

# Robust Persistence Diagrams using Reproducing Kernels

SIDDHARTH VISHWANATH<sup>†</sup>, KENJI FUKUMIZU<sup>‡\*</sup>,  
SATOSHI KURIKI<sup>‡\*</sup>, AND BHARATH SRIPERUMBUDUR<sup>†\*</sup>

<sup>†</sup>*Department of Statistics, The Pennsylvania State University, University Park, PA, USA*

<sup>‡</sup>*The Institute of Statistical Mathematics, Tokyo, Japan*

June 18, 2020

## Abstract

Persistent homology has become an important tool for extracting geometric and topological features from data, whose multi-scale features are summarized in a persistence diagram. From a statistical perspective, however, persistence diagrams are very sensitive to perturbations in the input space. In this work, we develop a framework for constructing robust persistence diagrams from superlevel filtrations of robust density estimators constructed using reproducing kernels. Using an analogue of the influence function on the space of persistence diagrams, we establish the proposed framework to be less sensitive to outliers. The robust persistence diagrams are shown to be consistent estimators in bottleneck distance, with the convergence rate controlled by the smoothness of the kernel—this in turn allows us to construct uniform confidence bands in the space of persistence diagrams. Finally, we demonstrate the superiority of the proposed approach on benchmark datasets.

## 1 Introduction

Given a set of points  $\mathbb{X}_n = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n\}$  observed from a probability distribution  $\mathbb{P}$  on an input space  $\mathcal{X} \subseteq \mathbb{R}^d$ , understanding the shape of  $\mathbb{X}_n$  sheds important insights on low-dimensional geometric and topological features which underlie  $\mathbb{P}$ , and this question has received increasing attention in the past few decades. To this end, Topological Data Analysis (TDA), with a special emphasis on persistent homology (Edelsbrunner et al., 2000; Zomorodian and Carlsson, 2005), has become a mainstay for extracting the shape information from data. In statistics and machine-learning, persistent homology has facilitated the development of novel methodology (e.g., Chazal et al. 2013; Chen et al. 2019; Brüel-Gabrielsson et al. 2018), which has been widely used in a variety of applications dealing with massive, unconventional forms of data (e.g., Gameiro et al. 2015; Bendich et al. 2016; Xu et al. 2019).

Informally speaking, persistent homology detects the presence of topological features across a range of resolutions by examining a nested sequence of spaces, typically referred to as a *filtration*. The filtration encodes the birth and death of topological features as the resolution varies, and is presented in the form of a concise representation—a persistence diagram or barcode. In the context of data-analysis, there are two different methods for obtaining filtrations. The first is computed from the pairwise Euclidean distances of  $\mathbb{X}_n$ , such as the Vietoris-Rips, Čech, and Alpha filtrations

---

\*Authors arranged alphabetically

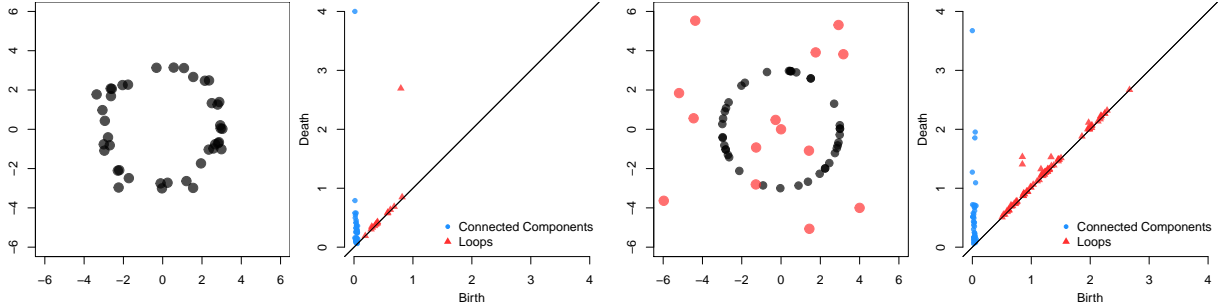


Figure 1: (Left)  $\mathbb{X}_n$  is sampled from a circle with small perturbations to each point. The persistence diagram detects the presence of the loop, as guaranteed by the stability of persistence diagrams Chazal et al. (2016); Cohen-Steiner et al. (2007). (Right)  $\mathbb{X}_n$  is sampled from a circle but with just a few outliers. The resulting persistence diagram changes dramatically — the persistence of the main loop plummets, and other spurious loops appear, as elaborated in Section 2.

Edelsbrunner et al. (2000). The second approach is based on choosing a function on  $\mathcal{X}$  that reflects the density of  $\mathbb{P}$  (or its approximation based on  $\mathbb{X}_n$ ), and, then, constructing a filtration. While the two approaches explore the topological features governing  $\mathbb{P}$  in different ways, in essence, they generate equivalent insights.

Despite obvious advantages, the adoption of persistent homology in mainstream statistical methodology is still limited. An important limitation among others, in the statistical context, is that the resulting persistent homology is highly sensitive to outliers. While the stability results of Chazal et al. (2016); Cohen-Steiner et al. (2007) guarantee that small perturbations on all of  $\mathbb{X}_n$  induce only small changes in the resulting persistence diagrams, a more pathological issue arises when a small fraction of  $\mathbb{X}_n$  is subject to very large perturbations. Figure 1 illustrates how inference from persistence diagrams can change dramatically when  $\mathbb{X}_n$  is contaminated with only a few outliers. Another challenge is the mathematical difficulty in performing sensitivity analysis in a formal statistical context. Since the space of persistence diagrams have an unusual mathematical structure, it falls victim to issues such as non-uniqueness of Fréchet means and bounded curvature of geodesics (Mileyko et al., 2011; Turner et al., 2014; Divol and Lacombe, 2019). With this background, the *central objective* of this paper is to develop outlier robust persistence diagrams, develop a framework for examining the sensitivity of the resulting persistence diagrams to noise, and establish statistical convergence guarantees. To the best of our knowledge, not much work has been carried out in this direction, except for Bendich et al. (2011) where robust persistence diagrams are constructed from Vietoris-Rips or Čech filtrations on  $\mathbb{X}_n$  by replacing the Euclidean distance with diffusion distance. However, no sensitivity analysis of the resultant diagrams are carried out in (Bendich et al., 2011) to demonstrate their robustness.

**Contributions.** The main contributions of this work are threefold. 1) We propose robust persistence diagrams constructed from filtrations induced by an RKHS-based robust KDE (kernel density estimator) Kim and Scott (2012) of the underlying density function of  $\mathbb{P}$  (Section 3). While this idea of inducing filtrations by an appropriate function—(Fasy et al., 2014; Chazal et al., 2017; Phillips et al., 2015) use KDE, distance-to-measure (DTM) and kernel distance (KDist), respectively—has already been explored, we show the corresponding persistence diagrams to be less robust compared to our proposal. 2) In Section 4.1, we generalize the notions of *influence function* and *gross error sensitivity*—which are usually defined for normed spaces—to the space of persistence diagrams, which lack the vector space structure. Using these generalized notions, we investigate the sensitivity of persistence diagrams constructed from filtrations induced by different functions (e.g., KDE, robust KDE, DTM) and demonstrate the robustness of the proposed method, both theoretically (Remark 4.1)

and numerically (Section 5). 3) We establish the statistical consistency of the proposed robust persistence diagrams and provide uniform confidence bands by deriving exponential concentration bounds for the uniform deviation of the robust KDE (Section 4.2).

**Definitions and Notations.** For a metric space  $(\mathcal{X}, \varrho)$ , the ball of radius  $r$  centered at  $\mathbf{x} \in \mathcal{X}$  is denoted by  $B_{\mathcal{X}}(\mathbf{x}, r)$ .  $L_p(\mathcal{X}, \mu)$  is the Banach space of functions of  $p^{\text{th}}$ -power  $\mu$ -integrable functions with norm  $\|\cdot\|_p$ , where  $\mu$  is a Borel measure defined on  $\mathcal{X}$ .  $\mathcal{P}(\mathbb{R}^d)$  is the set of all Borel probability measures on  $\mathbb{R}^d$ , and  $\mathcal{M}(\mathbb{R}^d)$  denotes the set of probability measures on  $\mathbb{R}^d$  with compact support and *tame*, bounded density function. For bandwidth  $\sigma > 0$ ,  $\mathcal{H}_\sigma$  denotes a reproducing kernel Hilbert space (RKHS) with  $K_\sigma : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  as its reproducing kernel. We assume that  $K_\sigma$  is radial, i.e.,  $K_\sigma(\mathbf{x}, \mathbf{y}) = \sigma^{-d} \psi(\|\mathbf{x} - \mathbf{y}\|_2 / \sigma)$  with  $\psi(\|\cdot\|_2)$  being a probability density function on  $\mathbb{R}^d$ , where  $\|\mathbf{x}\|_2^2 = \sum_{i=1}^d x_i^2$  for  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ . We denote  $\|K_\sigma\|_\infty \doteq \sup_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^d} K_\sigma(\mathbf{x}, \mathbf{y}) = \sigma^{-d} \psi(0)$ . For  $\mathbb{P} \in \mathcal{P}(\mathbb{R}^d)$ ,  $\mu_\mathbb{P} \doteq \int K_\sigma(\cdot, \mathbf{y}) d\mathbb{P}(\mathbf{y}) \in \mathcal{H}_\sigma$  is called the mean embedding of  $\mathbb{P}$ , and  $\mathcal{D}_\sigma \doteq \{\mu_\mathbb{P} : \mathbb{P} \in \mathcal{P}(\mathbb{R}^d)\}$  is the space of all mean embeddings.  $\delta_{\mathbf{x}}$  denotes a Dirac measure at  $\mathbf{x}$ .

## 2 Persistent Homology: Preliminaries

We present the necessary background on persistent homology for completeness. See Chazal and Michel (2017); Wasserman (2018) for a comprehensive introduction.

