

Universally consistent estimation of the reach

Alejandro Cholaquidis¹, Ricardo Fraiman¹ and Leonardo Moreno²

¹ Centro de Matemáticas, Facultad de Ciencias, Universidad de la República, Uruguay.

² Instituto de Estadística, Departamento de Métodos Cuantitativos, FCEA, Universidad de la República, Uruguay.

Abstract

The reach of a set $M \subset \mathbb{R}^d$, also known as condition number when M is a manifold, was introduced by Federer in 1959. The reach is a central concept in geometric measure theory, set estimation, manifold learning, among others areas. We introduce a universally consistent estimate of the reach, just assuming that the reach is positive. Under an additional assumption we provide rates of convergence. We also show that it is not possible to determine, based on a finite sample, if the reach of the support of a density is zero or not. We provide a small simulation study and a bias correction method for the case when M is a manifold.

1 Introduction

The reach of a set $M \subset \mathbb{R}^d$, denoted by $\text{reach}(M)$, is a key concept in geometric measure theory; see Federer (1969); Rataj and Zajíček (2017); Rataj and Zähle (2001), and the references therein. It is also of importance in set estimation. We first focus on its relevance as a geometric concept, and then on the importance of shape constraints in set estimation, and in particular the positive reach condition.

Is defined as the largest distance from which any point outside M has a unique nearest point in M ; see Figure 1 (see also Definition 1). It can be proved that it is infinity for convex set. Positive reach impose some regularity conditions on the set. For instance, its boundary has Lebesgue measure zero, as it follows from the fact that the class of sets with positive reach is a subclass of the cone-convex sets (see Cuevas, Fraiman and Pateiro-López (2012) and Proposition 2 in Cholaquidis et al. (2014)). Moreover, the volume (i.e., its d -dimensional Lebesgue measure) of the set of points at distance $t \leq \text{reach}(M)$ is a polynomial of degree d on t for all $t \in [0, \text{reach}(M)]$ (see Federer (1959)). If $\text{reach}(M) > 0$, then it is also possible to define its second fundamental form; see Fu (1989). The positivity of the reach of the boundary, ∂M , of M allows its Minkowsky content to be defined (see Ambrosio, Colesanti and Villa (2008)), which is a notion of surface area suitable for estimation purposes (see Cuevas et al (2013)). It have been proved (see for instance Proposition 6.1 in Niyogi et al. (2008)) that the reach is an upper bound for the inverse of the curvature of an arc-length parametrized geodesic. In addition, as mentioned in Aamari et al (2021) “it prevents quasi self-intersection at scales smaller than the reach” (see Theorem 3.4 in Aamari et al (2019))

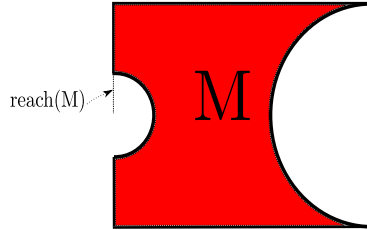


Figure 1: The reach of a set M , $\text{reach}(M)$ is the largest distance from which there exists a unique nearest point on M .

When the set is a smooth manifold, the reach parameter is also known as “condition number” (see for instance Niyogi et al. (2008)). It can also be proved that \mathcal{C}^2 compact manifolds with empty or \mathcal{C}^2 boundary, have positive reach; see Thäle (2008).

The reach has gained importance in the last two decades in the areas of statistics known as set estimation, manifold learning and persistent homology (see for instance Arias-Castro et al (2020); Cuevas, Fraiman and Pateiro-López (2016), and Horobeț et al. (2019), respectively). Given an unknown set $M \subset \mathbb{R}^d$ (not necessarily convex), set estimation deals with the problem of the estimation of M from a random sample $\{X_1, \dots, X_n\} \subset M$, as well as several geometrically important functionals associated with M . For instance, its Lebesgue measure, the Minkowsky content of its boundary, among others.

The starting and simplest problem is when M is the support of a distribution. This problem has been addressed by several authors where the ground-breaking Devroye-Wise estimator is a milestone (see Devroye and Wise (1980)), due to its universal consistency and simplicity. To obtain rates of convergence for this estimator, some geometric restrictions on the set are required (see Cuevas and Rodríguez-Casal (2004)). The best attainable rate for this estimator is of the order $\mathcal{O}(\log(n)/n)^{1/d}$.

To improve the rates of convergence, it is necessary to impose stronger shape restrictions than the one required in Cuevas and Rodríguez-Casal (2004) on the set. A well-known shape restriction is r -convexity. A set is said to be r -convex if it equals the complement of open balls of radius r not meeting the set. r -convexity is one of the most studied shape restriction in set estimation (see for instance, Pateiro-López and Rodríguez-Casal (2012); Rodríguez-Casal (2007); Rodríguez-Casal and Saavedra-Nieves (2016), and the references therein).

Positive reach is a stronger condition than r -convexity (see Cuevas, Fraiman and Pateiro-López (2012)). Therefore, when a ball of positive radius rolls inside the set (i.e there exists $\lambda > 0$ such that for all $x \in \partial M$, there exists $y \in M$ such that $x \in \partial \mathcal{B}(y, \lambda) \subset M$), and the set has reach r , the estimator proposed in Rodríguez-Casal and Saavedra-Nieves (2021) estimates the reach. Also, if the set and its complement are r -convex, then the set has positive reach (see Lemma A.0.7 in Pateiro-Lopez (2008) together with Walther, G. (1999)).

A second stage in set estimation is to estimate the level sets of the distribution (see Cadre, B. (2006); Cuevas et al (2006)), as well as the boundary of the set (see Cuevas and Rodríguez-Casal (2004)). In this setting, Cuevas and Fraiman (1997) proposes to estimate the support of the distribution by means of a kernel-based density estimator.

Next, and getting closer to our problem, the interest is on some functional of the set, such as the d -dimensional Lebesgue measure of the set, as well as the measure of its boundary; see for instance Cuevas et al (2007).

Recently, the study of statistical methods for manifold valued data (known as manifold learning) has gained attention, due to its application in dimension reduction (among others). The aim is to recover a lower dimensional structure from the data, see for instance Aamari and Levrard (2018); Fefferman et al. (2016); Genovese et al. (2012a,b); Niyogi et al. (2008) and the references therein, or a functional of it, see for instance Aaron and Cholaquidis (2020); Aaron, et al (2017). Several classical problems have been tackled in this setting, such as density estimation. Here again, the reach plays a key role as a shape restriction. In Aamari et al (2019), an estimator of the reach is proposed for manifold valued data with “the key assumption that both the point cloud and the tangent spaces were jointly observed”. For this “oracle framework [...] it is showed to achieve uniform expected loss bounds over a \mathcal{C}^3 -like model” and “upper and lower bounds on the minimax rate for estimating the reach

are obtained". Also for manifold valued data a different estimator of the reach was introduced recently in Berenfeld et al (2022), where the manifold is assumed to be at least of class C^3 , and it must be previously estimated with the manifold estimator proposed in Aamari and Levrard (2019). The rates of convergence obtained (in probability and in L^1) are of order $(\log(n)/n)^{k/(2d)}$ if the manifold $M \subset \mathbb{R}^d$ is of class C^k , for $k \geq 4$, and $(\log(n)/n)^{1/d}$, for $k = 3$. The two aforementioned estimators requires the manifold to have no boundary. This assumption is not required in our proposal, which is, up to our knowledge, the only consistent estimator proposed in this setup. In Aamari et al (2022) an estimator of the reach is also proposed. They obtain better convergence rates, but with stronger hypotheses and a non-computable estimator.

