Section 4.

The Lipschitz property of the transformation  $C \mapsto f_C$  will be established  
with respect to the standard  $L_1$  metric between densities and also for the  
so-called Wasserstein (or Kantorovich) metric defined, for two cumulative

distribution functions on the real line  $F$  and  $G$ , by

$$d_W(F, G) = \int_{\mathbb{R}} |F(x) - G(x)| dx = \int_0^1 |F^{-1}(t) - G^{-1}(t)| dt,$$

220 where  $F^{-1}, G^{-1}$  denote the corresponding quantile functions. This metric  
 221 has a number of interesting properties and applications. It has been some-  
 222 times called “the earth mover distance”, due to its connections with the  
 223 transportation problem; see Villani, C. (2003). In Rubner et al. (2000) and  
 224 Ling and Okada (2007) can be found some details on the use of this distance  
 225 in image retrieval. Of course, when  $F$  and  $G$  are absolutely continuous (as  
 226 it will always be the case in what follows),  $d_W$  can also be interpreted as a  
 227 distance between the density functions.

228 The following result can be seen as a statement of “compatibility” be-  
 229 tween the distances  $d_1(f, g) = \int_0^1 |f - g| d\mu$  or  $d_W$  (defined in the space of  
 230 densities on  $[0, 1]$ ) and the “natural” distance  $\tilde{d}_{HH}$  defined in our space of  
 231 shapes. The whole point is to replace, in practice, the use of  $\tilde{d}_{HH}$  (whose  
 232 effective calculation is cumbersome) by the more convenient distances  $d_1$  or  
 233  $d_W$ . In principle, the intuitive interpretation of  $d_1(f, g)$  (as the area of the  
 234 region between  $f$  and  $g$ ) is perhaps more direct but, as we have already  
 235 mentioned,  $d_W$  is also used in image analysis, Rubner et al. (2000). Our  
 236 experimental results, see Section 5 and the Supplementary Material doc-  
 237 ument, show a very similar behaviour for both distances with perhaps a  
 238 slightly better performance for  $d_1$ .

239 **Theorem 1.** *Let  $\mathcal{D}$  be the space of probability density functions (with respect*  
 240 *to the Lebesgue measure) on  $[0, 1]$ . Then*

241 (a) *The transformation  $T : \mathcal{C}_1 \rightarrow \mathcal{D}$  given by  $T(C) = f_C$  fulfills the*  
 242 *Lipschitz condition with respect to the  $L_1$  metric, that is,  $d_1(f_C, f_D) \leq$*   
 243  *$md_{HH}(C, D)$ , for some constant  $m > 0$ .*

244 (b) *Also, if we denote by  $F_C$  and  $F_D$  the cumulative distribution functions*  
 245 *of  $Y_C$  and  $Y_D$ , respectively, we have that  $d_W(F_C, F_D) \leq \frac{m}{2}d_{HH}(C, D)$ ,*  
 246 *where  $m$  is the same constant of statement (a).*

247 (c) *The transformation  $T$  induces another transformation  $\tilde{T}([C]) = f_C$ ,*  
 248 *defined in the quotient space, which is also Lipschitz, with constants*  
 249  *$m$  and  $m/2$  respectively, for both considered metrics.*

250 *Proof.* (a) From the relation between the  $L_1$  metric and the total variation  
 251 distance,

$$\int |f_C - f_D| d\mu = 2 \sup_A |P_C(A) - P_D(A)|, \quad (6)$$

252 where  $P_C$  and  $P_D$  are the probability measures associated with  $f_C$  and  $f_D$   
 253 and the supremum is taken on  $\mathcal{B} = \mathcal{B}([0, 1])$ , the Borel sets of  $[0, 1]$  on the



254 elements  $C, D$  chosen to represent  $[C]$  and  $[D]$ . Now, observe that for all  
 255  $A \in \mathcal{B}$ , and using the notation introduced in expression (5),

$$P_C(A) = \mathbb{P}(Y_C \in A) = \mathbb{P}(Y_C \in A | X_1, X_2 \in C \cap D) \mathbb{P}(X_1, X_2 \in C \cap D) \\ + \mathbb{P}(Y_C \in A | X_1 \text{ or } X_2 \notin C \cap D) \mathbb{P}(X_1 \text{ or } X_2 \notin C \cap D),$$

256 where  $X_1, X_2$  are iid uniformly distributed on  $C$ . A similar expression holds  
 257 for  $P_D(A)$ , except that  $C$  is replaced with  $D$  and  $X_1, X_2$  are replaced with  
 258  $X_1^*, X_2^*$ , iid uniform on  $D$ , that is,

$$P_D(A) = \mathbb{P}(Y_D \in A) = \mathbb{P}(Y_D \in A | X_1^*, X_2^* \in C \cap D) \mathbb{P}(X_1^*, X_2^* \in C \cap D) \\ + \mathbb{P}(Y_D \in A | X_1^* \text{ or } X_2^* \notin C \cap D) \mathbb{P}(X_1^* \text{ or } X_2^* \notin C \cap D),$$

259 Note that  $\mathbb{P}(Y_C \in A | X_1, X_2 \in C \cap D) = \mathbb{P}(Y_D \in A | X_1^*, X_2^* \in C \cap D)$ .  
 260 Therefore,

$$|P_C(A) - P_D(A)| \leq \mathbb{P}(Y_C \in A | X_1, X_2 \in C \cap D) \mathbb{P}(X_1 \text{ or } X_2 \notin C \cap D) \\ + \mathbb{P}(Y_C \in A | X_1, X_2 \in C \cap D) \mathbb{P}(X_1^* \text{ or } X_2^* \notin C \cap D) \\ + \mathbb{P}(Y_C \in A | X_1 \text{ or } X_2 \notin C \cap D) \mathbb{P}(X_1 \text{ or } X_2 \notin C \cap D) \\ + \mathbb{P}(Y_D \in A | X_1^* \text{ or } X_2^* \notin C \cap D) \mathbb{P}(X_1^* \text{ or } X_2^* \notin C \cap D).$$

261 For the first term in the right-hand side of  $|P_C(A) - P_D(A)|$  we have,

$$\mathbb{P}(Y_C \in A | X_1, X_2 \in C \cap D) \mathbb{P}(X_1 \text{ or } X_2 \notin C \cap D) \\ \leq \mathbb{P}(X_1 \text{ or } X_2 \in C \setminus D) \leq 2\mathbb{P}(X_1 \in C \setminus D) \leq \frac{2}{a} \mu(C \setminus D),$$

262 where  $a$  is the minimal area of the elements of  $\mathcal{C}$  defined in condition (i).  
 263 The same holds for the third term. Similarly, we have that the second and  
 264 fourth terms in  $|P_C(A) - P_D(A)|$  are smaller than  $\frac{2}{a} \mu(D \setminus C)$ . Hence,

$$\sup_A |P_C(A) - P_D(A)| \leq \frac{4}{a} \mu(C \Delta D), \quad (7)$$

265 where  $C \Delta D$  stands for the symmetric difference  $C \Delta D = (C \setminus D) \cup (D \setminus C)$ .