**Persistent Homology.** Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  be a non-negative function on the metric space  $(\mathcal{X}, d)$ . At level  $r > 0$ , the *sublevel* set  $\mathcal{X}_r = \phi^{-1}([0, r]) = \{\mathbf{x} \in \mathcal{X} : \phi(\mathbf{x}) \leq r\}$  encodes the topological information in  $\mathcal{X}$ . For  $0 \leq r < s \leq \infty$ , the sublevel sets are nested, i.e.,  $\mathcal{X}_r \subseteq \mathcal{X}_s$ . The sequence  $\{\mathcal{X}_r\}_{0 \leq r \leq \infty}$  is a nested sequence of topological spaces, called a *filtration*, denoted by  $\text{Sub}(\phi)$ , and  $\phi$  is called the *filter function*. As the level  $r$  varies, the evolution of the topology is captured in the filtration. Roughly speaking, new cycles (i.e., connected components, loops, voids and higher order analogues) can appear or existing cycles can merge. A new  $k$ -dimensional feature is said to be born at  $b \in \mathbb{R}$  when a nontrivial  $k$ -cycle appears in  $\mathcal{X}_b$ . The same  $k$ -cycle dies at level  $d > b$  when it disappears in all  $\mathcal{X}_{d+\epsilon}$  for  $\epsilon > 0$ . Persistent homology,  $PH_*(\phi)$ , is an algebraic module which tracks the *persistence pairs*  $(b, d)$  of births  $b$  and deaths  $d$  across the entire filtration  $\text{Sub}(\phi)$ . Mutatis mutandis, a similar notion holds for superlevel sets  $\mathcal{X}^r = \phi^{-1}([r, \infty))$ , inducing the filtration  $\text{Sup}(\phi)$ . For  $r < s$ , the inclusion  $\mathcal{X}^r \supseteq \mathcal{X}^s$  is reversed and a cycle born at  $b$  dies at a level  $d < b$ , resulting in the persistence pair  $(d, b)$  instead. Figure 2 shows three connected components in the superlevel set for  $r = 8$ . The components were born as  $r$  swept through the blue points, and die when  $r$  approaches the red points. We refer the reader to Appendix B for more details. Figure 2 is described in more detail in Figure 9.

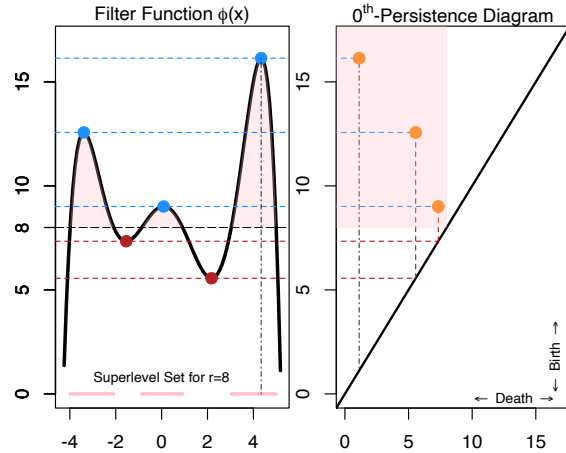


Figure 2:  $\text{Dgm}(\text{Sup}(\phi))$  for  $\phi : \mathbb{R} \rightarrow \mathbb{R}$

**Persistence Diagrams.** By collecting all persistence pairs, the persistent homology is concisely represented as a persistence diagram  $\text{Dgm}(\text{Sub}(\phi)) \doteq \{(b, d) \in \mathbb{R}^2 : 0 \leq b < d \leq \infty\}$ . A similar definition carries over to  $\text{Dgm}(\text{Sup}(\phi))$ , using  $(d, b)$  instead. See Figure 2 for an illustration. When

the context is clear, we drop the reference to the filtration and simply write  $\text{Dgm}(\phi)$ . The  $k^{\text{th}}$  persistence diagram is the subset of  $\text{Dgm}(\phi)$  corresponding to the  $k$ -dimensional features. The space of persistence diagrams is the locally-finite multiset of points on  $\Omega = \{(x, y) : 0 \leq x < y \leq \infty\}$ , endowed with the family of  $p$ -Wasserstein metrics  $W_p$ , for  $1 \leq p \leq \infty$ . We refer the reader to Edelsbrunner and Harer (2010); Divol and Lacombe (2019) for a thorough introduction.  $W_\infty$  is commonly referred to as the *bottleneck distance*.

**Definition 2.1.** *Given two persistence diagrams  $D_1$  and  $D_2$ , the bottleneck distance is given by*

$$W_\infty(D_1, D_2) = \inf_{\gamma \in \Gamma} \sup_{p \in D_1 \cup \Delta} \|p - \gamma(p)\|_\infty,$$

where  $\Gamma = \{\gamma : D_1 \cup \Delta \rightarrow D_2 \cup \Delta\}$  is the set of all bijections from  $D_1$  to  $D_2$ , including the diagonal  $\Delta = \{(x, y) \in \mathbb{R}^2 : 0 \leq x = y \leq \infty\}$  with infinite multiplicity.

An assumption we make at the outset is that the filter function  $f$  is *tame*. Tameness is a metric regularity condition which ensures that the number of points on the persistence diagrams are finite, and, in addition, the number of nontrivial cycles which share identical persistence pairings are also finite. Tame functions satisfy the celebrated stability property w.r.t. the bottleneck distance.

**Proposition 2.1** (Stability of Persistence Diagrams Cohen-Steiner et al., 2007; Chazal et al., 2016). *Given two tame functions  $f, g : \mathcal{X} \rightarrow \mathbb{R}$ ,*

$$W_\infty(\text{Dgm}(f), \text{Dgm}(g)) \leq \|f - g\|_\infty.$$

The space of persistence diagrams is, in general, challenging to work with. However, the stability property provides a handle on the persistence space through the function space of filter functions.

### 3 Robust Persistence Diagrams

Given  $\mathbb{X}_n = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n\} \subseteq \mathbb{R}^d$  drawn iid from a probability distribution  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$  with density  $f$ , the corresponding persistence diagram can be obtained by considering a filter function  $\phi_n : \mathbb{R}^d \rightarrow \mathbb{R}$ , constructed from  $\mathbb{X}_n$  as an approximation to its population analogue,  $\phi_{\mathbb{P}} : \mathbb{R}^d \rightarrow \mathbb{R}$ , that carries the topological information of  $\mathbb{P}$ . Commonly used  $\phi_{\mathbb{P}}$  include the (i) kernelized density,  $f_\sigma$ , (ii) Kernel Distance (KDist),  $d_{\mathbb{P}}^{K_\sigma}$ , and (iii) distance-to-measure (DTM),  $d_{\mathbb{P},m}$ , which are defined as:

$$f_\sigma(\mathbf{x}) \doteq \int_{\mathcal{X}} K_\sigma(\mathbf{x}, \mathbf{y}) d\mathbb{P}(\mathbf{y}) ; \quad d_{\mathbb{P}}^{K_\sigma} \doteq \|\mu_{\delta_{\mathbf{x}}} - \mu_{\mathbb{P}}\|_{\mathcal{H}_\sigma} ; \quad d_{\mathbb{P},m}(\mathbf{x}) \doteq \sqrt{\frac{1}{m} \int_0^m F_{\mathbf{x}}^{-1}(u) du},$$

where  $F_{\mathbf{x}}(t) = \mathbb{P}(\|\mathbf{X} - \mathbf{x}\|_2 \leq t)$  and  $\sigma, m > 0$ . For these  $\phi_{\mathbb{P}}$ , the corresponding empirical analogues,  $\phi_n$ , are constructed by replacing  $\mathbb{P}$  with the empirical measure,  $\mathbb{P}_n \doteq \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{X}_i}$ . For example, the empirical analogue of  $f_\sigma$  is the familiar kernel density estimator (KDE),  $f_\sigma^n = \frac{1}{n} \sum_{i=1}^n K_\sigma(\cdot, \mathbf{X}_i)$ . While KDE and KDist capture the topological information of  $\text{supp}(\mathbb{P})$  by approximating the density  $f$  (the sublevel sets of KDist are simply rescaled versions of the superlevel sets of KDE (Phillips et al., 2015; Chazal et al., 2017)), DTM, on the other hand, approximates the distance function to  $\text{supp}(\mathbb{P})$ .

Since  $\phi_n$  is based on  $\mathbb{P}_n$ , it is sensitive to outliers in  $\mathbb{X}_n$ , which, in turn affect the persistence diagrams (as illustrated in Figure 1). To this end, in this paper, we propose *robust persistence diagrams* constructed using superlevel filtrations of a robust density estimator of  $f$ , i.e., the filter function,  $\phi_n$

is chosen to be a robust density estimator of  $f$ . Specifically, we use the robust KDE,  $f_{\rho,\sigma}^n$ , introduced by Kim and Scott (2012) as the filter function, which is defined as a solution to the following M-estimation problem:

$$f_{\rho,\sigma}^n \doteq \arg \inf_{g \in \mathcal{G}} \int_{\mathcal{X}} \rho(\|\Phi_{\sigma}(\mathbf{y}) - g\|_{\mathcal{H}_{\sigma}}) d\mathbb{P}_n(\mathbf{y}), \quad (1)$$

where  $\rho : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  is a robust loss function,  $\Phi_{\sigma}(\mathbf{x}) = K_{\sigma}(\cdot, \mathbf{x}) \in \mathcal{H}_{\sigma}$  is the feature map associated with  $K_{\sigma}$  for a fixed  $\sigma > 0$ , and  $\mathcal{G} = \mathcal{H}_{\sigma} \cap \mathcal{D}_{\sigma} = \mathcal{D}_{\sigma}$  is the hypothesis class. Observe that when  $\rho(z) = \frac{1}{2}z^2$ , the unique solution to Eq. (1) is given by the KDE,  $f_{\sigma}^n$ . Therefore, a robust KDE is obtained by replacing the square loss with a *robust loss*, which satisfies the following assumptions. These assumptions, which are similar to those of Kim and Scott (2012); Vandermeulen and Scott (2013) guarantee the existence and uniqueness of  $f_{\rho,\sigma}^n$  (if  $\rho$  is convex; Kim and Scott 2012), and are satisfied by most robust loss functions, including the Huber loss,  $\rho(z) = \frac{1}{2}z^2 \mathbb{1}\{z \leq 1\} + (z - \frac{1}{2}) \mathbb{1}\{z > 1\}$  and the Charbonnier loss,  $\rho(z) = \sqrt{1 + z^2} - 1$ .

(A1)  $\rho$  is strictly-increasing and  $M$ -Lipschitz, with  $\rho(0) = 0$ .