In what follows, we will study the problem of estimating the reach of a set, looking for universally consistent estimators, that is, assuming only that the set has positive reach. The rate of convergence obtained depends on the Hausdorff distance between the sample and the set. No assumptions are made on the distribution of the sample.

Our estimation procedure is based on an equivalent definition of reach given in Boissonnat et al (2019), which provides a new nice geometrical interpretation of the reach. In Theorem 2, we prove the universal consistency of the estimator. With an additional assumption in Corollary 1, we derive a convergence rate for the proposed reach estimator. In Section 4, we prove that it is not possible to determine based on a finite sample if the reach is zero or not. In Section 5 we report the results of a small simulation study, and in Subsection 5.2, we introduce a bias correction method for the case where M is a manifold.

2 Main definitions and geometric results

We start by fixing some notation to be used throughout the manuscript. Given a set $M \subset \mathbb{R}^d$, we denote by ∂M and $\text{int}(M)$, the boundary and interior of M , respectively. In what follows, we assume that M is compact.

We denote by $\|\cdot\|$ the Euclidean norm in \mathbb{R}^d . A closed ball of radius $\varepsilon > 0$ centred at x is denoted by $\mathcal{B}(x, \varepsilon)$, and an open ball is denoted by $\mathring{\mathcal{B}}(x, \varepsilon)$. Given $\varepsilon > 0$ and a set $A \subset \mathbb{R}^d$, $B(A, \varepsilon)$ denotes the parallel set $B(A, \varepsilon) = \{x \in \mathbb{R}^d: d(x, A) \leq \varepsilon\}$, where $d(x, A) = \inf\{\|x - a\|: a \in A\}$. The d -dimensional Lebesgue measure on \mathbb{R}^d of a set M is denoted by $|M|_d$. Given two non-empty compact sets $A, C \subset \mathbb{R}^d$, the Hausdorff distance between A and C is defined as

$$d_H(A, C) = \max \left\{ \max_{a \in A} d(a, C), \max_{c \in C} d(c, A) \right\}.$$

Given a continuous curve $\gamma: [0, T] \rightarrow M$, we define its length as

$$l(\gamma) = \sup_P \sum_i \|\gamma(t_{i+1}) - \gamma(t_i)\|,$$

where the supremum is over all finite partitions of $[0, T]$. Given $x, y \in M$ we define the geodesic distance between them as $d_M(x, y) = \inf_\gamma l(\gamma)$, where the infimum is over all continuous curves joining x and y . In what follows, we assume that M is geodesically convex; that is, for any two points $x, y \in M$ there exists a geodesic connecting them, with length $d_M(x, y)$ (see Bernstein et al (2000)).

Following the notation in Federer (1959), let $\text{Unp}(M)$ be the set of points $x \in \mathbb{R}^d$ with a unique closest point on M .

Definition 1. For $x \in M$, let $\underline{\text{reach}}(M, x) = \sup\{r > 0 : \mathring{\mathcal{B}}(x, r) \subset \underline{\text{Unp}}(M)\}$. The reach of M is defined by $\underline{\text{reach}}(M) = \inf\{\underline{\text{reach}}(M, x) : x \in M\}$, and M is said to be of positive reach if $\underline{\text{reach}}(M) > 0$.

Theorem 1 in Boissonnat et al (2019), states that when M is closed $\text{reach}(M)$ equals

$$\sup\left\{r > 0, \forall a, b \in M, \|a - b\| < 2r \Rightarrow d_M(a, b) \leq 2r \arcsin\left(\frac{\|a - b\|}{2r}\right)\right\}. \quad (1)$$

Another important geometric restriction, which is required to get the convergence rate of the estimator in the iid case, is the standardness, see Cuevas and Rodriguez-Casal (2004).

Definition 2. A set $M \subset \mathbb{R}^d$ is said to be standard with respect to a Borel measure ν in a point x if there exists $\lambda > 0$ and $\eta > 0$ such that

$$\nu(\mathcal{B}(x, \varepsilon) \cap M) \geq \eta |\mathcal{B}(x, \varepsilon)|_d, \quad 0 < \varepsilon \leq \lambda. \quad (2)$$

A set $M \subset \mathbb{R}^d$ is said to be standard if (2) hold for all $x \in M$.

Let $\mathcal{X}_n = \{X_1, \dots, X_n\} \subset M$ be a finite set, and $G_n = (\mathcal{X}_n, E_n)$ be a graph with vertex in \mathcal{X}_n and edges E_n .

Define as in Bernstein et al (2000)

$$d_{G_n}(x, y) = \min_P \sum_{i=1}^{p-1} \|X_{j_{i+1}} - X_{j_i}\|,$$

where $P = (X_{j_1}, \dots, X_{j_p}) \subset \mathcal{X}_n$ varies over all paths along the edges of G_n connecting $x = X_{j_1}$ to $y = X_{j_p}$.

Definition 3. Let $\{\varepsilon_n\}_n$ a sequence of strictly positive real numbers such that $\varepsilon_n \rightarrow 0$, let us define the graph $G_n = (\mathcal{X}_n, E_n)$ such that $(X_i, X_j) \in E_n$ if and only if $\|X_i - X_j\| \leq \varepsilon_n$.

The following theorem, whose proof is given in the Appendix, states that d_{G_n} approximates d_M . The proof follows closely the ideas used in Bernstein et al (2000). It is included because we state our results in terms of the reach of the set, and not the curvature, as it is done in Bernstein et al (2000), and the points where there are changes are clarified.

Theorem 1. With the notation introduced before, let M be a compact, geodesically convex set, such that its reach, r_0 , is strictly positive. Let $\{\varepsilon_n\}_n$ be a sequence of strictly positive real numbers such that $\tau_n = d_H(\mathcal{X}_n, M)/\varepsilon_n \rightarrow 0$, where $\varepsilon_n \rightarrow 0$, and n is large enough to guarantee $\varepsilon_n < 2r_0$. Then, for all $x, y \in M$,

$$\left(1 - \frac{1}{24} \left(\frac{\pi \varepsilon_n}{2r_0}\right)^2\right) d_M(x, y) \leq d_{G_n}(x, y) \leq (1 + 4\tau_n) d_M(x, y), \quad (3)$$

where $d_{G_n}(x, y)$ is the distance in the graph build previously, but including x and y as vertices.