266 Let us now prove that

$$\mu(C \Delta D) \leq 2m_1 d_{HH}(C, D), \quad (8)$$

267 where  $m_1$  is the constant introduced in the definition on  $\mathcal{C}_1$ . To see this,  
 268 put  $d_{HH}(C, D) = r$  and take  $x \in C \setminus D$ . We must have  $x \in B(D, r) \setminus D$   
 269 which entails  $x \in B(\partial D, r) \subset B(\partial C, 2r)$ . Similarly, if  $x \in D \setminus C$  we have  
 270  $x \in B(C, r) \setminus C$  so that  $x \in B(\partial C, r)$ .

Thus, using assumption (iii) we have obtained that

$$\mu(C \Delta D) \leq \mu(B(\partial C, 2r)) \leq 2m_1 r = 2m_1 d_{HH}(C, D).$$

271 This, together with (6), (7) and (8) proves the first statement (a).

272

273 (b) This directly follows from Theorem 4 in Gibbs and Su (2002). Ac-  
274 cording to this result, if we consider probability measures defined on a space  
275  $\Omega$  with finite diameter,  $\text{diam}(\Omega)$ , we have  $d_W \leq \text{diam}(\Omega) \cdot d_{TV}$ . In our case,  
276 all the considered distributions are defined on the unit interval. This, to-  
277 gether with  $2d_{TV} = d_1$  leads to statement (b).

278

279 (c) This statement follows from parts (a) and (b) combined with the  
280 expression (4) of the quotient metric.  $\square$

281 **Remark 1.** *The search for a Lipschitz-type as that in Theorem 1 is quite*  
282 *natural in those situations where a set (or a shape) is replaced with a more*  
283 *convenient auxiliary function. For example, a result in a similar spirit can*  
284 *be found in Cabo and Baddeley (1995, Th. 5.4) but these authors consider*  
285 *the so-called covariogram function, instead of the interpoint distance density,*  
286 *and the distance  $d_{HH}$  is replaced with another metric defined in terms of*  
287 *the so-called “linear scan transform”.*

288 The covariogram of a bounded Borel set  $A \subset \mathbb{R}^d$  is defined by  $K_A(y) =$   
289  $\mu(A \cap T_y A)$ , where  $y \in \mathbb{R}^d$ ,  $T_y A = A - y = \{a - y : a \in A\}$  and  $\mu$  is  
290 the Lebesgue measure in  $\mathbb{R}^d$ . This function is useful in different problems of  
291 stochastic geometry and stereology. Some references are Cabo and Baddeley  
292 (1995, 2003) and Galerne (2011). Using some results in these papers it  
293 is easy to prove (see the Supplementary Material document for details)  
294 that the random interpoint distance  $Y_C$  of a bounded Borel set  $C$  in the  
295 plane has a continuous density  $f_C$  with  $f_C(0) = 0$  and  $f_C(\rho_C) = 0$ , where  
296  $\rho_C = \text{diam}(C)$ .

#### 297 4. The identifiability problem

298 In order to implement the idea of identifying a shape  $[C]$  with the cor-  
299 responding interpoint distance density  $f_C$ , we must still overcome a further  
300 obstacle. Even if we restrict to the space of shapes  $[C]$  with  $C \in \mathcal{C}_1$  (where  
301 the continuity of the transformation  $[C] \mapsto f_C$  is warranted) one might have  
302 that  $f_C = f_D$  for  $[C] \neq [D]$ . This follows as a consequence of a counterex-  
303 ample, due to Mallows and Clark (1970) [inspired by a question posed by  
304 Blaschke], showing two non-congruent polygons,  $C$  and  $D$  with the same  
305 *chord length* distribution. The chord length is the length of the segment  
306 intercepted in  $C$  by a random chord. Since the chord length distribution  
307 determines uniquely the interpoint distance distribution [see, Matern (1986,  
308 p. 25)] the mentioned counterexample applies also to the interpoint distance  
309 distribution.

310 The interpoint distance has been also used (with applications to crystal-  
311 lography and DNA mapping) in finite sets of points; see Caelli (1980) and  
312 Lemke et al. (2003) for further counterexamples, references and insights.

313 Thus, in summary, the interpoint distance distribution has not full ca-  
 314 pacity to discriminate shapes. Hence, we should further restrict our shape  
 315 space to those sets  $[C]$  such that  $C$  lives in an appropriate subset  $\mathcal{C}_2 \subset \mathcal{C}_1$   
 316 fulfilling the identifiability condition

$$(iv) \text{ For all } C, D \in \mathcal{C}_2 \text{ with } [C] \neq [D] \text{ we have } Y_C \stackrel{d}{\neq} Y_D, \quad (9)$$

317 where  $Y_C$  and  $Y_D$  denote the interpoint distances (5) on  $C$  and  $D$  and the  
 318 notation  $\stackrel{d}{\neq}$  means that both variables are not identically distributed.

319 Some identifiability problems similar to (9) have been considered in  
 320 the stochastic geometry literature under different conditions. For example,  
 321 Matheron (1986) formulated the following conjecture: *Every planar convex*  
 322 *body is determined within all planar convex bodies by its covariogram, up to*  
 323 *translations and reflections.* This conjecture was completely solved, in the  
 324 affirmative by Averkov and Bianchi (2009).

325 In the following subsection we will show that the analogous problem (9)  
 326 for the interpoint distance distribution can be solved under quite general  
 327 conditions, which do not require convexity.

#### 328 4.1. Interpoint distances and polynomial area

329 The main geometric assumption we will use to guarantee identifiability  
 330 is defined as follows.

331 **Definition 1.** *A set  $C \subset \mathbb{R}^2$  is said to have inner polynomial area* if there  
 332 exist constant  $R = R(C) > 0$  and  $L = L(C) > 0$  such that

$$\mu(I_r(C)) = \mu(C) - L(C)r + \pi r^2, \text{ for } 0 \leq r < R, \quad (10)$$

333 where  $I_r(C)$  denotes the *inner parallel set*  $I_r(C) = \{x \in C : B(x, r) \subset C\}$ .

334 For example, the circle  $C = B(0, m)$  fulfills (10) with  $L(C) = 2\pi m$ ,  
 335  $R < m$  and  $\mu(C) = \pi m^2$ .