(A2)  $\rho'(x)$  is continuous and bounded with  $\rho'(0) = 0$ .

(A3)  $\varphi(x) = \rho'(x)/x$  is bounded,  $L$ -Lipschitz and continuous, with  $\varphi(0) < \infty$ .

(A4)  $\rho''$  exists, with  $\rho''$  and  $\varphi$  nonincreasing.

Unlike for squared loss, the solution  $f_{\rho,\sigma}^n$  cannot be obtained in a closed form. However, it can be shown to be the fixed point of an iterative procedure, referred to as KIRWLS algorithm (Kim and Scott, 2012). The KIRWLS algorithm starts with initial weights  $\{w_i^{(0)}\}_{i=1}^n$  such that  $\sum_{i=1}^n w_i^{(0)} = 1$ , and generates the iterative sequence of estimators  $\{f_{\rho,\sigma}^{(k)}\}_{k \in \mathbb{N}}$  as

$$f_{\rho,\sigma}^{(k)} = \sum_{i=1}^n w_i^{(k-1)} K_{\sigma}(\cdot, \mathbf{X}_i) \quad ; \quad w_i^{(k)} = \frac{\varphi(\|\Phi_{\sigma}(\mathbf{X}_i) - f_{\rho,\sigma}^{(k)}\|_{\mathcal{H}_{\sigma}})}{\sum_{j=1}^n \varphi(\|\Phi_{\sigma}(\mathbf{X}_j) - f_{\rho,\sigma}^{(k)}\|_{\mathcal{H}_{\sigma}})}.$$

Intuitively, note that if  $\mathbf{X}_i$  is an outlier, then the corresponding weight  $w_i$  is small (since  $\varphi$  is nonincreasing) and therefore less weightage is given to the contribution of  $\mathbf{X}_i$  in the density estimator. Hence, the weights serve as a measure of *inlyingness*—smaller (*resp.* larger) the weights, lesser (*resp.* more) inlying are the points. When  $\mathbb{P}_n$  is replaced by  $\mathbb{P}$ , the solution of Eq. (1) is its population analogue,  $f_{\rho,\sigma}$ . Although  $f_{\rho,\sigma}$  does not admit a closed form solution, it can be shown (Kim and Scott, 2012) that there exists a non-negative real-valued function  $w_{\sigma}$  satisfying  $\int_{\mathbb{R}^d} w_{\sigma}(\mathbf{x}) d\mathbb{P}(\mathbf{x}) = 1$  such that

$$f_{\rho,\sigma} = \int_{\mathbb{R}^d} K_{\sigma}(\cdot, \mathbf{x}) w_{\sigma}(\mathbf{x}) d\mathbb{P}(\mathbf{x}) = \int_{\mathbb{R}^d} \frac{\varphi(\|\Phi_{\sigma}(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_{\sigma}})}{\int_{\mathbb{R}^d} \varphi(\|\Phi_{\sigma}(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_{\sigma}}) d\mathbb{P}(\mathbf{y})} K_{\sigma}(\cdot, \mathbf{x}) d\mathbb{P}(\mathbf{x}), \quad (2)$$

where  $w_{\sigma}$  acts as a population analogue of the weights in KIRWLS algorithm.

To summarize our proposal, the fixed point of the KIRWLS algorithm, which yields the robust density estimator  $f_{\rho,\sigma}^n$ , is used as the filter function to obtain a robust persistence diagram of  $\mathbb{X}_n$ . On the computational front, note that  $f_{\rho,\sigma}^n$  is computationally more complex than the KDE,  $f_{\sigma}^n$ , requiring  $O(n\ell)$  computations compared to  $O(n)$  of the latter, with  $\ell$  being the number of iterations required to reach the fixed point of KIRWLS. However, once these filter functions are computed, the corresponding persistence diagrams have similar computational complexity as both require computing superlevel sets, which, in turn, require function evaluations that scale as  $O(n)$  for both  $f_{\rho,\sigma}^n$  and  $f_{\sigma}^n$ .

## 4 Theoretical Analysis of Robust Persistence Diagrams

In this section, we investigate the theoretical properties of the proposed robust persistence diagrams. First, in Section 4.1, we examine the sensitivity of persistence diagrams to outlying perturbations through the notion of *metric derivative* and compare the effect of different filter functions. Next, in Section 4.2, we establish consistency and convergence rates for the robust persistence diagram to its population analogue. These results allow to construct uniform confidence bands for the robust persistence diagram. The proofs of the results are provided in Section 6.

### 4.1 A measure of sensitivity of persistence diagrams to outliers

The influence function and gross error sensitivity are arguably the most popular tools in robust statistics for diagnosing the sensitivity of an estimator to a single adversarial contamination Hampel et al. (2011); Huber (2004). Given a statistical functional  $T : \mathcal{P}(\mathcal{X}) \rightarrow (V, \|\cdot\|_V)$ , which takes an input probability measure  $\mathbb{P} \in \mathcal{P}(\mathcal{X})$  on the input space  $\mathcal{X}$  and produces a statistic  $\mathbb{P} \mapsto T(\mathbb{P})$  in some normed space  $(V, \|\cdot\|_V)$ , the *influence function* of  $\mathbf{x} \in \mathcal{X}$  at  $\mathbb{P}$  is given by the Gâteaux derivative of  $T$  at  $\mathbb{P}$  restricted to the space of signed Borel measures with zero expectation:

$$IF(T; \mathbb{P}, \mathbf{x}) \doteq \frac{\partial}{\partial \epsilon} T\left((1 - \epsilon)\mathbb{P} + \epsilon\delta_{\mathbf{x}}\right)\Big|_{\epsilon=0} = \lim_{\epsilon \rightarrow 0} \frac{T((1 - \epsilon)\mathbb{P} + \epsilon\delta_{\mathbf{x}}) - T(\mathbb{P})}{\epsilon},$$

and the *gross error sensitivity* at  $\mathbb{P}$  is given by  $\Gamma(T; \mathbb{P}) \doteq \sup_{\mathbf{x} \in \mathcal{X}} \|IF(T; \mathbb{P}, \mathbf{x})\|_V$ . However, a persistence diagram (which is a statistical functional) does not take values in a normed space, and, therefore, the notion of influence function has to be generalized to metric spaces through the concept of a metric derivative: Given a complete metric space  $(X, d_X)$  and a curve  $s : [0, 1] \rightarrow X$ , the *metric derivative* at  $\epsilon = 0$  is given by

$$|s'| (0) \doteq \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} d_X(s(0), s(\epsilon)).$$

Using this generalization, we have the following definition, which allows to examine the influence an outlier has on the persistence diagram obtained from a filtration.

**Definition 4.1.** *Given a probability measure  $\mathbb{P} \in \mathcal{P}(\mathbb{R}^d)$  and a filter function  $\phi_{\mathbb{P}}$  depending on  $\mathbb{P}$ , the persistence influence of a perturbation  $\mathbf{x} \in \mathbb{R}^d$  on  $\text{Dgm}(\phi_{\mathbb{P}})$  is defined as*

$$\Psi(\phi_{\mathbb{P}}; \mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} W_{\infty}(\text{Dgm}(\phi_{\mathbb{P}^{\epsilon}}), \text{Dgm}(\phi_{\mathbb{P}})),$$

where  $\mathbb{P}^{\epsilon} \doteq (1 - \epsilon)\mathbb{P} + \epsilon\delta_{\mathbf{x}}$ , and the gross-influence is defined as  $\Gamma(\phi_{\mathbb{P}}) = \sup_{\mathbf{x} \in \mathbb{R}^d} \Psi(\phi_{\mathbb{P}}; \mathbf{x})$ .

The following result (proved in Section 6.1) bounds the persistence influence for the persistence diagram induced by the filter function  $f_{\rho, \sigma}$ , which is the population analogue of robust KDE.

**Theorem 4.1.** *For a loss  $\rho$  satisfying (A1)–(A3), and  $\sigma > 0$ , the persistence influence of  $\mathbf{x} \in \mathbb{R}^d$  on  $\text{Dgm}(f_{\rho, \sigma})$  satisfies*

$$\Psi(f_{\rho, \sigma}; \mathbf{x}) \leq \|K_{\sigma}\|_{\infty}^{\frac{1}{2}} \rho' \left( \|\Phi_{\sigma}(\mathbf{x}) - f_{\rho, \sigma}\|_{\mathcal{H}_{\sigma}} \right) \left( \int_{\mathbb{R}^d} \zeta \left( \|\Phi_{\sigma}(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_{\sigma}} \right) d\mathbb{P}(\mathbf{y}) \right)^{-1}, \quad (3)$$

where  $\zeta(z) = \varphi(z) - z\varphi'(z)$ .

**Remark 4.1.** We make the following observations from Theorem 4.1.

(i) Choosing  $\rho(z) = \frac{1}{2}z^2$  and noting that  $\varphi(z) = \rho''(z) = 1$ , Eq. (3) yields a bound for the persistence influence of the KDE as

$$\Psi(f_\sigma; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \|\Phi_\sigma(\mathbf{x}) - f_\sigma\|_{\mathcal{H}_\sigma}.$$

On the other hand, for robust loss functions, the term involving  $\rho'$  is bounded because of (A2), making them less sensitive to very large perturbations. In fact, for nonincreasing  $\varphi$ , it can be shown (see Section 6.2) that

$$\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} w_\sigma(\mathbf{x}) \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma},$$

where, in contrast to KDE, the measure of inlyingness,  $w_\sigma$ , weighs down extreme outliers.

(ii) For the generalized Charbonnier loss (a robust loss function), given by  $\rho(z) = (1 + z^2)^{\alpha/2} - 1$  for  $1 \leq \alpha \leq 2$ , the persistence influence satisfies

$$\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \left(1 + \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}^2\right)^{\frac{\alpha-1}{2}} \left(1 + \int_{\mathbb{R}^d} \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}^2 d\mathbb{P}(\mathbf{y})\right)^{\frac{1-\alpha}{2}}.$$

Note that for  $\alpha = 1$ , the bound on the persistence influence  $\Psi(f_{\rho,\sigma}; \mathbf{x})$  does not depend on how extreme the outlier  $\mathbf{x}$  is. Similarly, for the Cauchy loss, given by  $\rho(z) = \log(1 + z^2)$ , we have

$$\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \left(1 + \int_{\mathbb{R}^d} \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}^2 d\mathbb{P}(\mathbf{y})\right).$$

This shows that for large perturbations, the gross error sensitivity for the Cauchy and Charbonnier losses are far more stable than that of KDE. This behavior is also empirically illustrated in Figure 3. The experiment is detailed in Section 5.2.