3 Reach estimation

Given a set $\mathcal{X}_n = \{X_1, \dots, X_n\} \subset M$, and a sequence of positive real numbers $\{\varepsilon_n\}_n$, the plug-in estimator of $\text{reach}(M)$ based on the graph G_n defined in the previous section, is

$$\hat{r}_n = \sup\left\{r > 0, \forall X_i \neq X_j \in \mathcal{X}_n, \|X_i - X_j\| < 2r \Rightarrow d_{G_n}(X_i, X_j) \leq 2r(1 + \varepsilon_n^2) \arcsin\left(\frac{\|X_i - X_j\|}{2r}\right)\right\}.$$

In Figure 2, it is shown as a green dotted line $d_{G_n}(a, b)$, the distance between two points $a, b \in M$ (with $\|a - b\| < 2r$) in the graph G_n built as before for a given ε_n . As a black solid line, it is shown an arc in the circle of radius r joining a and b ($a, b \in \mathcal{X}_n$), whose length is $2r \arcsin(\|a - b\|/2r) \geq d_{G_n}(a, b)$.

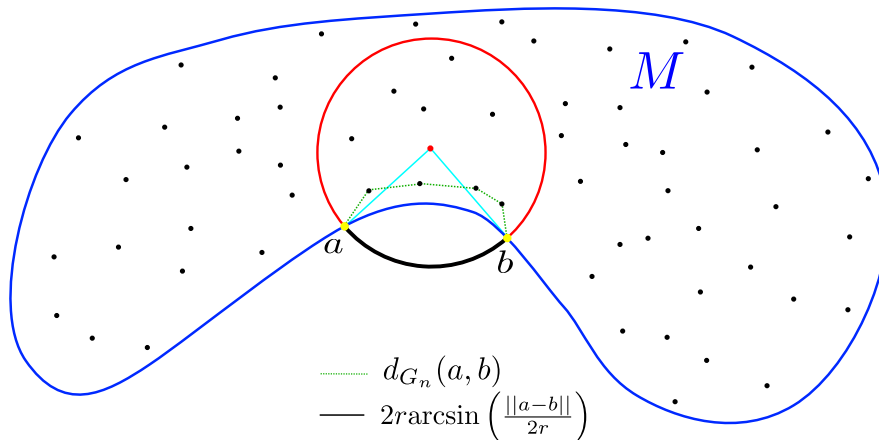


Figure 2: As a green dotted line, it is shown $d_{G_n}(a, b)$. As a black solid line, it is shown as the arc of length $2r \arcsin(\|a - b\|/2r)$ joining two points $a, b \in \mathcal{X}_n$.

The following theorem states that the estimator \hat{r}_n is, for n large enough, bounded from below by the reach r_0 of M , and bounded from above by $r_0/(1 - \varepsilon_n)$.

Theorem 2. *Under the hypotheses of Theorem 1, for all n large enough, we have*

$$r_0 \leq \hat{r}_n \leq \frac{r_0}{1 - \varepsilon_n}, \quad (4)$$

where ε_n is a sequence of strictly positive real numbers such that $\varepsilon_n \rightarrow 0$ and $\delta_n/\varepsilon_n^3 \rightarrow 0$, being $\delta_n = d_H(\mathcal{X}_n, M)$.

If we assume that $\text{reach}(M) > 0$, and M is standard (see Definition 2), we get the following corollary, regarding the convergence rate of \hat{r}_n , which is a consequence of Theorem 3 in Cuevas and Rodriguez-Casal (2004) and Theorem 2.

Corollary 1. *Let X_1, X_2, \dots be a sequence of iid observation drawn from a distribution P_X on \mathbb{R}^d . Assume that the support M of P_X is compact, standard with respect to P_X , and has reach $r_0 > 0$. Then, with probability one, for n large enough,*

$$r_0 \leq \hat{r}_n \leq r_0 + c \left(\frac{\log(n)}{n} \right)^{\frac{1}{3d}} \beta_n, \quad (5)$$

for $\varepsilon_n = c(\log(n)/n)^{\frac{1}{3d}} \beta_n$, being c any constant larger than $(2/(\eta\omega_d))^{1/d}$, where $\omega_d = |\mathcal{B}(0, 1)|_d$, and η is the standardness constant, and $\beta_n \rightarrow +\infty$ any sequence.

Remark 1. *The standardness hypothesis (defined in (2)) in Corollary 1 is not fulfilled when M is a dt -dimensional Riemannian manifold, $dt < d$, and P_X is a distribution supported on M . In that case, instead of assumption (2), the required condition for P_X and M is the following: there exists $\lambda > 0$ and $\eta > 0$ such that for all $x \in M$ $P_X(\mathcal{B}(x, \varepsilon)) \geq \eta |\mathcal{B}(x, \varepsilon)|_d$ for all $0 < \varepsilon < \lambda$, $\mathcal{B}(x, \varepsilon)$ being the ball w.r.t. the Riemannian metric, and $|\cdot|_d$ being the dt -dimensional Lebesgue measure. With the same ideas used to prove Theorem 3 in Cuevas and Rodriguez-Casal (2004), it can be proved that $d_H(\mathcal{X}_n, M) = \mathcal{O}((\log(n)/n)^{1/dt})$, and then the rate (5) is improved up to*

$$r_0 \leq \hat{r}_n \leq r_0 + c(\log(n)/n)^{1/(3dt)} \beta_n,$$

for $\varepsilon_n = c(\log(n)/n)^{\frac{1}{3d}}\beta_n$, being c any constant larger than $(2/(\eta\omega_d))^{1/d}$, where ω_d is the d -dimensional Lebesgue measure of a ball of radius one in \mathbb{R}^d , η is the standardness constant, and $\beta_n \rightarrow +\infty$ is any sequence of positive real numbers.

4 Non-estimability

In this section we will prove that it is not possible to determine, based on a finite sample, if the reach of the support, $\text{supp}(f)$, of a density f is zero or not. This is equivalent to showing that the functional

$$\alpha(f) = \begin{cases} 1 & \text{if } \text{reach}(\text{supp}(f)) = 0 \\ 0 & \text{otherwise,} \end{cases}$$

cannot be consistently estimated from a sequence of iid observations. This is the case of several different problems. The most simple and well-known is that one cannot determine if the mean of a distribution is finite or not based on a finite sample. However, as in our case, once assuming that the mean is finite, it can be consistently estimated. A fundamental contribution on this general problem is given in Le Cam and Schwartz (1960), where they provide necessary and sufficient conditions for the existence of consistent estimators.

To prove that $\alpha(f)$ cannot be consistently estimated, we will make use Lemma 1.1 in Fraiman and Meloche (1999).