336 **Remark 2.** *It is clear that, if (10) holds, the quantity  $L(C)$  could be ob-*  
 337 *tained as a sort of inner Minkowski content,  $L_0^-(C)$  defined in a similar way*  
 338 *to outer version  $L_0^+(C)$  given in (3). Moreover, if the ordinary (two-sided)*  
 339 *Minkowski content,  $L_0(C)$  does exist [see (2)] then condition (10) clearly*  
 340 *entails  $L(C) = L_0(C) = L_0^+(C)$ .*

341 Now, our goal is to motivate this definition in a twofold way. First,  
 342 we will relate it to some relevant mathematical concepts. Second, we will  
 343 exhibit a broad class of sets satisfying (10). For this purpose, it will be  
 344 useful to recall some notions, due to Federer (1959), from geometric mea-  
 345 sure theory: the *reach* of a closed set is defined as the supremum,  $\text{reach}(C)$ ,  
 346 of those values such that any point  $x$  whose distance to  $C$  is smaller than  
 347  $\text{reach}(C)$  has only one closest point on  $C$ . This concept leads to a valuable  
 348 generalization of the notion of convex set, which can be interpreted also as

349 a geometric smoothness condition (not directly relying on differentiability  
 350 assumptions). Figure 1 illustrates the nice intuitive meaning of this notion.  
 351 It can be shown that  $C$  is convex if and only if  $\text{reach}(C) = \infty$ . Accord-  
 352 ing to a result proved by Federer (1959) [which is a generalization of the  
 353 classical Steiner’s formula for convex sets], the sets of positive reach have a  
 354 polynomial volume. More precisely [Federer (1959), Ths. 5.6 and 5.19]:

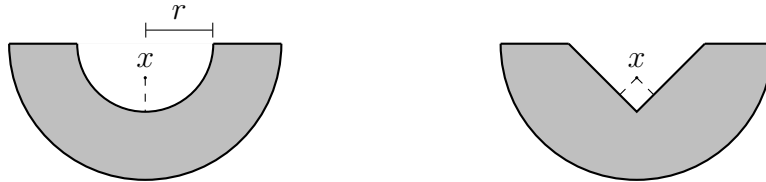


Figure 1: The set  $C$  in the left has positive reach  $r$  (any  $x$  whose distance to  $C$  is smaller than  $r$  has only one closest point on  $C$ ). The set  $C$  in the right has not positive reach.

355

356 *If  $S \subset \mathbb{R}^d$  is a compact set with  $r_0 = \text{reach}(S) > 0$ , then there exist*  
 357 *unique values  $\Phi_0(S), \dots, \Phi_d(S)$  over such that*

$$\mu(B(S, r)) = \sum_{i=0}^d r^{d-i} \omega_{d-i} \Phi_i(S), \text{ for } 0 \leq r < r_0, \quad (11)$$

358 *where  $\omega_j$  is the  $j$ -dimensional measure of a unit ball in  $\mathbb{R}^j$ .*

359 **Remark 3.** *The above result has some connections with other important*  
 360 *geometric notions. Some are almost immediate: for example, if  $S$  is a com-*  
 361 *compact set with positive reach, then  $\Phi_d(S) = \mu(S)$  and the outer Minkowski*  
 362 *content defined in (2) always exists and corresponds to the first-degree term*  
 363 *in (11). Another, not so obvious, deep geometric connection of (11) is as*  
 364 *follows: the coefficient  $\Phi_0(S)$  coincides with the Euler characteristic of  $S$ .*  
 365 *This is an integer-valued topological invariant with deep geometric implica-*  
 366 *tions, far beyond the scope of this paper; see, e.g., Hatcher (2002) for details.*  
 367 *In the following remark we show an example which, in addition to recall the*  
 368 *intuitive meaning of  $\Phi_0(S)$ , will also serve for further generalizations.*

369 *On the other hand, note that  $\text{reach}(S) = r_0 > 0$  is just a sufficient*  
 370 *condition for polynomial volume in the interval  $[0, r_0)$ . Many other sets,*  
 371 *which do not satisfy  $\text{reach}(S) > 0$  (such as that of the right panel in Figure*  
 372 *1), might fulfill a polynomial volume property of type (11).*

373 **Remark 4.** *Let us consider the annulus  $D = B(0, M) \setminus \text{int}(B(0, m))$ , with*  
 374  *$m < M$ . A direct calculation shows that  $\mu(B(D, r)) = 2\pi(M+m)r + \pi(M^2 -$*   
 375  *$m^2)$ . Moreover, it is clear that  $\text{reach}(D) = m$ . As a conclusion, the annulus*  
 376  *$D$  fulfills  $\Phi_0(D) = 0$  in (11). By the way, the same holds for any set, of*  
 377 *positive reach, homeomorphic to the annulus (as the Euler characteristic is*  
 378 *a topological invariant).*

379 Now, we are ready to show that in fact (10) applies to a broad class  
 380 of sets under a quite general condition (expressed in terms of the classical  
 381 positive reach property).

382 **Proposition 2.** *The class  $\mathcal{P}(R)$  of sets which fulfill condition (10) contains*  
 383 *all regular sets  $C$  such that for some closed ball  $B_1$ , with  $C \subset \text{int}(B_1)$ , the*  
 384 *set  $E = B_1 \setminus \text{int}(C)$  has positive reach  $R$  and it is homeomorphic to an*  
 385 *annulus (as that considered in Remark 4).*

*Proof.* Note that  $\mu(B(E, r)) = \mu(E) + \mu(B(B_1, r)) - \mu(B_1) + \mu(C) - \mu(I_r(C))$ .  
 Now,  $E$  has positive reach  $R$  and, by (11),  $\mu(B(E, r)) = rL_0^+(E) + \mu(E)$ .  
 Note also that  $\Phi_0(E) = 0$  since  $B_1 \setminus \text{int}(C)$  is homeomorphic to an annulus  
 $D$  (for which  $\Phi_0(D) = 0$ , according to Remark 4). Therefore,

$$\mu(I_r(C)) = \mu(C) - L(C)r + \pi r^2, \text{ with } L(C) = L_0^+(E) - L_0(B_1).$$

386 □

387 As a conclusion, we have that the class of sets fulfilling (10) includes  
 388 many relevant sets found in practice. See Berrendero et al. (2014) for further  
 389 information and statistical applications of the notion of polynomial volume.

390 We are now ready to establish the main result of this section which  
 391 provides a large class  $\mathcal{R}$  of sets which can be distinguished from each other  
 392 according to the distribution of the respective interpoint distances. In other  
 393 words, if  $C, D \in \mathcal{R}$  then  $f_C \neq f_D$ , where  $f_C$  denotes the density function of  
 394 the interpoint distance  $Y_C$ .

395 **Theorem 2.** (a) *Suppose that  $C$  is a compact set in  $\mathbb{R}^2$  fulfilling condition*  
 396 *(10) of inner polynomial area. Denote by  $Y_C$  the interpoint distance in  $C$ .*  
 397 *Then*

$$\mathbb{P}(Y_C \leq \rho) = \frac{\pi\rho^2}{\mu(C)} - \frac{\pi\rho^3 L(C)}{\mu(C)^2} + \frac{\pi^2\rho^4}{\mu(C)^2} + \frac{1}{\mu(C)^2} \int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx, \quad (12)$$

398 for  $\rho > 0$  be small enough so that  $\rho < R$  in (10) and  $I_\rho(C) \neq \emptyset$ , where  
 399  $I_\rho(C)$  denotes the inner parallel set  $I_\rho(C) = \{x \in C : B(x, \rho) \subset C\}$ .

400 (b) Let  $C, D$  be compact sets, with diameter 1, in  $\mathbb{R}^2$  fulfilling the poly-  
 401 nomial inner area condition (10). If  $\mu(C) \neq \mu(D)$ , then the respective  
 402 interpoint distance have different distributions, that is,  $Y_C \stackrel{d}{\neq} Y_D$ .