(iii) For the DTM function, it can be shown that

$$\Psi(d_{\mathbb{P},m}; \mathbf{x}) \leq \frac{2}{\sqrt{m}} \sup \left\{ \left| f(\mathbf{x}) - \int_{\mathbb{R}^d} f(\mathbf{y}) d\mathbb{P}(\mathbf{y}) \right| : \|\nabla f\|_{L_2(\mathbb{P})} \leq 1 \right\}. \quad (4)$$

While  $d_{\mathbb{P},m}$  cannot be compared to both  $f_\sigma$  and  $f_{\rho,\sigma}$ , as it captures topological information at a different scale, determined by  $m$ , we point out that when  $\text{supp}(\mathbb{P})$  is compact,  $\Psi(d_{\mathbb{P},m}; \mathbf{x})$  is not guaranteed to be bounded, unlike in  $\Psi(f_{\rho,\sigma}; \mathbf{x})$ . We refer the reader to Section 6.2 for more details.

It follows from Remark 4.1 that as  $\sigma \rightarrow 0$ , the persistence influence of both the KDE and robust KDE behave as  $O(\sigma^{-d})$ , showing that the robustness of robust persistence diagrams manifests only in cases where  $\sigma > 0$ . However, robustness alone has no bearing if the robust persistence diagram and the persistence diagram from the KDE are fundamentally different, i.e., they estimate different quantities as  $\sigma \rightarrow 0$ . The following result (proved in Section 6.3) shows that as  $\sigma \rightarrow 0$ ,  $\text{Dgm}(f_{\rho,\sigma})$  recovers the same information as that in  $\text{Dgm}(f_\sigma)$ , which is same as  $\text{Dgm}(f)$ , where  $f$  is the density of  $\mathbb{P}$ .

**Theorem 4.2.** For a strictly-convex loss  $\rho$  satisfying (A1)–(A4), and  $\sigma > 0$ , suppose  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$  with density  $f$ , and  $f_{\rho,\sigma}$  is the robust KDE. Then

$$W_\infty(\text{Dgm}(f_{\rho,\sigma}), \text{Dgm}(f)) \rightarrow 0 \quad \text{as } \sigma \rightarrow 0.$$



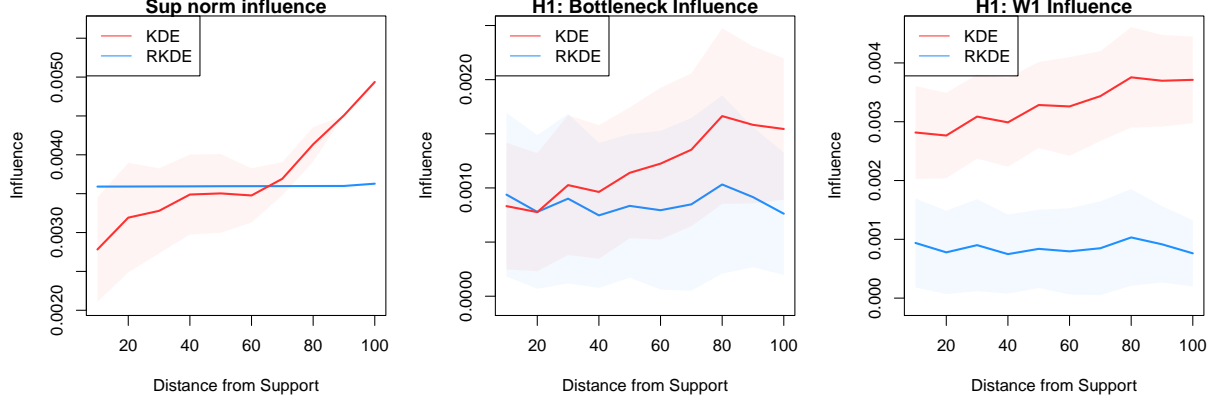


Figure 3: Points  $\mathbb{X}_n$  are sampled from  $\mathbb{P}$  with nontrivial 1<sup>st</sup>-order homological features and outliers  $\mathbb{Y}_m$  are added at a distance  $r$  from the support of  $\mathbb{P}$ . (Left) The average  $L_\infty$  distance between the density estimators computed using  $\mathbb{X}_n$  and  $\mathbb{X}_n \cup \mathbb{Y}_m$  as  $r$  increases. (Center) The average  $W_\infty$  distance between the corresponding persistence diagrams for the 1<sup>st</sup>-order homological features. (Right) The  $W_1$  distance (defined in Eq. B.1 in Appendix B) between the same persistence diagrams. The results show that the outliers  $\mathbb{Y}_m$  have little influence on the persistence diagrams from the robust KDEs. In contrast, as the outliers become more extreme (i.e.,  $r$  increases) their influence on the persistence diagrams from the KDE becomes more prominent.

Suppose  $\mathbb{P} = (1 - \pi)\mathbb{P}_0 + \pi\mathbb{Q}$ , where  $\mathbb{P}_0$  corresponds to the true signal which we are interested in studying, and  $\mathbb{Q}$  manifests as some ambient noise with  $0 < \pi < \frac{1}{2}$ . In light of Theorem 4.2, by letting  $\sigma \rightarrow 0$ , along with the topological features of  $\mathbb{P}_0$ , we are also capturing the topological features of  $\mathbb{Q}$ , which may obfuscate any statistical inference made using the persistence diagrams. In a manner, choosing  $\sigma > 0$  suppresses the noise in the resulting persistence diagrams, thereby making them more stable. On a similar note, the authors in Fasy et al. (2014) state that for a suitable bandwidth  $\sigma > 0$ , the level sets of  $f_\sigma$  carry the same topological information as  $\text{supp}(\mathbb{P})$ , despite the fact that some subtle details in  $f$  may be omitted. In what follows, we consider the setting where robust persistence diagrams are constructed for a fixed  $\sigma > 0$ .

## 4.2 Statistical properties of robust persistence diagrams from samples

Suppose  $\text{Dgm}(f_{\rho,\sigma}^n)$  is the robust persistence diagram obtained from the robust KDE on a sample  $\mathbb{X}_n$  and  $\text{Dgm}(f_{\rho,\sigma})$  is its population analogue obtained from  $f_{\rho,\sigma}$ . The following result (proved in Section 6.4) establishes the consistency of  $\text{Dgm}(f_{\rho,\sigma}^n)$  in the  $W_\infty$  metric.

**Theorem 4.3.** *For convex loss  $\rho$  satisfying (A1)–(A4), and fixed  $\sigma > 0$ , suppose  $\mathbb{X}_n$  is observed iid from a distribution  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$  with density  $f$ . Then*

$$W_\infty(\text{Dgm}(f_{\rho,\sigma}^n), \text{Dgm}(f_{\rho,\sigma})) \xrightarrow{p} 0 \text{ as } n \rightarrow \infty.$$

We present the convergence rate of the above convergence in Theorem 4.4, which depends on the smoothness of  $\mathcal{H}_\sigma$ . In a similar spirit to Fasy et al. (2014), this result paves the way for constructing uniform confidence bands. Before we present the result, we first introduce the notion of *entropy numbers* associated with an RKHS.



**Definition 4.2** (Entropy Number). *Given a metric space  $(T, d)$  the  $n^{\text{th}}$  entropy number is defined as*

$$e_n(T, d) \doteq \inf \left\{ \epsilon > 0 : \exists \{t_1, t_2, \dots, t_{2^{n-1}}\} \subset T \text{ such that } T \subset \bigcup_{i=1}^{2^{n-1}} B_d(t_i, \epsilon) \right\}.$$

*Further, if  $(V, \|\cdot\|_V)$  and  $(W, \|\cdot\|_W)$  are two normed spaces and  $L : V \rightarrow W$  is a bounded, linear operator, then  $e_n(L) = e_n(L : V \rightarrow W) \doteq e_n(L(B_V), \|\cdot\|_W)$ , where  $B_V$  is a unit ball in  $V$ .*

Loosely speaking, entropy numbers are related to the eigenvalues of the integral operator associated with the kernel  $K_\sigma$ , and measure the capacity of the RKHS in approximating functions in  $L_2(\mathbb{R}^d)$ . In our context, the entropy numbers will provide useful bounds on the covering numbers of sets in the hypothesis class  $\mathcal{G}$ . We refer the reader to Steinwart and Christmann (2008) for more details. With this background, the following theorem (proved in Section 6.5) provides a method for constructing uniform confidence bands for the persistence diagram constructed using the robust KDE on  $\mathbb{X}_n$ .

**Theorem 4.4.** *For convex loss  $\rho$  satisfying (A1)–(A4), and fixed  $\sigma > 0$ , suppose the kernel  $K_\sigma$  satisfies  $e_n(\text{id} : \mathcal{H}_\sigma \rightarrow L_\infty(\mathcal{X})) \leq a_\sigma n^{-\frac{1}{2p}}$ , where  $a_\sigma > 1$ ,  $0 < p < 1$  and  $\mathcal{X} \subset \mathbb{R}^d$ . Then, for a fixed confidence level  $0 < \alpha < 1$ ,*

$$\sup_{\mathbb{P} \in \mathcal{M}(\mathcal{X})} \mathbb{P}^{\otimes n} \left\{ W_\infty(\text{Dgm}(f_{\rho, \sigma}^n), \text{Dgm}(f_{\rho, \sigma})) > \frac{2M \|K_\sigma\|_\infty^{\frac{1}{2}}}{\mu} \left( \xi(n, p) + \delta \sqrt{\frac{2 \log(1/\alpha)}{n}} \right) \right\} \leq \alpha,$$

where  $\xi(n, p)$  is given by

$$\xi(n, p) = \begin{cases} \gamma \frac{a_\sigma^p}{(1-2p)} \cdot \frac{1}{\sqrt{n}} & \text{if } 0 < p < \frac{1}{2}, \\ \gamma C \sqrt{a_\sigma} \cdot \frac{\log(n)}{\sqrt{n}} & \text{if } p = \frac{1}{2}, \\ \gamma \frac{p \sqrt{a_\sigma}}{2p-1} \cdot \frac{1}{n^{1/4p}} & \text{if } \frac{1}{2} < p < 1, \end{cases}$$

for fixed constants  $\gamma > \frac{12}{\sqrt{\log 2}}$ ,  $C > 3 - \log(9a_\sigma)$  and  $\mu = 2 \min \left\{ \varphi(2 \|K_\sigma\|_\infty^{\frac{1}{2}}), \rho''(2 \|K_\sigma\|_\infty^{\frac{1}{2}}) \right\}$ .