Lemma 1 (Lemma 1.1 in Fraiman and Meloche (1999)). *Let \mathcal{T} a subset of densities equipped with some norm $\|\cdot\|$ that makes $(\mathcal{T}, \|\cdot\|)$ a complete metric space. Assume that $\|f\|_1 \leq c\|f\|$ for some constant c and all $f \in \mathcal{T}$. Let $\Phi : \mathcal{T} \rightarrow [-K, K]$ be any bounded characteristic of the densities in \mathcal{T} . If Φ is consistently estimable on \mathcal{T} , then there exists a dense subset of points in \mathcal{T} at which Φ is continuous with respect to the topology induced by the $\|\cdot\|$ norm. Therefore, if Φ is discontinuous at every point in \mathcal{T} , it is not consistently estimable in \mathcal{T} .*

In our case, \mathcal{T} is the set of densities (w.r.t the Lebesgue measure) endowed with the L^1 norm, $\|\cdot\|_1$, which is, by Scheffe's lemma, a complete space.

Theorem 3. *Let f be a density such that $\underline{\text{reach}}(\text{supp}(f)) > 0$, then for all $\delta > 0$ there exists f_δ a density such that $\|f - f_\delta\|_1 < \delta$ and $\underline{\text{reach}}(\text{supp}(f_\delta)) = 0$. Moreover, if $\underline{\text{reach}}(\text{supp}(f)) = 0$, then for all $\varepsilon > 0$ there exists f_ε such that $\|f_\varepsilon - f\|_1 < \varepsilon$ and $\underline{\text{reach}}(\text{supp}(f_\varepsilon)) > 0$. Thus, the functional α is discontinuous at any density in \mathcal{T} , and therefore from Lemma 1, it is not estimable.*

Proof. If $\text{reach}(\text{supp}(f)) > 0$ then from Propositions 1 and 2 in Cuevas, Fraiman and Pateiro-López (2012) $\text{supp}(f)$ fulfills the exterior rolling ball condition (see section 2 in Cuevas, Fraiman and Pateiro-López (2012)). It is easy to see that this imply that $\text{supp}(f)$ is ρ, h -cone-convex for some $\rho \in (0, \pi/4]$ and $h > 0$ (see Definition 4 in Cholaquidis et al. (2014)), then from Proposition 2 in Cholaquidis et al. (2014) it follows that $|\partial\text{supp}(f)|_d = 0$, from where it follows that $\text{int}(\text{supp}(f)) \neq \emptyset$. Let $x \in \text{int}(\text{supp}(f))$, then we can remove from $\text{int}(\text{supp}(f))$ an arbitrary small open cone T centred at x , whose interior is included in $\text{int}(\text{supp}(f))$. Define $f_\delta = cf\mathbb{1}_{T^c}$ as being c a normalizing constant, and T^c the complement of the set T . Clearly $\text{reach}(\text{supp}(f_\delta)) = 0$. Let us assume now that $\text{reach}(\text{supp}(f)) = 0$. Let $\varepsilon > 0$ and $K_1 = K_1(\varepsilon)$ a compact set such that $\int_{K_1^c} f(x)dx < \varepsilon/4$. Let $\kappa > 0$ small enough such that

$$\int f(x)\mathbb{1}_{f^{-1}([0, \kappa])}dx < \varepsilon/4.$$

Table 1: Mean, median and standard deviation over 100 replication, of \hat{r}_n for different values of the reach r and n .

n	Inner radius						ε_n
	$r=0.25$			$r=0.5$			
	mean	median	sd	mean	median	sd	
500	0.262	0.265	0.005	0.503	0.504	0.002	0.44
750	0.258	0.261	0.003	0.502	0.503	0.002	0.41
1000	0.256	0.258	0.002	0.502	0.502	0.001	0.40
1250	0.256	0.258	0.002	0.502	0.500	0.001	0.39
1500	0.255	0.255	0.002	0.501	0.500	0.001	0.37

Define $K = K_1 \cap \overline{f^{-1}([\kappa, +\infty))}$ and $\xi(x) = d(x, K^c)\mathbb{I}_K(x)$. Let γ small enough such that $|K \setminus \xi^{-1}([\gamma, +\infty))|_d < \varepsilon/2$. Let g a \mathcal{C}^2 function such that $\sup_{x \in K} |\xi(x) - g(x)| < \gamma/2$. Let $S_\gamma := g^{-1}([\gamma/2, \infty))$ then $S_\gamma \subset K$, is compact and \mathcal{C}^2 , and then its boundary is has positive reach, (see Thäle (2008)). The density $cf\mathbb{I}_{S_\gamma}$ fulfills the desired properties, where the constant c is chosen to integrate 1. \square

5 A small simulation study

5.1 Example 1

In this first example, to study the performance of \hat{r}_n , we considered the two dimensional sets,

$$M_r = \{(x, y) \in \mathbb{R}^2 : r^2 < x^2 + y^2 < r^2 + 1/\pi\},$$

for $r \in \{0.25, 0.5\}$. Then, $\text{reach}(M_r) = r$. For each r fixed, we drawn a sample of n points uniformly distributed on M_r , for $n \in \{500, 750, 1000, 1250, 1500\}$. The values of \hat{r}_n are computed using $\varepsilon_n = (\max_i \min_{j \neq i} \|X_i - X_j\|)^{1/2}$, for a constant $c > (4/\pi)^{1/2}$. The whole procedure is replicated 100 times. The mean, median and standard deviation of these replications are shown in Table 1. As is observed in both scenarios, the performance of our estimator improves when the sample sizes increases. In addition, the deviation is smaller for $r = 0.5$ than for $r = 0.25$. There is a small overestimation of the reach, expected according to Theorem 1, which is smaller for $r = 0.5$.

5.2 Bias correction

The overestimation of the reach is clearly established by Theorem 1, and is also confirmed by the results in Table 1. This problem becomes worse in the case of manifold-valued data, where we conjecture that it is due to the fact that d_{G_n} tends to underestimate d_M when M is a lower-dimensional manifold. Thus, we consider a bias correction, following the proposal given in Arias-Castro et al (2019) for volume estimation, adapted for our problem. It worth to be mention that the proposal in Arias-Castro et al (2019) has been shown to be minimax optimal under the assumption that the data is uniformly distributed on the manifold. However, the consistency of the bias correction goes far beyond the scope of this manuscript. We show the performance through a simulation study.

We split the set \mathcal{X}_n into two subsets \mathcal{X}_1 and \mathcal{X}_2 of sizes n_1 and n_2 , respectively.

1. Calculate \hat{r}_{n_1} based on the set \mathcal{X}_1 .

2. Calculate \hat{p}_n , defined by

$$\binom{n_2}{2}^{-1} \left| \left\{ (X_i, X_j) \in \mathcal{X}_2 : i \neq j, \|X_i - X_j\| < 2\hat{r}_{n_1} \text{ but } d_{G_{n_2}}(X_i, X_j) > 2\hat{r}_{n_1}(1 - \varepsilon_{n_2}^2) \arcsin\left(\frac{\|X_i - X_j\|}{2\hat{r}_{n_1}}\right) \right\} \right|.$$

where $|\cdot|$ denotes the cardinality of the set.