403 *Proof.* (a) Let  $X_1, X_2$  be iid random variables uniformly distributed on  $C$ .  
 404 Denote by  $P_C$  the probability distribution uniform on  $C$ .

$$\begin{aligned} \mathbb{P}(Y_C \leq \rho) &= \int_C \mathbb{P}(X_1 \in B(x, \rho)) dP_C(x) = \int_C P_C(B(x, \rho)) dP_C(x) \\ &= \int_{I_\rho(C)} P_C(B(x, \rho)) dP_C(x) + \int_{C \setminus I_\rho(C)} P_C(B(x, \rho)) dP_C(x) \\ &= \frac{1}{\mu(C)^2} \int_{I_\rho(C)} \mu(B(x, \rho)) dx + \frac{1}{\mu(C)^2} \int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx \end{aligned}$$

$$\begin{aligned}
&= \pi\rho^2 \frac{\mu(I_\rho(C))}{\mu(C)^2} + \frac{1}{\mu(C)^2} \int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx \\
&= \pi\rho^2 \frac{\mu(C) - L(C)\rho + \pi\rho^2}{\mu(C)^2} + \frac{1}{\mu(C)^2} \int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx \\
&= \frac{\pi\rho^2}{\mu(C)} - \frac{\pi\rho^3 L(C)}{\mu(C)^2} + \frac{\pi^2\rho^4}{\mu(C)^2} + \frac{1}{\mu(C)^2} \int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx
\end{aligned}$$

405 (b) This result readily follows from (a). First note that the integral  
406  $\int_{C \setminus I_\rho(C)} \mu(B(x, \rho) \cap C) dx$  in the last term of (12) is of order  $\rho^3$  as  $\rho \rightarrow 0$   
407 since the integrand is of type  $O(\rho^2)$  and the measure of the integration set  
408 is  $O(\rho)$ , from the polynomial area assumption. Therefore the main term in  
409 (12) is  $\frac{\pi\rho^2}{\mu(C)}$ . Now, If  $\mu(C) \neq \mu(D)$ , the main terms  $\frac{\pi\rho^2}{\mu(C)}$  in the respective  
410 expressions (12) for the distribution functions of  $Y_C$  and  $Y_D$  are different.  
411 Hence, these distribution functions must be different for  $\rho$  small enough.  $\square$

## 412 5. An application to fish family identification from otolith images

413 The AFORO database (<http://www.icm.csic.es/aforo/>) offers an  
414 open online catalogue of fish otolith images. As defined by Tuset et al.  
415 (2008), otoliths are “acellular concretions of calcium carbonate and other  
416 inorganic salts that develop over a protein matrix in the inner ear of ver-  
417 tebrates”. The application of otoliths research has developed significantly  
418 over the last years, see Begg et al. (2005). Fish species identification, age  
419 and growth determination or stock and hatchery management are some of  
420 the most common and important applications of otolith data.

421 The AFORO database contains at present more than 4500 high res-  
422 olution images corresponding to 1382 species and 216 families from the  
423 Mediterranean Sea and the Antarctic, Atlantic, Indic and Pacific Oceans.  
424 For this study, we have considered fishes belonging to three families: *Solei-*  
425 *dae*, *Labridae* and *Scombridae*. There are important features of otoliths  
426 that can be used for species identification. The otolith shape (outline), the  
427 inner groove and the otolith margins, among others, are important char-  
428 acteristics in the morphological description of otoliths. According to the  
429 characterization in Tuset et al. (2008), the terms that better describe the  
430 shape of the otolith’s outline in the family *Soleidae* are discoidal, elliptic  
431 and bullet-shaped (and intermediate shapes between these three). For the  
432 family *Labridae*, the otolith’s outlines are mainly cuneiform, oval and rect-  
433 angular (and intermediate shapes). For the family *Scombridae*, the otoliths  
434 are characterized by their serrate margins. See Figure 2 for examples of  
435 otoliths from these three families.

436  
437 *Interpoint distance: estimated distribution and density functions.* We have  
438 240 high resolution images of otoliths and their corresponding contours (70  
439 *Soleidae*, 125 *Labridae* and 45 *Scombridae*). For the practical implementa-  
440 tion of the method in this example, we need to generate pairs of uniform

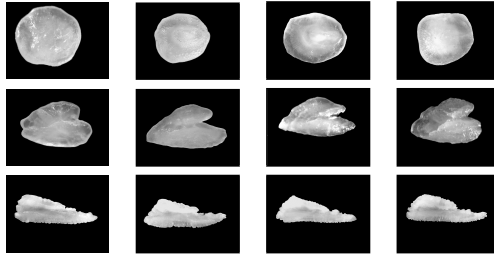


Figure 2: High resolution images of otoliths. First row: *Soleidae*. Second row: *Labridae*. Third row: *Scombridae*.

441 points within the otoliths (area in black in the filled-in contour images,  
 442 see Supplementary material). For this purpose, we can use the standard  
 443 acceptance-rejection method, generating uniform points on a rectangle con-  
 444 taining the otolith and accepting those points belonging to the black area.  
 445 This procedure will be slow on images with a small percentage of black pix-  
 446 els with respect to the bounding rectangle. Another possibility, faster than  
 447 the acceptance-rejection method, is to select pixels in black randomly and,  
 448 for each pixel, generate a uniformly distributed random point within that  
 449 pixel. Other issues about sampling generation in more general situations,  
 450 such as 3D shapes, are discussed in Osada et al. (2002). For each otolith,  
 451 we compute the empirical cumulative distribution function of the interpoint  
 452 distance using the distances (rescaled by the estimated diameter) between  
 453 50000 pairs of random points on the otolith. Figure 3 shows the empirical  
 454 cumulative distribution functions (left) and the estimated interpoint dis-  
 455 tance densities (right) corresponding to the 240 otoliths (*Soleidae*, *Labridae*  
 456 and *Scombridae* in dark, medium and light gray, respectively).

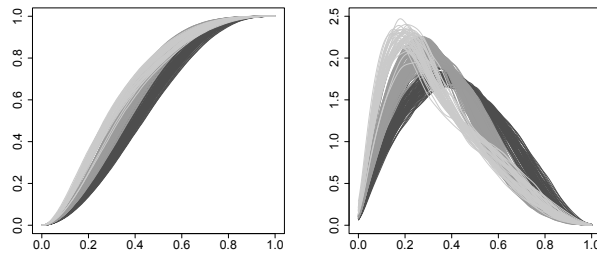


Figure 3: Left, empirical distribution functions of the interpoint distance on the otoliths. Right, estimated densities. In dark gray, *Soleidae*. In medium gray, *Labridae*. In light gray, *Scombridae*.

457  
 458 *Hierarchical clustering.* First, we apply an agglomerative hierarchical clus-  
 459 tering procedure for each pair of families, considering both the  $L_1$  distance  
 460 between densities and the Wasserstein distance between cumulative dis-  
 461 tribution functions as the dissimilarity criterion. As linkage method, we  
 462 have considered single-linkage, complete-linkage and average-linkage. For

463 the sake of brevity, we only discuss here the average-linkage method, which  
 464 gives the best results.