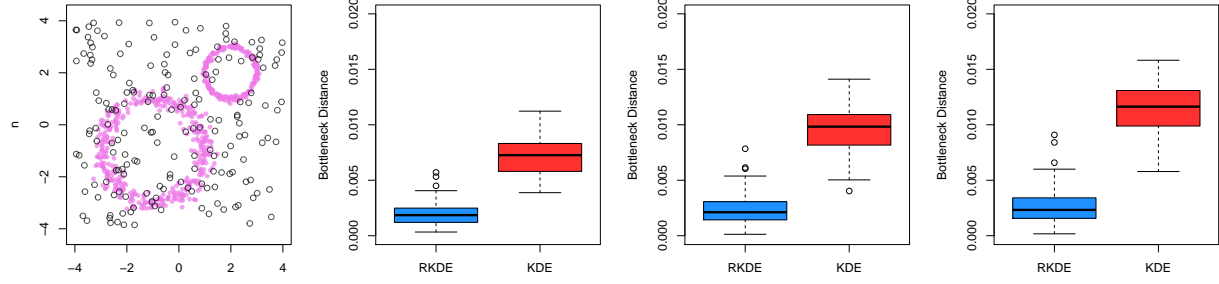
**Remark 4.2.** *We highlight some salient observations from Theorem 4.4.*

(i) *If  $\text{diam}(\mathcal{X}) = r$ , and the kernel  $K_\sigma$  is  $m$ -times differentiable, then from (Steinwart and Christmann, 2008, Theorem 6.26), the entropy numbers associated with  $K_\sigma$  satisfy*

$$e_n(\text{id} : \mathcal{H}_\sigma \rightarrow L_\infty(\mathcal{X})) \leq cr^m n^{-\frac{m}{d}}.$$

*In light of Theorem 4.4, for  $p = \frac{d}{2m}$ , we can make two important observations. First, as the dimension of the input space  $\mathcal{X}$  increases, we have that the rate of convergence decreases; which is a direct consequence from the curse of dimensionality. Second, for a fixed dimension of the input space, the parameter  $p$  in Theorem 4.4 can be understood to be inversely proportional to the smoothness of the kernel. Specifically, as the smoothness of the kernel increases, the rate of convergence is faster, and we obtain sharper confidence bands. This makes a case for employing smoother kernels.*

(ii) *A similar result is obtained in (Fasy et al., 2014, Lemma 8) for persistence diagrams from the KDE, with a convergence rate  $O_p(n^{-1/2})$ , where the proof relies on a simple application of Hoeffding's inequality, unlike the sophisticated tools the proof of Theorem 4.4 warrants for the robust KDE.*



(a)  $\mathbb{X}_n$  (in  $\bullet$ ) and  $\mathbb{Y}_m$  (in  $\circ$ ) (b)  $\pi = 20\%$ ,  $p = 4 \times 10^{-60}$  (c)  $\pi = 30\%$ ,  $p = 2 \times 10^{-72}$  (d)  $\pi = 40\%$ ,  $p = 2.5 \times 10^{-75}$

Figure 4: (a) A realization of  $\mathbb{X}_n \cup \mathbb{Y}_m$ . (b, c, d) As the noise level  $\pi$  increases, boxplots for  $W_\infty(\mathbf{D}_{\rho,\sigma}, \mathcal{D}_\sigma^\#)$  in blue and  $W_\infty(\mathbf{D}_\sigma, \mathcal{D}_\sigma^\#)$  in red show that the robust persistence diagram recovers the underlying signal better.

## 5 Experiments

We illustrate the performance of robust persistence diagrams in machine learning applications through synthetic and real-world experiments. In all the experiments, the kernel bandwidth  $\sigma$  is chosen as the median distance of each  $\mathbf{x}_i \in \mathbb{X}_n$  to its  $k^{\text{th}}$ -nearest neighbour using the Gaussian kernel with the Hampel loss (similar setting as in Kim and Scott, 2012)—we denote this bandwidth as  $\sigma(k)$ . Since DTM is closely related to the  $k$ -NN density estimator (Biau et al., 2011), we choose the DTM smoothing parameter as  $m(k) = k/n$ .

### 5.1 Bottleneck Simulation

The objective of this experiment is to assess how the robust KDE persistence diagram compares to the KDE persistence diagram in recovering the topological features of the underlying signal.  $\mathbb{X}_n$  is observed uniformly from two circles and  $\mathbb{Y}_m$  is sampled uniformly from the enclosing square such that  $m = 200$  and  $m/n = \pi \in \{20\%, 30\%, 40\%\}$ —shown in Figure 4 (a). For each noise level  $\pi$ , and for each of  $N = 100$  realizations of  $\mathbb{X}_n$  and  $\mathbb{Y}_m$ , the robust persistence diagram  $\mathbf{D}_{\rho,\sigma}$  and the KDE persistence diagram  $\mathbf{D}_\sigma$  are constructed from the noisy samples  $\mathbb{X}_n \cup \mathbb{Y}_m$ . In addition, we compute the KDE persistence diagram  $\mathcal{D}_\sigma^\#$  on  $\mathbb{X}_n$  alone as a proxy for the target persistence diagram one would obtain in the absence of any contamination. The bandwidth  $\sigma(k) > 0$  is chosen for  $k = 5$ . For each realization  $i$ , bottleneck distances  $U_i = W_\infty(\mathbf{D}_{\rho,\sigma}, \mathcal{D}_\sigma^\#)$  and  $V_i = W_\infty(\mathbf{D}_\sigma, \mathcal{D}_\sigma^\#)$  are computed for 1<sup>st</sup>-order homological features. The boxplots and  $p$ -values for the one-sided hypothesis test  $H_0 : U - V = 0$  vs.  $H_1 : U - V < 0$  are reported in Figures 4 (b, c, d). The results demonstrate that the robust persistence diagram is noticeably better in recovering the true homological features, and in fact demonstrates superior performance when the noise levels are higher.

### 5.2 Persistence-Influence Experiment

Points  $\mathbb{X}_n$  are sampled from an annular region inside  $[-5, 5]^2$  along with some uniform noise in the ambient space, corresponding to the black points in Figure 5 (a).  $\mathbb{X}_n$  has interesting 1<sup>st</sup>-order homological features. We compute the robust KDE  $f_{\rho,\sigma}^n$  and the KDE  $f_\sigma^n$  on the points  $\mathbb{X}_n$  along with the corresponding persistence diagrams  $\text{Dgm}(f_{\rho,\sigma}^n)$  and  $\text{Dgm}(f_\sigma^n)$ . Outliers  $\mathbb{Y}_m$  are added to the original points at a distance  $r$  from the origin, the number of points roughly equal to  $r$ . Figure 5 (a) depicts these outliers in orange when  $r = 20$ .

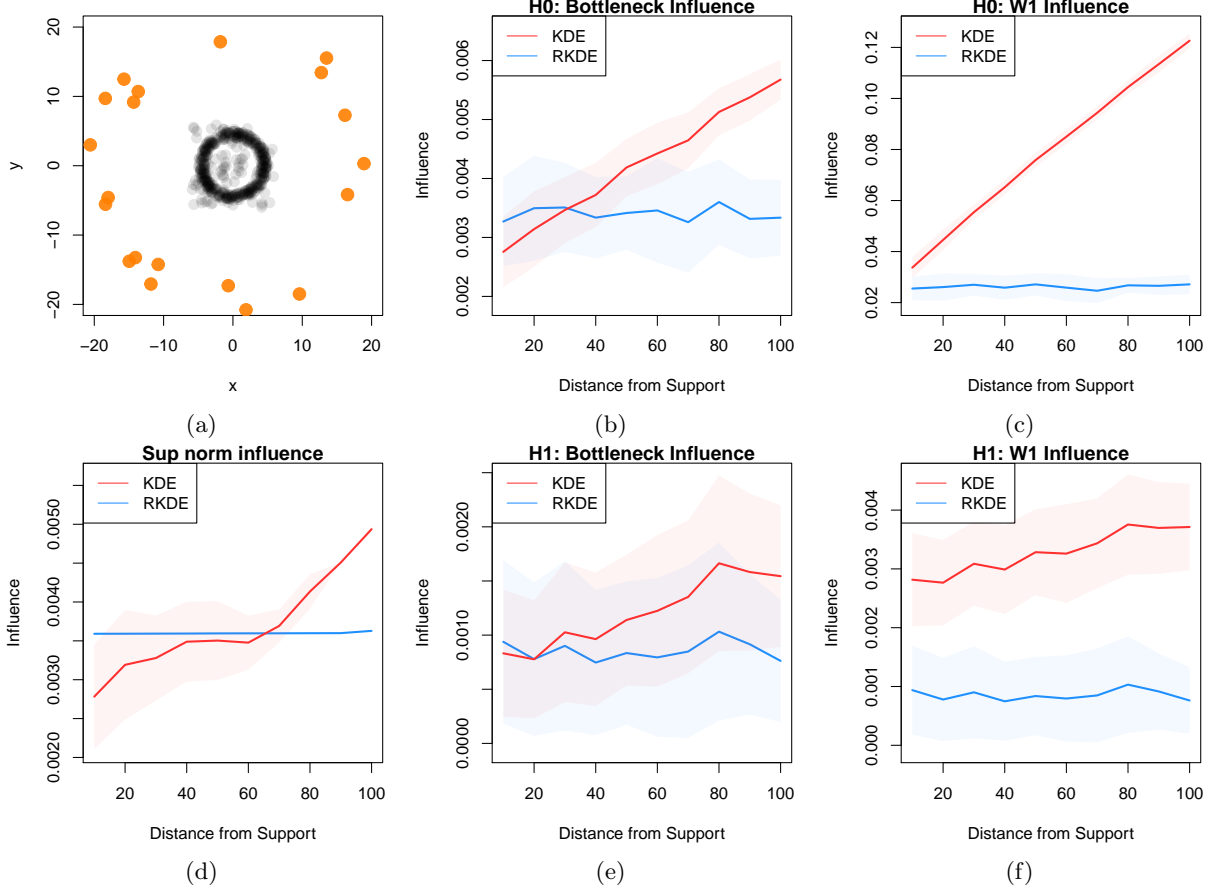


Figure 5: (a) An example of  $\mathbb{X}_n$  in blue and the contamination  $\mathbb{Y}_m$  when  $r = 10$ . (d) The  $L_\infty$  influence of  $\mathbb{Y}_m$  on the KDE and robust KDE. (b, e) The bottleneck influence of  $\mathbb{Y}_m$ . (c, f) The 1-Wasserstein influence of  $\mathbb{Y}_m$  as the distance  $r$  increases.