3. Output: $\hat{r}_n = [(1 - \hat{p}_n) \vee 1/2]\hat{r}_{n_1}$.

5.2.1 Example 2

In this example, the set is $M = \{(x, y) \in \mathbb{R}^2 : x^2 + 4y^2 = 1, x \geq 0\}$, which is a manifold with boundary; see Figure 3. It is easy to see that $\text{reach}(M) = 0.25$. We compare the performance of our estimator and the estimator proposed in Aamari et al (2019). Its first approach requires to know the tangent spaces at the sample points. More precisely, given an iid sample $\mathcal{X}_n \subset M$, M being a d -dimensional manifold in \mathbb{R}^D , the authors proposes the estimator r given by

$$r = \inf_{X_i \neq X_j \in \mathcal{X}_n} \frac{\|X_j - X_i\|^2}{2d(X_j - X_i, T_{X_i}M)}, \quad (6)$$

where $T_{X_i}M$ is the d affine linear space, tangent to M at X_i .

If M is unknown, then $T_{X_i}M$ must be estimated. Proposition 6.1 in Aamari et al (2019) provides bounds for the stability of the estimator τ when the tangent spaces are perturbed. Let $TM = \{\widehat{T_{X_i}M}\}_{X_i \in \mathcal{X}_n}$ be a family of d affine linear spaces, let

$$r(\mathcal{X}_n, T) = \inf_{X_i \neq X_j \in \mathcal{X}_n} \frac{\|X_j - X_i\|^2}{2d(X_j - X_i, \widehat{T_{X_i}M})}. \quad (7)$$

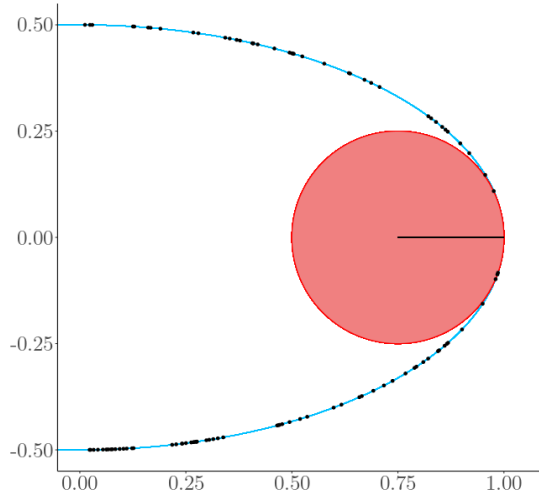


Figure 3: The half-ellipse $M = \{(x, y) \in \mathbb{R}^2 : x^2 + 4y^2 = 1, x \geq 0\}$, whose reach is 0.25, and a sample on M .

To compare with the proposals given in Aamari et al (2019), we consider $S_1 = \{S_{1, X_i}\}_{X_i \in \mathcal{X}_n}$, defined for each $X_i \in \mathcal{X}_n$, as follows: S_{1, X_i} is the estimation of $T_{X_i}M$, by means of local principal components, see Aamari et al (2021).

We draw samples of size n in the half-ellipse, for $n \in \{400, 600\}$. The samples were generated as follows: first we draw a sample (X, Y) with standard normal bivariate distribution, then we take $(|X|, Y)$. Finally we project $Z = (|X|, Y)$ onto the ellipse considering $Z/\|Z\|_e$ with $\|(x, y)\|_e = \sqrt{x^2 + 4y^2}$.

We perform 100 replications. For each of them we calculate 1) $r_{1,n} = r(\mathcal{X}_n, S_1)$ based on equation (7), 2) our estimator, $\hat{r}_{1,n}$, without bias correction, and 3) $\hat{r}_{2,n}$, the estimator with the bias correction proposed in Section 5.2. We take $\varepsilon_n = (\max_i \min_{j \neq i} \|X_i - X_j\|)^{1/2}$, while we follow the suggestion in Aamari et al (2019) considering only those pairs of points whose distance is at least δ with δ of order $\log(n)/n$.

The boxplots for these estimators are given in Figure 4. The average differences between the true tangent subspaces TM and their approximations S_1 (i.e., $\text{Error}(k) = \max_{X_i \in \mathcal{X}_n} \|T_{X_i}M - S_{1, X_i}\|_{op}$, where k stands for the k -th replicate) is 0.078 and 0.054 for $n = 400$ and $n = 600$ respectively. Figure 4 shows how the bias is reduced when we use the bias-corrected estimator $\hat{r}_{2,n}$, which also performs better than $r_{1,n}$.

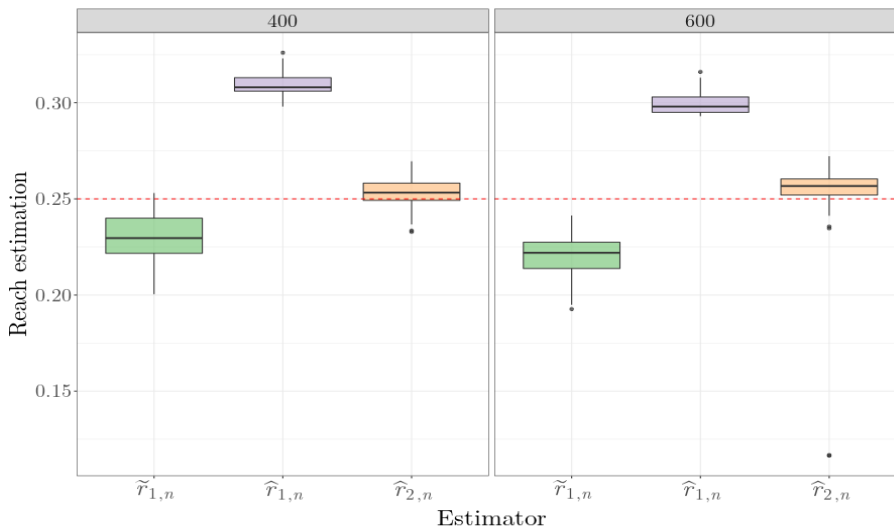


Figure 4: Estimators boxplots for $r_{1,n} = r(\mathcal{X}_n, S_1)$, $\hat{r}_{1,n}$ without bias correction and $\hat{r}_{2,n}$ the estimator with the bias correction.

6 Concluding remarks

We deal with an important problem in set estimation related to geometric measure theory: the estimation of the reach, r_0 , of a set $M \subset \mathbb{R}^d$, which includes the case where M is a manifold. A universally consistent estimator of r_0 , in the sense that no assumptions are required except for r_0 being positive for is proposed. We show that based on a finite sample with density f w.r.t. the Lebesgue measure, it is not possible to determine if the reach of the support of f is zero or not. The consistency result is related to the convergence to zero of $d_H(M, \mathcal{X}_n)$. We conjecture that given any consistent estimator of the reach, r_n , and given any sequence $\beta_n \rightarrow 0$, it is possible to find a set (depending on β_n), with positive reach, for which, r_n converge to the reach at a rate slower than β_n .