465 Let us first discuss the results on the dataset consisting of *Soleidae* and  
 466 *Labridae* otoliths (dataset A). Figure 4 shows the dendrogram based on the  
 467  $L_1$  distance between the estimated densities. We can consider the otoliths  
 468 divided in two big groups (represented in dark and light gray). We ob-  
 469 serve, see Table 1 (left), that one cluster is dominated by *Soleidae* otoliths  
 470 (94.29% of *Soleidae* otoliths belong to cluster 1) and the other contains  
 471 mainly *Labridae* otoliths (98.40% of *Labridae* otoliths belong to cluster 2).  
 472 The results of the clustering procedure based on the Wasserstein distance  
 473 between distribution functions are quite similar, see Table 1 (right).

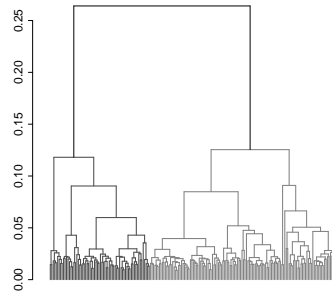


Figure 4: Dendrogram using the  $L_1$  distance between interpoint distance densities for the dataset consisting of *Soleidae* and *Labridae* otoliths (dataset A). The tree is cut into two groups, represented in dark and light gray.

Table 1: Hierarchical clustering on three datasets of otoliths. For each dissimilarity criterion, count and row percent of the true family labels versus the group labels for a partition into two clusters.

		$L_1$ distance		Wasserstein distance	
		Cluster 1	Cluster 2	Cluster 1	Cluster 2
Dataset A	<i>Soleidae</i>	66	4	67	3
		94.29%	5.71%	95.71%	4.29%
	<i>Labridae</i>	2	123	2	123
		1.60%	98.40%	1.60%	98.40%
Dataset B	<i>Soleidae</i>	69	1	69	1
		98.57%	1.43%	98.57%	1.43%
	<i>Scombridae</i>	0	45	0	45
		0.00%	100.00%	0.00%	100.00%
Dataset C	<i>Labridae</i>	123	2	123	2
		98.40%	1.60%	98.40%	1.60%
	<i>Scombridae</i>	2	43	2	43
		4.44%	95.56%	4.44%	95.56%

474 Now, let us consider the dataset consisting of *Soleidae* and *Scombridae*



475 otoliths (dataset B). We apply again an agglomerative hierarchical cluster-  
476 ing procedure using both the  $L_1$  distance and the Wasserstein distance as  
477 the dissimilarity criterion. We split the corresponding dendrograms into  
478 two groups. The results are summarized in Table 1 (dataset B). We found  
479 that all but one of the *Soleidae* otoliths belong to the first cluster and all  
480 the *Scombridae* otoliths belong to the other cluster.

481 Finally, we consider the complete dataset consisting of otoliths from the  
482 three families and apply the agglomerative hierarchical clustering procedure  
483 using the  $L_1$  distance. If we cut the corresponding tree into three groups, we  
484 obtain that 94.29% of *Soleidae* otoliths belong to the first cluster, 96.80%  
485 of *Labridae* otoliths belong to the second cluster and 95.56% of *Scombridae*  
486 otoliths belong to the third cluster. The dendrogram and the complete  
487 table of results based on the  $L_1$  distance and the Wasserstein distance can  
488 be found in the Supplementary Material.

489  
490 *k-means clustering.* Now, we investigate the performance of the  $k$ -means  
491 clustering algorithm. We apply the  $k$ -means algorithm to each pair of fam-  
492 ilies of otoliths ( $k = 2$ ). Here we briefly describe the results based on the  
493  $L_1$  distance (the complete table of results based on the  $L_1$  distance and  
494 the Wasserstein distance is provided as Supplementary Material). For the  
495 dataset consisting of *Soleidae* and *Labridae* images, we obtain a 96.92% of  
496 correctly clustered otoliths. For the dataset consisting of *Soleidae* and *Scom-*  
497 *bridae* images, we obtain a 99.13% of correctly clustered otoliths. For the  
498 dataset consisting of *Labridae* and *Scombridae* images, we obtain a 97.64%  
499 of correctly clustered otoliths.

500  
501 *Final remarks.* (a) We observe that both clustering methods (hierachical  
502 clustering and  $k$ -means) perform reasonably well.

503 We would also like to note that the main reason to choose the families  
504 *Soleidae*, *Labridae* and *Scombridae* was that the AFORO database contains  
505 a large number of images of each of these families. At the beginning of the  
506 study, we had also considered two other large families: *Gobiidae* and *Ser-*  
507 *ranidae* (see the Supplementary Material for examples of otoliths in these  
508 two families). As might be expected, the clustering methods did not per-  
509 form well, for example, for the dataset containing *Gobiidae* and *Soleidae*  
510 otoliths since the shape of some of the *Gobiidae* otoliths resembles that of  
511 the *Soleidae* otoliths. The same occurs for the dataset containing *Serranidae*  
512 and *Labridae* otoliths.

513 (b) As a referee pointed out to us, the use of interpoint distance distri-  
514 butions can be extended to more general (not necessarily planar) situations.  
515 Thus, otholits are in fact three-dimensional structures, one might consider  
516 also the 3D extension of our technique. Likewise, one might think of in-  
517 corporating possibly non-uniform choices of the random points defining the  
518 interpoint distances. This would entail additional theoretical and computa-  
519 tional challenges; see Tebaldi et al. (2011) for computational aspects related

520 to interpoint distance distributions.

## 521 **6. Discussion. Connections with FDA**

522 The study of those problems where the “sample elements” and/or the  
523 target “parameters” are members of an infinite-dimensional space is today  
524 a mainstream topic in statistical research. Of course, the classical nonpara-  
525 metric curve estimation theory (developed since the 1960’s) is an impor-  
526 tant precedent but perhaps the excellent book by Grenander (1981) is one  
527 of the pioneering references in putting together these ideas in a more or  
528 less systematic fashion. As it often happens in the beginnings of a new  
529 scientific theory, the terminologies are not unified. Grenander’s proposal  
530 *abstract inference*, has been later be replaced by the non-exactly equiva-  
531 lent, *infinite-dimensional statistics* (Bongiorno et al. (2014)) or *functional*  
532 *statistics*. Recently, the overview paper Marron and Alonso (2014) pro-  
533 poses the name Object Oriented Data Analysis (OODA) to refer to “*sta-*  
534 *tistical analysis of populations of complex objects*”; In that paper, classical  
535 Kendall’s Shape Analysis (SA) is explicitly included in the OODA frame-  
536 work, alongside *Functional Data Analysis (FDA)*, the study of statistical  
537 methods (regression, classification, principal components, etc.) suitable for  
538 those situations in which the sample data  $x_1, \dots, x_n$  are functions, typically  
539 (but not necessarily) depending of one real variable,  $x_i : [a, b] \rightarrow \mathbb{R}$ .

540 If we take the number of publications as a hint of the popularity of a  
541 scientific topic, FDA is perhaps the most successful chapter in the field of  
542 infinite-dimensional statistics. Since the popular textbook by Ramsay and  
543 Silverman (1997), several other well-known monographs have contributed  
544 to the popularization of FDA; see Ferraty and Vieu (2006), Ferraty and  
545 Romain (2011) and Horváth and Kokoszka (2012), among others. See also,  
546 Cuevas (2014) for a recent overview.