The robust KDE  $f_{\rho,\sigma}^{n+m}$  and  $f_\sigma^{n+m}$  are now computed on the composite sample  $\mathbb{X}_n \cup \mathbb{Y}_m$  along with the persistence diagrams  $\text{Dgm}(f_{\rho,\sigma}^{n+m})$  and  $\text{Dgm}(f_\sigma^{n+m})$ . The bandwidth  $\sigma(k)$  is chosen for  $k = 5$ . We then compute the  $L_\infty$  influence of  $\mathbb{Y}_m$  i.e.,  $\|f^{n+m} - f^n\|_\infty$  as shown in Figure 5 (d). Additionally for each of the  $0^{th}$ -order and  $1^{st}$ -order persistence diagrams, we compute the persistence influence of  $\mathbb{Y}_m$ , i.e.,  $W_\infty(\text{Dgm}(f^{n+m}), \text{Dgm}(f^n))$  as shown in Figures 5 (b, e), and the 1-Wasserstein influence, i.e.,  $W_1(\text{Dgm}(f^{n+m}), \text{Dgm}(f^n))$  as shown in Figures 5 (c, f). We refer the reader to Eq. (B.1) in Appendix B for the definition of  $W_1$  metric.

For each value of  $r$ , we generate 100 such samples and report the average in Figure 5. Figure 5 (d, e, f) correspond to the experiment illustrated in Figure 3. The results indicate that the robust persistence diagrams,  $\text{Dgm}(f_{\rho,\sigma}^n)$ , are relatively unperturbed when the outliers are added. It exhibits stability even as  $r$  become very large. The KDE persistence diagrams,  $\text{Dgm}(f_\sigma^n)$ , on the other hand, are unstable as the outlying noise becomes more extreme.

As discussed in the Remark 4.1 (iii), the persistence influence for DTM has a much weaker bound as the outliers become more extreme, and in general is not guaranteed to be bounded. In Figure 6 we illustrate the results from the same experiment when the persistence diagrams from DTM is contrasted with the persistence diagrams from the KDE. This analysis is for the same data as that used in Figure 3. We remark that even though DTM is highly sensitive to extreme outliers, DTM

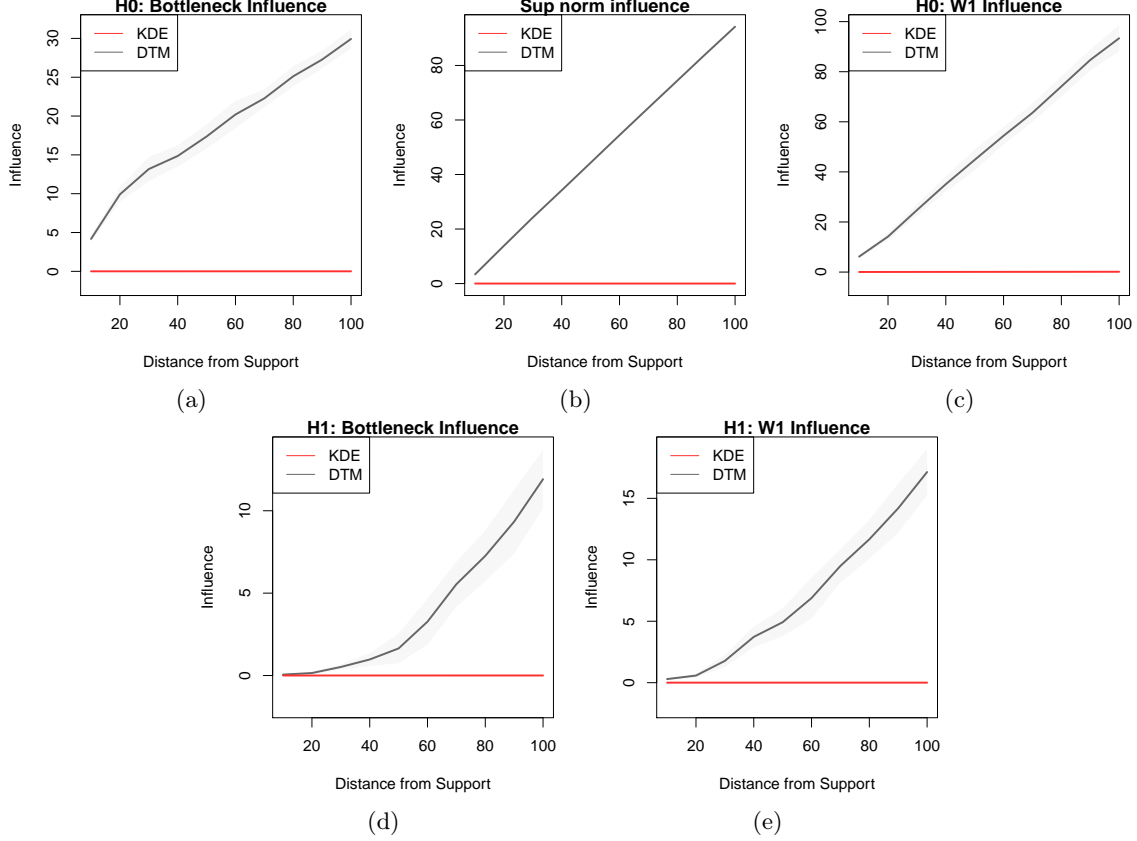


Figure 6: For the same data in Figure 5, (a, d) depicts the bottleneck influence for the DTM in contrast to the KDE – the red line is the same as the one from Figure 5 (b, e). Similarly, in (c, e) we see the  $W_1$  persistence influence of  $\mathbb{Y}_m$  for the DTM in contrast to the KDE. (b) shows the  $L_\infty$  influence of  $\mathbb{Y}_m$  on the DTM. The robust KDE lines are omitted from all plots as it appears to almost merge with the KDE at this scale.

based filtrations have other remarkable properties, as described in Chazal et al. (2017). They are very useful for analyzing persistent homology when one has access to just a single collection of points  $\mathbb{X}_n$ . For DTM the smoothing parameter is chosen as  $m(k) = k/n$  with  $k = 5$ .

### 5.3 Random Circles

The objective of this simulation is to evaluate the performance of persistence diagrams in a supervised learning task. We select circles  $\mathbb{S}_1, \mathbb{S}_2, \dots, \mathbb{S}_N$  randomly in  $\mathbb{R}^2$  with centers inside  $[0, 2]^2$ , with the number of such circles,  $N$  uniformly sampled from  $\{1, 2, \dots, 5\}$ . Conditional on  $N = N$ ,  $\mathbb{X}_n$  is sampled uniformly from  $\mathbb{S}_1, \dots, \mathbb{S}_N$  with 50% noise in the enclosing square. Two such point clouds are shown in Figure 7 (a, b). Persistence diagrams  $\text{Dgm}(f_\sigma^n)$  and  $\text{Dgm}(f_{\rho, \sigma}^n)$  are constructed for bandwidth  $\sigma(k)$  selected from  $k = 5, 7$ , and vectorized in the form of persistence images  $\text{lmg}(f_\sigma^n, h)$ , and  $\text{lmg}(f_{\rho, \sigma}^n, h)$  for varying bandwidths  $h$  (Adams et al., 2017). With  $N$  as the response and the persistence images as the input, results from a support vector regression, averaged over 50 random splits, is shown in Figure 7 (c, d). For a fixed  $h$  the robust persistence diagram seems to always contain more predictive information, as observed in the envelope it forms in Figure 7 (c, d).

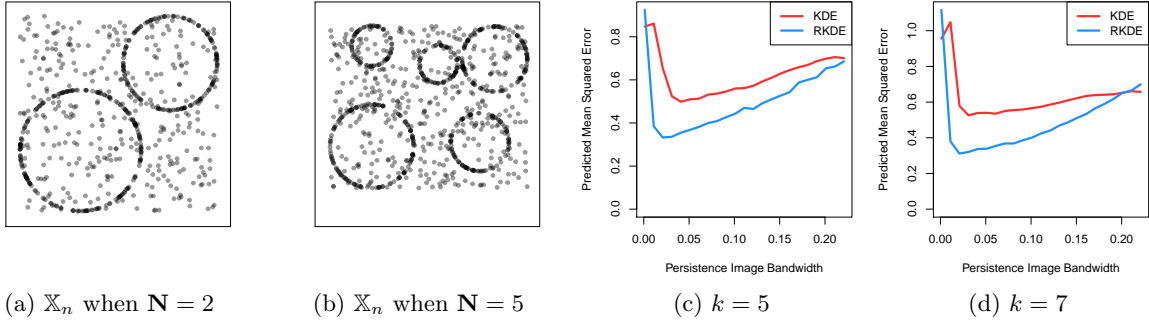


Figure 7: (a, b) A realization  $\mathbb{X}_n$  when  $N = 2$  and  $N = 5$ . (c, d) The predicted mean-squared error vs. the persistence image bandwidth for persistence diagrams in support vector regression.