However, under the weak additional assumption of standardness, we provide rates of convergence for the proposed estimator. For the case where M is a manifold, we adapt our procedure by adding a bias correction method. An alternative way to deal with this problem is to consider the estimator $r_n = (1 - \alpha\varepsilon_n)\hat{r}_n$, but the optimal choice of the parameter α should be addressed.

The results obtained in our small simulation section are promising. However, a more extended simulation study, to compare with other proposals in the literature is an interesting problem to be addressed.

Appendix

6.1 Proof of Theorem 2

Proof. To prove the first inequality in (4), let $X_i, X_j \in \mathcal{X}_n$ such that $\|X_i - X_j\| < 2r_0$. Let n be large enough such that $(1 + 4\tau_n) \leq (1 + \varepsilon_n^2)$. From (1)

$$\begin{aligned} d_{G_n}(X_i, X_j) &\leq (1 + 4\tau_n)d_M(X_i, X_j) \leq 2r_0(1 + 4\tau_n)\arcsin\left(\frac{\|X_i - X_j\|}{2r_0}\right) \\ &\leq 2r_0(1 + \varepsilon_n^2)\arcsin\left(\frac{\|X_i - X_j\|}{2r_0}\right). \end{aligned}$$

Then $\hat{r}_n \geq r_0$. To prove the second inequality in (4), we have to prove that for all $x, y \in M$ such that $\|x - y\| \leq 2\hat{r}_n(1 - \varepsilon_n)$,

$$d_M(x, y) \leq 2\hat{r}_n(1 - \varepsilon_n)\arcsin\left(\frac{\|x - y\|}{2\hat{r}_n(1 - \varepsilon_n)}\right).$$

Since $r_0 \leq \hat{r}_n$, it is enough to verify previous inequality for all $x, y \in M$ such that $\|x - y\| \geq 2r_0(1 - \varepsilon_n)$. This condition will be used later on in the proof, to guarantee that the value of n for which previous inequality holds, does not depend on the pair x, y .

Denote by

$$\gamma_n = \left(1 - \frac{1}{24} \left(\frac{\pi\varepsilon_n}{2\text{reach}(M)}\right)^2\right).$$

Let $X_i, X_j \in \mathcal{X}_n$ be the closest points, wrt the euclidean distance, to x and y , respectively. From $\delta_n/\varepsilon_n^3 \rightarrow 0$, it follows that, for n large enough,

$$\|X_i - X_j\| \leq \|X_i - x\| + \|x - y\| + \|y - X_j\| \leq 2\delta_n + 2\hat{r}_n(1 - \varepsilon_n) < 2\hat{r}_n.$$

Then, for n large enough

$$d_{G_n}(X_i, X_j) \leq 2\hat{r}_n(1 + \varepsilon_n^2)\arcsin\left(\frac{\|X_i - X_j\|}{2\hat{r}_n}\right). \quad (8)$$

From Equation (3), (recall that $d_{G_n}(x, y)$ is the distance in the graph, including x and y as vertices)

$$\gamma_n d_M(x, y) \leq d_{G_n}(x, y) \leq d_{G_n}(x, X_i) + d_{G_n}(X_i, X_j) + d_{G_n}(X_j, y). \quad (9)$$

Since $\|x - X_i\| < \delta_n < \varepsilon_n$, $d_{G_n}(x, X_i) = \|x - X_i\| < \delta_n$. In the same manner, $d_{G_n}(y, X_j) < \delta_n$. Since $\gamma_n \rightarrow 1$ we can take n large enough such that $2\delta_n/\gamma_n < 3\delta_n$, then from (8) and (9),

$$d_M(x, y) \leq 3\delta_n + 2\hat{r}_n \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin\left(\frac{\|X_i - X_j\|}{2\hat{r}_n}\right).$$

Let x, y such that $\|x - y\| < 2\hat{r}_n(1 - \varepsilon_n)$, then

$$\begin{aligned} \frac{\|X_i - X_j\|}{2\hat{r}_n} &\leq \frac{\|x - y\| + 2\delta_n}{2\hat{r}_n} = \frac{\|x - y\|}{2\hat{r}_n(1 - \varepsilon_n)} + \frac{\|x - y\|}{2\hat{r}_n} \left(1 - \frac{1}{1 - \varepsilon_n}\right) + \frac{\delta_n}{\hat{r}_n} = \\ &\frac{\|x - y\|}{2\hat{r}_n(1 - \varepsilon_n)} - \varepsilon_n \frac{\|x - y\|}{2\hat{r}_n(1 - \varepsilon_n)} + \frac{\delta_n}{\hat{r}_n}. \end{aligned}$$

Using that $\arcsin(a - b) \leq \arcsin(a) - \arcsin'(a - b)b$, for all $a, b > 0$ such that $0 \leq a - b < 1$, (which follows from the fact that $\arcsin'(t) > 1$ for all $0 < t < 1$), for $a = \|x - y\|/(2\hat{r}_n(1 - \varepsilon_n))$ and $b = \varepsilon_n a - \delta_n/\hat{r}_n$, we get

$$\begin{aligned} d_M(x, y) &\leq 2\hat{r}_n \frac{(1 + \varepsilon_n^2)}{\gamma_n} \left[\arcsin(a) - \arcsin'(a - b)b \right] + 3\delta_n = \\ &2\hat{r}_n(1 - \varepsilon_n)\arcsin(a) + 2\hat{r}_n \left(\frac{(1 + \varepsilon_n^2)}{\gamma_n} - 1 + \varepsilon_n \right) \arcsin(a) \\ &\quad - 2\hat{r}_n \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a - b) \left(\varepsilon_n a - \frac{\delta_n}{\hat{r}_n} \right) + 3\delta_n. \end{aligned} \quad (10)$$

Using that $a < 1$, $\hat{r}_n \geq r_0 > 0$, and $\gamma_n \rightarrow 1$, we get that, for n large enough,

$$2\hat{r}_n \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a - b) \frac{\delta_n}{\hat{r}_n} \leq 3\arcsin'(1 - \varepsilon_n + \delta_n/r_0)\delta_n \leq C_1 \frac{\delta_n}{\sqrt{\varepsilon_n}} \leq \varepsilon_n^2, \quad (11)$$

C_1 being a positive constant. In addition, for n large enough, using that $\hat{r}_n \leq \text{diam}(M)$, $\arcsin(a) \leq \pi/2$, and $\gamma_n \rightarrow 1$,

$$\hat{r}_n \left(\frac{(1 + \varepsilon_n^2)}{\gamma_n} - 1 \right) \arcsin(a) \leq \frac{\text{diam}(M)\pi}{4} (1 - \gamma_n + \varepsilon_n^2) \leq C\varepsilon_n^2. \quad (12)$$