547 We think that Marron and Alonso (2014) make a good point in bringing  
548 together shape analysis and FDA as two particular instances of OODA. In  
549 fact, the conceptual relation between both topics is quite obvious at a formal  
550 level, since shapes can be ultimately identified with functions of some kind  
551 (or equivalence classes of functions). However, the connection holds true  
552 from, at least, two other more relevant aspects:

553 (a) We have shown that (under some restrictions) shapes can be identi-  
554 fied with *density functions* (those of the corresponding interpoint distance  
555 distributions). Hence, following our approach, a statistical problem with  
556 shapes can be recast as a FDA problem in which the available data are  
557 density functions. See Delicado (2011) for an account of this topic. Many  
558 interesting issues can be considered in such a setup: for example, principal  
559 components analysis and other techniques of dimension reduction.

560 (b) Still, considering SA from the FDA point of view suggest to study  
561 the adaptation of the increasing literature on FDA *variable selection* (or  
562 *feature selection*), to the SA framework; see, for example Berrendero et al.

563 (2015) and references therein for some recent theoretical and practical in-  
564 sights on this subject. In particular, it seems worthwhile to analyze the  
565 possible connections between some of these variable selection and the clas-  
566 sical landmarks theory in shape analysis.

## 567 Acknowledgement

568 This work has been partially supported by Spanish Grants MTM2013-  
569 44045-P (Berrendero and Cuevas) and MTM2013-41383-P (Pateiro-López).

## 570 Supplementary material

571 The “Supplementary material” document provides additional figures and  
572 tables for Section 5. It includes also a short discussion on the relation  
573 between the covariogram function and the interpoint distances distribution.

## 574 References

- 575 Ambrosio, L., Colesanti, A. and Villa, E. (2008). Outer Minkowski content for  
576 some classes of closed sets. *Math. Ann.* 342, 727–748.
- 577 Amit, Y., Grenander, U. and Piccioni, M. (1991). Structural image restoration  
578 trough deformable templates. *J. Amer. Statist. Assoc.* 86, 376–387.
- 579 Averkov, G. Bianchi, G. (2009). Confirmation of Matheron’s conjecture on the  
580 covariogram of a planar convex body. *J. Eur. Math. Soc.* 11, 1187-1202.
- 581 Begg, G. A., Campana, S. E., Fowler, A. J., and Suthers, I. M. (2005) Otolith  
582 research and application: current directions in innovation and implementation.  
583 *Marine and Freshwater Research*, 56, 477-483.
- 584 Berrendero, J.R., Cholaquidis, A., Cuevas, A. and Fraiman, R. (2014). A geo-  
585 metrically motivated parametric model in manifold estimation. *Statistics* 48,  
586 983–1004.
- 587 Berrendero, J.R., Cuevas, A. and Torrecilla, J.L. (2015). Variable selection in  
588 functional data classification: a maxima-hunting proposal. *Statistica Sinica*, to  
589 appear.
- 590 Bonetti, M. and Pagano, M. (2005). The interpoint distance distribution as a  
591 descriptor of point patterns, with an application to cluster detection. *Statistics*  
592 *in Medicine* 24, 753-773.
- 593 Bongiorno, E.G., Goia, A., Salinelli, E. and Vieu, P. (2014). *Contributions in*  
594 *infinite-dimensional statistics and related topics.*. S.E. Esculapio, Bologna.
- 595 Burago, D., Burago, Y. and Ivanov, S. (2002). *A Course in Metric Geometry.*  
596 American Mathematical Society.
- 597 Cabo, A.J. and Baddeley, A.J. (1995). Line transects, covariance functions and  
598 set convergence. *Adv. Appl. Prob.* 27, 585-605.
- 599 Cabo, A.J. and Baddeley, A.J. (2003). Estimation of mean particle volume using  
600 the set covariance function. *Adv. Appl. Prob.* 35, 27-46.

- 601 Caelli, T. (1980). On generating spatial configurations with identical interpoint  
602 distance distributions. *Combinatorial mathematics, VII* (Proc. Seventh Aus-  
603 tralian Conf., Univ. Newcastle, Newcastle, 1979), pp. 69-75, Lecture Notes in  
604 Math., 829, Springer, Berlin.
- 605 Cagliari, F., Di Fabio, B. and Landi, C. (2014). The natural pseudo-distance as  
606 a quotient pseudo-metric, and applications. *Forum Mathematicum*, to appear.  
607 DOI: 10.1515/forum-2012-0152.
- 608 Cuevas, A. (2014). A partial overview of the theory of statistics with functional  
609 data. *J. Statist. Plann. Inference* 147, 1–23.
- 610 Cuevas, A., Fraiman, R. and Rodríguez-Casal, A. (2007). A nonparametric ap-  
611 proach to the estimation of lengths and surface areas. *Ann. Statist.* 35, 1031-  
612 1051.
- 613 Delicado, P. (2011). Dimensionality reduction when data are density functions.  
614 *Comp. Stat. Data Anal.* 55, 401–420.
- 615 Federer, H. (1959). Curvature measures. *Trans. Amer. Math. Soc.* 93, 418–491.
- 616 Ferraty, F. and Romain, Y., eds. (2011). *The Oxford Handbook of Functional*  
617 *Data Analysis*. Oxford University Press, Oxford.
- 618 Ferraty, F. and Vieu, P. (2006). *Nonparametric Functional Data Analysis: Theory*  
619 *and Practice*. Springer.
- 620 Folland, G.B. (1999). *Real Analysis. Modern Techniques and Their Applications*.  
621 Wiley, New York.
- 622 Galerne, B. (2011). Computation of the perimeter of measurable sets via their  
623 covariogram. Applications to random sets. *Image Anal. Stereol.* 30, 39-51.
- 624 Gibbs, A.L. and Su, F.E. (2002). On choosing and bounding probability metrics.  
625 *Int. Stat. Rev.*, 70, 3, 419-435.
- 626 Grenander, U. (1976). *Pattern Synthesis. Lectures in Pattern Theory. Volume 1*.  
627 Springer-Verlag, New York.
- 628 Grenander, U. (1981). *Abstract Inference*. Wiley, New York.
- 629 Hatcher, A. (2002). *Algebraic Topology*. Cambridge University Press.
- 630 Hobolt, A. and Vedel-Jensen, E.B. (2000). Modelling stochastic changes in curve  
631 shape, with an application to cancer diagnostics. *Adv. Appl. Prob.* 32, 344–362.
- 632 Hobolt, A., Pedersen, J. and Vedel-Jensen, E.B. (2003). A continuous parametric  
633 shape model. *Ann. Inst. Statist. Math.* 55, 227–242.
- 634 Horváth, L. and Kokoszka, P. (2012). *Inference for Functional Data with Appli-*  
635 *cations*. Springer, New York.
- 636 Kendall, D.G. (1989). A survey of the statistical theory of shape. *Statist. Sci.* 4,  
637 87–120.
- 638 Kendall, D.G., Barden, D., Carne, T.K. and Le, H. (1999). *Shape and Shape*  
639 *Theory*. Wiley.
- 640 Kendall, W.S. and Le, H. (2010). Statistical shape theory. In *New Perspectives in*  
641 *Stochastic Geometry*, W.S. Kendall and I. Molchanov, eds., pp. 348-373. Oxford  
642 University Press.
- 643 Kent, J.T. (1994). The complex Bingham distribution and shape analysis. *J.*  
644 *Royal Statist. Soc. B* 56, 285–299.