## 5.4 MPEG7

In this experiment, we examine the performance of persistence diagrams in a classification task using the MPEG7 dataset Latecki et al. (2000). For simplicity, we only consider five classes: *beetle*, *bone*, *spring*, *deer* and *horse*. We first extract the boundary of the images using a Laplace convolution, and sample  $\mathbb{X}_n$  uniformly from the boundary of each image, adding uniform noise ( $\pi = 15\%$ ) in the enclosing region. Persistence diagrams  $\text{Dgm}(f_\sigma^n)$  and  $\text{Dgm}(f_{\rho,\sigma}^n)$  from the KDE and robust KDE are constructed. In addition, owing to their ability to capture nuanced multi-scale features, we also construct  $\text{Dgm}(d_{n,m})$  from the DTM filtration. The smoothing parameters  $\sigma(k)$  and  $m(k)$  are chosen as earlier for  $k = 5$ . The persistence diagrams are normalized to have a max persistence  $\max\{|d - b| = 1 : (b, d) \in \text{Dgm}(\phi)\}$ , and then vectorized as persistence images for various bandwidths  $h$ . A linear SVM classifier is then trained on the resulting persistence images. In the first experiment we only consider the first three classes, and in the second experiment we consider all five classes. The results for the classification error, shown in Figure 8, demonstrate the advantage of the proposed method.

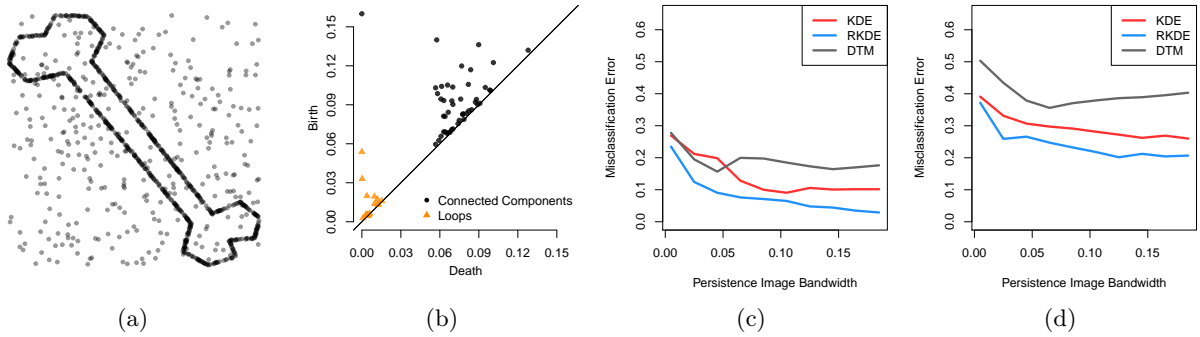


Figure 8: (a)  $\mathbb{X}_n$  is sampled from the image boundary of a *bone*, and uniform noise is added. (b) The resulting persistence diagram from the robust KDE. The persistence diagram picks up the 1<sup>st</sup>-order features near the joints of the cartoon bone. The misclassification error for the KDE, robust KDE and DTM as the persistence image bandwidth increases, (c) for the three-class classification and, (d) for the five-class classification.

## 6 Proofs

In what follows, for a fixed loss  $\rho$ , we will use the notation  $\ell_g(\cdot) = \ell(\cdot, g) = \rho(\|\Phi(\cdot) - g\|_{\mathcal{H}_\sigma})$  in order to emphasize the dependency of the loss on the choice of  $g \in \mathcal{G}$ . Borrowing some notation from empirical process theory, we define the empirical risk-functional in Eq. (1) as

$$\mathcal{J}_n(g) \doteq \mathbb{P}_n \ell_g = \sum_{i=1}^n \rho(\|\Phi_\sigma(\mathbf{X}_i) - g\|_{\mathcal{H}_\sigma}),$$

and, similarly, the population risk functional  $\mathcal{J}(g)$  is given by

$$\mathcal{J}(g) \doteq \mathbb{P} \ell_g = \int_{\mathbb{R}^d} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}).$$

### 6.1 Proof of Theorem 4.1

For  $\epsilon > 0$ , define the risk functional associated with  $\mathbb{P}_\mathbf{x}^\epsilon$  to be

$$\mathcal{J}_{\epsilon, \mathbf{x}}(g) = \mathbb{P}_\mathbf{x}^\epsilon \ell_g = (1 - \epsilon) \mathcal{J}(g) + \epsilon \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}),$$

and let  $f_{\rho, \sigma}^{\epsilon, \mathbf{x}} = \inf_{g \in \mathcal{G}} \mathcal{J}_{\epsilon, \mathbf{x}}(g)$  be its minimizer. From the stability result of Proposition 2.1 we have that

$$\Psi(f_{\rho, \sigma}; \mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} W_\infty(\text{Dgm}(f_{\rho, \sigma}^{\epsilon, \mathbf{x}}), \text{Dgm}(f_{\rho, \sigma})) \leq \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_\infty.$$

Using Propositions A.2 and A.3, we know that the sequence  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  is equi-coercive, and  $\mathcal{J}_{\epsilon, \mathbf{x}}$   $\Gamma$ -converges to  $\mathcal{J}$  as  $\epsilon \rightarrow 0$ . From the fundamental theorem of  $\Gamma$ -convergence (Braides, 2002) we have that  $\|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma} \rightarrow 0$ , and, consequently,  $\|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_\infty \rightarrow 0$  as  $\epsilon \rightarrow 0$ . Thus,

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_\infty = \left\| \lim_{\epsilon \rightarrow 0} \frac{f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}}{\epsilon} \right\|_\infty. \quad (5)$$

Let the limit in the right hand side of Eq. (5) be denoted by  $\dot{f}_{\rho, \sigma}$ . Although  $\dot{f}_{\rho, \sigma}$  does not admit a closed-form solution, from (Kim and Scott, 2012, Theorem 8) we have that  $\dot{f}_{\rho, \sigma}$  satisfies  $V = a \dot{f}_{\rho, \sigma} + B$ , where

$$\begin{aligned} V &= \varphi\left(\|\Phi_\sigma(\mathbf{x}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}\right) \cdot (\Phi_\sigma(\mathbf{x}) - f_{\rho, \sigma}), \\ a &= \int_{\mathbb{R}^d} \varphi\left(\|\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}\right) d\mathbb{P}(\mathbf{y}), \quad \text{and} \\ B &= \int_{\mathbb{R}^d} \left( \frac{\varphi'\left(\|\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}\right)}{\|\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}} \left\langle \dot{f}_{\rho, \sigma}, \Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma} \right\rangle_{\mathcal{H}_\sigma} \cdot (\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}) \right) d\mathbb{P}(\mathbf{y}). \end{aligned}$$

For brevity, we adopt the notation  $z(\mathbf{y}) = \|\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}$  and  $u(\cdot, \mathbf{y}) = \frac{\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}}{\|\Phi_\sigma(\mathbf{y}) - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}} \in \mathcal{H}_\sigma$ . Then note that  $a \in \mathbb{R}$  and  $B \in \mathcal{H}_\sigma$  are given by

$$\begin{aligned} a &= \int_{\mathbb{R}^d} \varphi(z(\mathbf{y})) d\mathbb{P}(\mathbf{y}), \quad \text{and} \\ B &= \int_{\mathbb{R}^d} z(\mathbf{y}) \varphi'(z(\mathbf{y})) \left\langle \dot{f}_{\rho, \sigma}, u(\cdot, \mathbf{y}) \right\rangle_{\mathcal{H}_\sigma} u(\cdot, \mathbf{y}) d\mathbb{P}(\mathbf{y}). \end{aligned}$$

Using the reverse triangle inequality we have

$$\|V\|_{\mathcal{H}_\sigma} \geq a \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma} - \|B\|_{\mathcal{H}_\sigma}. \quad (6)$$

We now look for an upper bound on  $\|B\|_{\mathcal{H}_\sigma}$ . By noting that

$$\left\langle \dot{f}_{\rho,\sigma}, u(\cdot, \mathbf{x}) \right\rangle_{\mathcal{H}_\sigma} \left\langle \dot{f}_{\rho,\sigma}, u(\cdot, \mathbf{y}) \right\rangle_{\mathcal{H}_\sigma} \left\langle u(\cdot, \mathbf{x}), u(\cdot, \mathbf{y}) \right\rangle_{\mathcal{H}_\sigma} \leq \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma}^2,$$

we have

$$\begin{aligned} \|B\|_{\mathcal{H}_\sigma}^2 &= \left\langle B, B \right\rangle_{\mathcal{H}_\sigma} \leq \iint z(\mathbf{x}) \varphi'(z(\mathbf{x})) z(\mathbf{y}) \varphi'(z(\mathbf{y})) \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma}^2 d\mathbb{P}(\mathbf{x}) d\mathbb{P}(\mathbf{y}) \\ &= \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma}^2 \left( \int_{\mathbb{R}^d} z(\mathbf{y}) \varphi'(z(\mathbf{y})) d\mathbb{P}(\mathbf{y}) \right)^2. \end{aligned}$$

Plugging this back into Eq. (6) we get

$$\begin{aligned} \|V\|_{\mathcal{H}_\sigma} &\geq \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma} \int_{\mathbb{R}^d} \varphi(z(\mathbf{y})) - z(\mathbf{y}) \varphi'(z(\mathbf{y})) d\mathbb{P}(\mathbf{y}) \\ &= \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma} \int_{\mathbb{R}^d} \zeta(z(\mathbf{y})) d\mathbb{P}(\mathbf{y}), \end{aligned} \quad (7)$$

where  $\zeta(z) = \varphi(z) - z\varphi'(z)$ . Similarly, by using the definition of  $\varphi$ , it follows that

$$\|V\|_{\mathcal{H}_\sigma} = \varphi \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) \cdot \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} = \rho' \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right).$$

Combining this with Eq. (7) we get

$$\left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma} \leq \frac{\rho' \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right)}{\int_{\mathbb{R}^d} \zeta \left( \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) d\mathbb{P}(\mathbf{y})}.$$

By noting that  $\left\| \dot{f}_{\rho,\sigma} \right\|_\infty \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \left\| \dot{f}_{\rho,\sigma} \right\|_{\mathcal{H}_\sigma}$  and  $\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \left\| \dot{f}_{\rho,\sigma} \right\|_\infty$ , the result follows. ■

## 6.2 Supplementary Results for the Persistence Influence

In this section, we collect the proofs for the results on persistence influence established in Remark 4.1 from Section 4.1. The following result shows that when  $\varphi$  is nonincreasing, the persistence influence in Eq. (3) can be written in a more succinct form.