C being a positive constant. Then, using (11) and (12) in (10), and the fact that $\delta_n/\varepsilon_n^3 \rightarrow 0$, we get that, for n large enough,

$$\begin{aligned} d_M(x, y) &\leq 2\hat{r}_n(1 - \varepsilon_n)\arcsin(a) + C'\varepsilon_n^2 + \\ &2\hat{r}_n\varepsilon_n \left[\arcsin(a) - \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a(1 - \varepsilon_n) + \delta_n/\hat{r}_n)a \right], \end{aligned}$$

C' being a positive constant. To prove the second inequality in (4), it is enough to prove that there exists $k > 0$ such that for n large enough,

$$\arcsin(a) - \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a(1 - \varepsilon_n) + \delta_n/\hat{r}_n)a < -k. \quad (13)$$

Observe that $k = 0$ is not enough to conclude because, if that is the case it may happen that

$$\arcsin(a) - \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a(1 - \varepsilon_n) + \delta_n/\hat{r}_n)a \rightarrow 0^-$$

as $n \rightarrow \infty$ at a faster rate than ε_n , which does not guarantee that $d_M(x, y) \leq 2\hat{r}_n(1 - \varepsilon_n)\arcsin(a)$ due to the $C'\varepsilon_n^2$ term.

From $\|x - y\| \geq r_0(1 - \varepsilon_n)$ and $\hat{r}_n \leq \text{diam}(M)$, we know that $r_0/\text{diam}(M) < a$, also $a < 1 + b$. Let us bound,

$$\begin{aligned} \arcsin(a) - \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a(1 - \varepsilon_n) + \delta_n/\hat{r}_n)a &\leq \\ \arcsin(a) - \frac{(1 + \varepsilon_n^2)}{\gamma_n} \arcsin'(a(1 - \varepsilon_n))a &=: f_n(a). \end{aligned}$$

Let us consider the function $g(t) = \arcsin(t) - \arcsin'(t)t$, which fulfills $g(0) = 0$ and $g'(t) < 0$ for all $0 < t < 1$. From $r_0/\text{diam}(M) < a$ it follows that $g(a) \leq g(r_0/\text{diam}(M)) < 0$. The functions f_n converges uniformly to g , on any closed interval containing a . This entails that $f_n(a) \leq g(r_0/\text{diam}(M))/2$, for n large enough, from where it follows (13).

□

6.2 Proof of Theorem 1

This proof is based on the following Lemma, which is a modified version of Main Theorem A in Bernstein et al (2000) for sets of positive reach. The proof follows the same ideas used to prove Main Theorem A, so we will provide a sketch.

Lemma 2. *Let $M \subset \mathbb{R}^d$ be compact, such that $\text{reach}(M) = r_0 > 0$. Let $\mathcal{X}_n \subset M$ a set of n points. We are given a graph G on \mathcal{X}_n and positive real numbers λ_1, λ_2 . We also refer to positive real numbers ε_{\min} , ε_{\max} , and δ . Suppose that,*

1. *The graph G contains all edges xy of length $\|x - y\| \leq \varepsilon_{\min}$.*
2. *All edges of G have length $\|x - y\| \leq \varepsilon_{\max}$.*
3. *The set \mathcal{X}_n fulfills $d_H(\mathcal{X}_n, M) \leq \delta$.*
4. *M is geodesically convex.*

Then provided that

5. $\varepsilon_{\max} < 2r_0$,
6. $\varepsilon_{\max} \leq (2/\pi)r_0\sqrt{24}\lambda_1$,
7. $\delta \leq \lambda_2\varepsilon_{\min}/4$,

it follows that inequalities

$$(1 - \lambda_1)d_M(x, y) \leq d_G(x, y) \leq (1 + \lambda_2)d_M(x, y) \quad (14)$$

are valid for all $x, y \in M$.

Proof. The second inequality in (14) is proven as the second inequality in Main Theorem A in Bernstein et al (2000), which is based on Theorem 2 in Bernstein et al (2000). The proof of Theorem 2 in Bernstein et al (2000) does not require smoothness assumptions on M . To prove the first inequality, let us consider a path $x_0x_1 \dots x_p$ on G connecting x and y (i.e., $x_0 = x$ and $x_p = y$.) From conditions 2, 5 and 7 $\|x_i - x_{i+1}\| \leq 2r_0$ and $\|x_i - x_{i+1}\| \leq (2/\pi)r_0\sqrt{24}\lambda_1$. From (1), we get that $2r_0 \sin(d_M(x, y)/2r_0) \leq \|x - y\|$, for all x, y such that $\|x - y\| \leq 2r_0$ (this is Lemma 3 in Bernstein et al (2000)). Thus

$$d_M(x_i, x_{i+1}) \leq (\pi/2)\|x_i - x_{i+1}\| \leq r_0\sqrt{24}\lambda_1.$$

We can then apply Corollary 4 in Bernstein et al (2000) with $\lambda = \lambda_1$ (its proof depends only on Lemma 3). Then, $d_M(x_i, x_{i+1}) \leq (1 - \lambda_1^{-1})\|x_i - x_{i+1}\|$. Lastly,

$$d_M(x, y) \leq (1 - \lambda)^{-1}\|x_0 - x_1\| + \dots + (1 - \lambda)^{-1}\|x_{p-1} - x_p\| = (1 - \lambda_1)^{-1}d_G(x, y).$$

□

Proof of Theorem 1

Proof. Let us verify condition 1 to 6 in Lemma 2. By construction of G_n , condition 1 and 2 are guaranteed with $\varepsilon_{\max} = \varepsilon_{\min} = \varepsilon_n$. Also $\delta_n = d_H(\mathcal{X}_n, M)$, guarantee condition 3. Conditions 4 and 5 are fulfilled by hypotheses. Lastly, let us define

$$\lambda_1 := \frac{1}{24} \left(\frac{\pi\varepsilon_n}{2r_0} \right)^2 \quad \text{and} \quad \lambda_2 := 4\tau_n,$$

then, conditions 6 and 7 are fulfilled, from where it follows (3). □

Acknowledgements

This research has been partially supported by grant FCE-1-2019-1-156054, ANII, Uruguay. The constructive comments and criticisms from two anonymous referees are gratefully acknowledged.