- 645 Kent, J.T. (1995). Current issues for statistical inference in shape analysis. In  
646 *Proceedings in Current Issues in Statistical Shape Analysis*, K.V. Mardia and  
647 C.A. Gill eds., pp. 167–175. Leeds University Press.
- 648 Lemke, P., Skiena, S.S. and Smith, W.D. (1995). Reconstructing sets from in-  
649 terpoint distances. In *Discrete and Computational Geometry. Algorithms and*  
650 *Combinatorics* Volume 25, B. Aronov, S. Basu, J. Pach y M. Sharir, eds., pp.  
651 597-631. Springer, New York.
- 652 Ling, H., Okada. K. (2007). An Efficient Earth Mover’s Distance Algorithm for  
653 Robust Histogram Comparison. *IEEE Transactions on Pattern Analysis and*  
654 *Machine Intelligence* 29, 840-853.
- 655 Lombarte, A., Chic, O., Parisi-Barabad, V., Olivella, R., Piera, J., García-  
656 Ladona, E. (2006). A web-based environment for shape analysis of fish otoliths.  
657 The AFORO database. *Sci. Mar.* 70, 147-152.
- 658 Mallows, C.L. and Clark, J.M.C. (1970). Linear-intercept distributions do not  
659 characterize plane sets. *J. Appl. Probability* 7, 240-244.
- 660 Marron, J.S. and Alonso, A.M. (2014). Overview of object oriented data analysis.  
661 *Biometrical Journal* 56, 732–753.
- 662 Martin, G.E. (1982). *Transformation Geometry. An Introduction to Symmetry.*  
663 Springer-Verlag. New York.
- 664 Matern, B. (1986). *Spatial Variation.* Lecture Notes in Statistics 36, Springer.  
665 New York.
- 666 Matheron, G. (1986). Le covariogramme gometrique des compacts convexes des  
667  $\mathbb{R}^2$ . *Technical report N- 2/86/G*, Centre de Gostatistique, Ecole Nationale Su-  
668 périeure des Mines de Paris.
- 669 Mattila, P. (1995). *Geometry of Sets and Measures in Euclidean Spaces: Fractals*  
670 *and Rectifiability.* Cambridge University Press. Cambridge.
- 671 Osada, R., Funkhouser, T., Chazelle, B. and Dobkin, D. (2002). Shape Distribu-  
672 tions. *ACM Transactions on Graphics* 21, 807-832.
- 673 Ramsay, J. O. and Silverman, B. W. (1997). *Functional Data Analysis.* Springer,  
674 New York.
- 675 Rubner, Y., Tomasi, C. and Guibas, L.J. (2000). The Earth Movers Distance as  
676 a metric for image retrieval. *Intl J. Computer Vision* 40, 99–121.
- 677 Storbeck, F. and Daan, B. (2001). Fish species recognition using computer vision  
678 and a neural network. *Fisheries Research* 51, 11-15.
- 679 Tebaldi, P., Bonetti, M., and Pagano, M. (2011). M statistic commands: inter-  
680 point distance distribution analysis. *The Stata Journal* 11, 271–289.
- 681 Tuset, V. M., Lombarte, A. and Assis, C. A. (2008) Otolith atlas for the western  
682 Mediterranean, north and central eastern Atlantic. *Scientia Marina*, 72, 7-198.
- 683 Villani, C. (2003). *Topics in Optimal Transportation.* Graduate Studies in Math-  
684 ematics, 58. American Mathematical Society, Providence, RI.

# Supplementary Material for “Shape classification based on interpoint distance distributions”

José R. Berrendero<sup>a</sup>, Antonio Cuevas<sup>a</sup>, Beatriz Pateiro-López<sup>b,\*</sup>

<sup>a</sup>*Departamento de Matemáticas, Universidad Autónoma de Madrid, Spain*

<sup>b</sup>*Departamento de Estadística e Investigación Operativa, Universidad de Santiago de Compostela, Spain*

---

---

## An application to fish family identification from otolith images

*Filled-in contour images.* The AFORO database contains both high resolution and filled-in contour images of otoliths, see Figure 1. The results in the study are obtained from the filled-in contour images.



Figure 1: High resolution image (left) and filled-in contour image (right) of a *Soleidae* otolith.

*Hierarchical clustering.* We consider the complete dataset consisting of otoliths from three families of fishes (*Soleidae*, *Labridae* and *Scombridae*) and apply an agglomerative hierarchical clustering procedure using both the  $L_1$  distance and the Wasserstein distance as dissimilarity criterion. In Figure 2 we show the dendrogram obtained using the  $L_1$  distance. If we cut the corresponding tree into three groups, we obtain that 94.29% of *Soleidae* otoliths belong to the first cluster, 96.80% of *Labridae* otoliths belong to the second cluster and 95.56% of *Scombridae* otoliths belong to the third cluster. See Table 1 for the complete table of results.

*k-means clustering.* In this section, we investigate the performance of the  $k$ -means clustering algorithm. We apply the  $k$ -means algorithm to each pair of families of otoliths ( $k=2$ ). We present the results based on the  $L_1$  distance between densities and the Wasserstein distance between distributions, see Table 2.

We observe that  $k$ -means performs reasonably well, except perhaps on the dataset consisting on *Labridae* and *Scombridae* otoliths (dataset C), see Table 2. The  $k$ -means algorithm highly depends on the initial centroids and this may be the reason for the not so good results in this dataset.

---

\*Corresponding address: Rúa Lope Gómez de Marzoa s/n. 15782 Santiago de Compostela. Spain

*Email address:* [beatriz.pateiro@usc.es](mailto:beatriz.pateiro@usc.es) (Beatriz Pateiro-López)

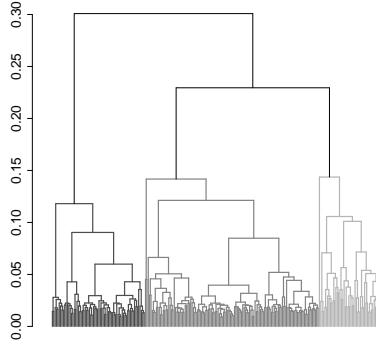


Figure 2: Dendrogram using the  $L_1$  distance between the estimated interpoint distance densities of the otoliths in the families *Soleidae*, *Labridae* and *Scombridae*. The tree is cut into three groups, represented in different tones of gray.

Table 1: Results of the hierarchical clustering procedure for *Soleidae*, *Labridae* and *Scombridae* otoliths. For each dissimilarity criterion, count and row percent of the true family labels versus the group labels for a partition into three clusters.

	$L_1$ distance			Wasserstein distance		
	Cluster 1	Cluster 2	Cluster 3	Cluster 1	Cluster 2	Cluster 3
<i>Soleidae</i>	66	4	0	58	12	0
	94.29%	5.71%	0.00%	82.85%	17.14%	0.00%
<i>Labridae</i>	2	121	2	0	123	2
	1.60%	96.80%	1.60%	0.00%	98.40%	1.60%
<i>Scombridae</i>	0	2	43	0	2	43
	0.00%	4.44%	95.56%	0.00%	4.44%	95.56%

Table 2: Results of the  $k$ -means algorithm ( $k = 2$ ). For each distance, contingency table (count and row percent) of the true family labels versus the group labels.

		$L_1$ distance		Wasserstein distance	
		Cluster 1	Cluster 2	Cluster 1	Cluster 2
Dataset A	<i>Soleidae</i>	68	2	65	5
		97.14%	2.86%	92.86%	7.14%
	<i>Labridae</i>	4	121	0	125
		3.20%	96.80%	0.00%	100.00%
Dataset B	<i>Soleidae</i>	69	1	67	3
		98.57%	1.43%	95.71%	4.29%
	<i>Scombridae</i>	0	45	0	45
		0.00%	100.00%	0.00%	100.00%
Dataset C	<i>Labridae</i>	123	2	101	24
		98.40%	1.60%	80.80%	19.20%
	<i>Scombridae</i>	2	43	5	40
		4.44%	95.56%	11.11%	88.89%

*Gobiidae and Serranidae otoliths.* At the beginning of the study, we had also considered two other large families: *Gobiidae* and *Serranidae* (see Figure 3 for examples of otoliths in these two families). As might be expected, the clustering methods did not perform well, for example, for the dataset containing *Gobiidae* and *Soleidae* otoliths. Note that the shape of some of the *Gobiidae* otoliths resembles that of the *Soleidae* otoliths. The same occurs for the dataset containing *Serranidae* and *Labridae* otoliths.

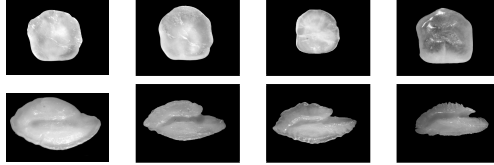


Figure 3: High resolution images of otoliths. First row: *Gobiidae*. Second row: *Serranidae*

### Interpoint distances and covariogram

In this section we will establish some simple relationships between the interpoint distance and the covariogram, a well-known tool in stochastic geometry. As a consequence, some properties of the interpoint distance distribution will result.

The covariogram of a bounded Borel set  $A \subset \mathbb{R}^d$  is defined by

$$K_A(y) = \mu(A \cap T_y A),$$

where  $y \in \mathbb{R}^d$ ,  $T_y A = A - y = \{a - y : a \in A\}$  and  $\mu$  denotes the Lebesgue measure in  $\mathbb{R}^d$ .

This function has proven to be useful in different problems of stochastic geometry and stereology. Some references are Cabo and Baddeley (1995, 2003) and Galerne (2011). First note that

$$K_A(y) = \int_{\mathbb{R}^d} \mathbb{I}_A(x) \mathbb{I}_A(x - y) dx,$$

so that,  $K_A$  can be alternatively expressed in terms of a convolution of two indicator functions,

$$K_A = \mathbb{I}_A * \mathbb{I}_{-A}, \quad (1)$$

where  $-A$  denotes the symmetric set  $-A = \{-x : x \in A\}$ .

Note that (1) is, up to a multiplicative constant, the density function of  $X_1 - X_2$ , where  $X_1, X_2$  are iid random variables uniform on  $A$ . As a conclusion,  $K_A$  fully determines the distribution of the interpoint distance  $Y_A$ .



Let us now briefly summarize some other relevant properties of this function; see, e.g. Lemmas 1.2, 1.3 and 1.4 in Cabo and Baddeley (1995) and Proposition 2 in Galerne (2011).

**Lemma 1.** *Let  $A \in \mathbb{R}^d$  be a bounded Borel set with covariogram  $K_A$ .*

(i) *For all  $y \in \mathbb{R}^d$ ,  $0 \leq K_A(y) \leq K_A(0) = \mu(A)$ . Moreover,  $K_A(y) = 0$  whenever  $\|y\| \geq \text{diam}(A)$ ,  $K_A(y) = K_A(-y)$  for all  $y \in \mathbb{R}^d$  and  $K_A$  is uniformly continuous on  $\mathbb{R}^d$ .*

(ii) *For any integrable  $f : [0, \infty) \rightarrow \mathbb{R}$ ,*

$$\int_A \int_A f(\|x - y\|) dx dy = \int_{\mathbb{R}^d} f(\|w\|) K_A(w) dw.$$

*This is the so-called ‘‘Borel’s overlap formula’’. Two interesting particular cases are obtained for  $f \equiv 1$  and  $f(t) = \mathbb{I}_{[0, \rho]}(t)/\mu(A)^2$ , leading respectively to*

$$\int_{\mathbb{R}^d} K_A(y) dy = \mu(A)^2 \tag{2}$$

*and*

$$\mathbb{P}\{Y_A \leq \rho\} = \frac{1}{\mu(A)^2} \int_{B(0, \rho)} K_A(y) dy, \text{ for } \rho > 0, \tag{3}$$

*where  $Y_A = \|X_1 - X_2\|$ ,  $X_1$  and  $X_2$  being independent random variables uniformly distributed on  $A$ .*

The following property of the interpoint distance distribution follows directly from Lemma 1.

**Proposition 1.** *Let  $X_1, X_2$  be independent random variables uniformly distributed on  $C$ . Denote  $Y_C = \|X_1 - X_2\|$ . Then,  $Y_C$  has a continuous density  $f_C$  with  $f_C(0) = 0$  and  $f_C(\rho_C) = 0$ , where  $\rho_C = \text{diam}(C)$ .*

*Proof.* Performing a change of variables to polar coordinates in (3) we have

$$\mathbb{P}\{\|X_1 - X_2\| \leq \rho\} = \frac{1}{\mu(C)^2} \int_0^\rho \int_0^{2\pi} r K_C(r \cos \theta, r \sin \theta) d\theta dr$$

Since  $K_C$  is continuous, we can differentiate under the integral sign to get that the distribution of the interpoint distance has the following continuous density

$$f_C(\rho) = \frac{1}{\mu(C)^2} \int_0^{2\pi} \rho K_C(\rho \cos \theta, \rho \sin \theta) d\theta, \text{ for all } \rho \in [0, 1].$$

In particular, for  $\rho = 0$  we get  $f_C(0) = 0$ . Also, from result (i) in Lemma 1,  $f_C(\rho_C) = 0$ .  $\square$

## References

- Cabo, A.J. and Baddeley, A.J. (1995). Line transects, covariance functions and set convergence. *Adv. Appl. Prob.* 27, 585-605.
- Cabo, A.J. and Baddeley, A.J. (2003). Estimation of mean particle volume using the set covariance function. *Adv. Appl. Prob.* 35, 27-46.
- Galerie, B. (2011). Computation of the perimeter of measurable sets via their covariogram. Applications to random sets. *Image Anal. Stereol.* 30, 39-51.