References

- Aamari, E., Aaron, C., and Levrard, C. (2021) Minimax Boundary Estimation and Estimation with Boundary. <https://arxiv.org/abs/2108.03135>
- Aamari E., Kim, J., Chazal, F., Michel, B., Rinaldo, A., and Wasserman, L. (2019). Estimating the reach of a manifold. *Electronic Journal of Statistics*, 13(1), 1359-1399.
- Aamari, E. and Levrard, C. (2019). Nonasymptotic rates for manifold, tangent space and curvature estimation. *The Annals of Statistics*, 47(1), 177-204.
- Aamari, E. and Levrard, C. (2018). Stability and minimax optimality of tangential Delaunay complexes for manifold reconstruction. *Discrete & Computational Geometry*, 59(4), 923–971.
- Aamari, E, Berenfeld, C. and Levrard, C. (2022). Optimal reach estimation and metric Learning. <https://arxiv.org/abs/2207.06074>
- Aaron, C. and Cholaquidis, A. (2020). On boundary detection. *In Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 56(3), 2028–2050.
- Aaron, C., Cholaquidis, A., and Cuevas, A. (2017). Detection of low dimensionality and data denoising via set estimation techniques. *Electronic Journal of Statistics*, 11(2), 4596-4628.
- Ambrosio, L., Colesanti, A. and Villa, E. (2008). Outer Minkowski content for some classes of closed sets. *Mathematische Annalen*, 342, 727–748.
- Arias-Castro, E., Javanmard, A. and Pelletier, B. (2020). Perturbation bounds for procrustes, classical scaling, and trilateration, with applications to manifold learning. *Journal of Machine Learning Research*, 21, 1–37.
- Arias-Castro, E., Pateiro-López, B. and Rodríguez-Casal, A. (2019). Minimax estimation of the volume of a set under the rolling ball condition. *Journal of the American Statistical Association*, 114(527), 1162–1173.
- Berenfeld, C., Harvey, J., Hoffmann, M., and Shankar, K. (2022). Estimating the reach of a manifold via its convexity defect function. *Discrete and Computational Geometry*, 67(2), 403-438.
- Bernstein, M., De Silva, V., Langford, J. C., and Tenenbaum, J. B. (2000). Graph approximations to geodesics on embedded manifolds. *Technical report, Department of Psychology, Stanford University*, 961–968.
- Boissonnat, J. D., Lieutier, A., and Wintraecken, M. (2019). The reach, metric distortion, geodesic convexity and the variation of tangent spaces. *Journal of Applied and Computational Topology*, 3(1), 29–58.
- Cadre, B. (2006) Kernel estimation of density level sets. *Journal of Multivariate Analysis*, 97(4), 999-1023.

- Cholaquidis, A., Cuevas, A. and Fraiman, R. (2014) On Poincaré cone property. *The Annals of Statistics*, 42, 255–284.
- Cholaquidis, A. and Fraiman, R. and Lugosi, G. and Pateiro-López, B. (2016). Set estimation from reflected Brownian motion. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5), 1057–1078
- Cuevas, A. and Rodríguez-Casal, A.(2004) On boundary estimation. *Advances in Applied Probability*, 36, 340–354.
- Cuevas, A., Fraiman, R. and Pateiro-López, B. (2012). On statistical properties of sets fulfilling rolling-type conditions. *Advances in Applied Probability*, 44, 311–329.
- Cuevas, A., González-Manteiga, W., and Rodríguez-Casal, A.(2006). Plug-in estimation of general level sets. *Australian New Zealand Journal of Statistics*, 48(1), 7-19.
- Cuevas, A. and Fraiman, R. (1997). A plug-in approach to support estimation. *The Annals of Statistics*, 25, 2300-2312.
- Cuevas, A., Fraiman, R. and Györfi, L. (2013). Towards a universally consistent estimator of the Minkowski content. *ESAIM: Probability and Statistics*, 17, 359–369.
- Cuevas, A., Fraiman, R., and Rodríguez-Casal, A. (2007). A nonparametric approach to the estimation of lengths and surface areas. *The Annals of Statistics*, 35(3), 1031-1051.
- Devroye, L. and Wise, G. (1980) Detection of abnormal behaviour via nonparametric estimation of the support. *SIAM Journal on Applied Mathematics*, 3, 480–488.
- Devroye, L. (1983) On arbitrarily slow rates of global convergence in density estimation. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 62(4), 475–483.
- Donoho, D. L.(1988) One-Sided Inference about Functionals of a Density. *The Annals of Statistics* 16(4), 1390–1420.
- Federer, H. (1959) Curvature measures. *Transactions of the American Mathematical Society*, 93, 418–491.
- Federer, H. (1969). Geometric measure theory. *Springer*.
- Fefferman, C., Mitter, S. and Narayanan, H. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4), 983–1049.
- Fraiman, R. and Meloche, J. (1999). Counting bumps. *Annals of the Institute of Statistical Mathematics*, 51(3), 541–569.
- Fu, J.H.G (1989) Curvature measures and generalized Morse theory. *Journal of Differential Geometry*, 30(3), 619–642.
- Genovese, C.R., Perone-Pacífico, M., Verdinelli, I. and Wasserman, L. (2012a). The geometry of nonparametric filament estimation. *Journal of the American Statistical Association*, 107, 788–799.
- Genovese, C.R., Perone-Pacífico, M., Verdinelli, I. and Wasserman, L. (2012b). Minimax manifold estimation. *Journal of Machine Learning Research*, 13, 1263–1291.

- Horobeţ, Emil and Weinstein, Madeleine (2019). Offset hypersurfaces and persistent homology of algebraic varieties. *Computer Aided Geometric Design*, 74, 101767.
- Le Cam, L. and Schwartz, L (1960). A necessary and sufficient condition for the existence of consistent estimates. *Annals of Mathematical Statistics*, 31, 140–150.
- Niyogi, P., Smale, S., and Weinberger, S.(2008) Finding the homology of submanifolds with high confidence from random samples. *Discrete & Computational Geometry*, 39, 419–441.
- Pateiro-López, B. and Rodríguez-Casal, A. (2009) Surface area estimation under convexity type assumptions. *Journal of Nonparametric Statistics*, 21(6), 729–741
- Pateiro-López, B. (2008). Set estimation under convexity type restrictions. Universidade de Santiago de Compostela.
- Rataj, J. and Zajíček, L. On the structure of sets with positive reach. (2017) *Mathematische Nachrichten*, 290(11-12), 1806–1829.
- Rataj, J., and Zähle, M. (2001). Curvatures and currents for unions of sets with positive reach, II. *Annals of Global Analysis and Geometry*, 20(1), 1-21.
- Rodríguez Casal, A.R. (2007) Set estimation under convexity type assumption. *Annales de l'Institut Henri Poincaré*, 43, 763–774.
- Rodríguez-Casal, A., and Saavedra-Nieves, P. (2021). Spatial distribution of invasive species: an extent of occurrence approach. *TEST*, 31, 416–441.
- Rodríguez-Casal, A., and Saavedra-Nieves, P. (2021). A fully data-driven method for estimating the shape of a point cloud. *ESAIM*, 20, 332-348.
- Thäle, C. (2008). 50 years sets with positive reach. A survey. *Surveys in Mathematics and its Applications*, 3, 123–165.
- Walther, G. (1999). On a generalization of Blaschke's rolling theorem and the smoothing of surfaces, *Mathematical Methods in the Applied Sciences*, 22, 301–316.