**Proposition 6.1.** *Under the conditions of Theorem 4.1, if  $\varphi$  is nonincreasing, then the persistence influence of  $\mathbf{x} \in \mathbb{R}^d$  on  $\text{Dgm}(f_{\rho,\sigma})$  satisfies*

$$\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} w_\sigma(\mathbf{x}) \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma},$$

where  $w_\sigma$  is the measure of inlyingness from Eq. (2).



*Proof.* From Theorem 4.1 we have that the persistence influence satisfies

$$\Psi(f_{\rho,\sigma}; \mathbf{x}) \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \rho' \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) \left( \int_{\mathbb{R}^d} \zeta \left( \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) d\mathbb{P}(\mathbf{y}) \right)^{-1}, \quad (8)$$

where  $\zeta(z) = \varphi(z) - z\varphi'(z)$ . When  $\varphi$  is nonincreasing, observe that  $z\varphi'(z) \leq 0$  for all  $0 \leq z < \infty$ . Consequently,  $\zeta$  can be bounded below by  $\varphi$ , and the r.h.s. in Eq. (8) can be bounded above by

$$\begin{aligned} \Psi(f_{\rho,\sigma}; \mathbf{x}) &\stackrel{(i)}{\leq} \|K_\sigma\|_\infty^{\frac{1}{2}} \frac{\rho' \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right)}{\int_{\mathbb{R}^d} \varphi \left( \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) d\mathbb{P}(\mathbf{y})} \\ &\stackrel{(ii)}{=} \|K_\sigma\|_\infty^{\frac{1}{2}} \frac{\varphi \left( \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right)}{\int_{\mathbb{R}^d} \varphi \left( \|\Phi_\sigma(\mathbf{y}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \right) d\mathbb{P}(\mathbf{y})} \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \\ &\stackrel{(iii)}{=} \|K_\sigma\|_\infty^{\frac{1}{2}} w_\sigma(\mathbf{x}) \|\Phi_\sigma(\mathbf{x}) - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}, \end{aligned}$$

where (i) follows from the fact that  $\zeta(z) \geq \varphi(z)$ , (ii) follows from the definition of  $\varphi$ , i.e.,  $\rho'(z) = z\varphi(z)$ , and (iii) follows from the definition of  $w_\sigma$  in Eq. (2), yielding the desired result.  $\blacksquare$

The following result establishes the bound for the distance-to-measure described in Eq. (4).

**Proposition 6.2.** *For  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$ , the persistence influence for the distance-to-measure function is given by*

$$\Psi(d_{\mathbb{P},m}; \mathbf{x}) \leq \frac{2}{\sqrt{m}} \sup \left\{ \left| f(\mathbf{x}) - \int_{\mathbb{R}^d} f(\mathbf{y}) d\mathbb{P}(\mathbf{y}) \right| : \|\nabla f\|_{L_2(\mathbb{P})} \leq 1 \right\}$$

where  $\|\nabla f\|_{L_2(\mathbb{P})}$  is a modified, weighted Sobolev norm (Villani, 2003; Peyre, 2018).

*Proof.* From (Chazal et al., 2011, Theorem 3.5) the following stability result holds:

$$\|d_{\mathbb{P},m} - d_{\mathbb{P}_\mathbf{x}^\epsilon,m}\|_\infty \leq \frac{1}{\sqrt{m}} W_2(\mathbb{P}, \mathbb{P}_\mathbf{x}^\epsilon).$$

From (Peyre, 2018, Theorem 1) we have that

$$W_2(\mathbb{P}, \mathbb{P}_\mathbf{x}^\epsilon) \leq 2 \|\mathbb{P} - \mathbb{P}_\mathbf{x}^\epsilon\|_{\dot{H}^{-1}(\mathbb{P})},$$

where the weighted, homogeneous Sobolev norm  $\|\cdot\|_{\dot{H}^{-1}(\mu)}$  for a signed measure  $\nu$  w.r.t. a positive measure  $\mu$  is given by

$$\|\nu\|_{\dot{H}^{-1}(\mu)} = \sup \left\{ \left| \int_{\mathbb{R}^d} f(\mathbf{x}) d\nu(\mathbf{x}) \right| : \|\nabla f\|_{L_2(\mu)} \leq 1 \right\}.$$

Observe that  $\mathbb{P} - \mathbb{P}_\mathbf{x}^\epsilon = \epsilon(\delta_\mathbf{x} - \mathbb{P})$  and since  $\|\cdot\|_{\dot{H}^{-1}(\mu)}$  defines a norm, we have that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \|d_{\mathbb{P},m} - d_{\mathbb{P}_\mathbf{x}^\epsilon,m}\|_\infty &\leq \frac{1}{\sqrt{m}} \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} W_2(\mathbb{P}, \mathbb{P}_\mathbf{x}^\epsilon) \\ &\leq \frac{2}{\sqrt{m}} \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \|\epsilon(\delta_\mathbf{x} - \mathbb{P})\|_{\dot{H}^{-1}(\mathbb{P})} \\ &= \frac{2}{\sqrt{m}} \|(\delta_\mathbf{x} - \mathbb{P})\|_{\dot{H}^{-1}(\mathbb{P})} \\ &= \frac{2}{\sqrt{m}} \sup \left\{ \left| f(\mathbf{x}) - \int_{\mathbb{R}^d} f(\mathbf{y}) d\mathbb{P}(\mathbf{y}) \right| : \|\nabla f\|_{L_2(\mathbb{P})} \leq 1 \right\}. \end{aligned}$$

From the stability for persistence diagrams, we have that

$$\Psi(d_{\mathbb{P},m}; \mathbf{x}) \leq \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \|d_{\mathbb{P},m} - d_{\mathbb{P}_{\mathbf{x}},m}\|_{\infty}$$

and the result follows. ■

### 6.3 Proof of Theorem 4.2

Using the triangle inequality we can break our problem down as follows

$$\|f_{\rho,\sigma}^n - f\|_{\infty} \leq \underbrace{\|f_{\sigma} - f\|_{\infty}}_{\text{(a)}} + \underbrace{\|f_{\rho,\sigma} - f_{\sigma}\|_{\infty}}_{\text{(b)}},$$

where,  $f_{\sigma} = \int_{\mathbb{R}^d} K_{\sigma}(\cdot, \mathbf{x}) d\mathbb{P}(\mathbf{x})$  is the population level KDE. For term (a), when  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$ , it is well known (Chen, 2017) that the approximation error for the KDE vanishes, i.e.,

$$\|f_{\sigma} - f\|_{\infty} \rightarrow 0,$$

as  $\sigma \rightarrow 0$ . So, it remains to verify that (b) vanishes, i.e.,  $\|f_{\rho,\sigma} - f_{\sigma}\|_{\infty} \rightarrow 0$ . With this in mind, consider the map  $T_{\sigma} : \mathcal{G} \rightarrow \mathcal{G}$  given by

$$T_{\sigma}(g) = \int_{\mathbb{R}^d} \frac{\varphi(\|\Phi_{\sigma}(\mathbf{x}) - g\|_{\mathcal{H}_{\sigma}})}{\int_{\mathbb{R}^d} \varphi(\|\Phi_{\sigma}(\mathbf{x}) - g\|_{\mathcal{H}_{\sigma}}) d\mathbb{P}(\mathbf{x})} \Phi_{\sigma}(\mathbf{x}) d\mathbb{P}(\mathbf{x}).$$

Our approach to verifying that (b) vanishes is similar to Vandermeulen and Scott (2013, Lemma 9), where we show that the map  $T_{\sigma}$  is a contraction map when restricted to the subspace

$$\mathcal{Q}_{\sigma} \doteq B_{\mathcal{H}_{\sigma}}(\mathbf{0}, \delta\nu_{\sigma}) \cap \mathcal{D}_{\sigma}.$$

A key difference is that we work with  $\|\cdot\|_{\infty}$ -norm, requiring us to obtain a sharper bound for the Lipschitz constant associated with the contraction.

For brevity, we adopt the notation  $m(\mathbf{x}, g) = \varphi(\|\Phi_{\sigma}(\mathbf{x}) - g\|_{\mathcal{H}_{\sigma}})$ . Kim and Scott (2012) show that  $f_{\rho,\sigma}$  is a fixed point of the map  $T_{\sigma}$ , i.e.,  $T_{\sigma}(f_{\rho,\sigma}) = f_{\rho,\sigma}$ , and that  $f_{\sigma}$  is the image of  $\mathbf{0}$  under  $T_{\sigma}$ , i.e.,  $T_{\sigma}(\mathbf{0}) = f_{\sigma}$ . Additionally, from Lemma A.2, we know that  $\|f_{\sigma}\|_{\mathcal{H}_{\sigma}} \leq \delta\nu_{\sigma}$ , for some  $0 < \delta < 1$ . Thus, we can rewrite  $f_{\rho,\sigma} - f_{\sigma} = T_{\sigma}(f_{\rho,\sigma}) - T_{\sigma}(\mathbf{0})$ .

Let  $g, h \in \mathcal{Q}_{\sigma}$ . Then we have that

$$\begin{aligned} T_{\sigma}(g) - T_{\sigma}(h) &= \int_{\mathbb{R}^d} \frac{m(\mathbf{x}, g)}{\int_{\mathbb{R}^d} m(\mathbf{y}, g) d\mathbb{P}(\mathbf{y})} \Phi_{\sigma}(\mathbf{x}) d\mathbb{P}(\mathbf{x}) - \int_{\mathbb{R}^d} \frac{m(\mathbf{u}, h)}{\int_{\mathbb{R}^d} m(\mathbf{v}, h) d\mathbb{P}(\mathbf{v})} \Phi_{\sigma}(\mathbf{u}) d\mathbb{P}(\mathbf{u}) \\ &= \frac{1}{\alpha\beta} \cdot \left( \beta \int_{\mathbb{R}^d} m(\mathbf{x}, g) \Phi_{\sigma}(\mathbf{x}) d\mathbb{P}(\mathbf{x}) - \alpha \int_{\mathbb{R}^d} m(\mathbf{u}, h) \Phi_{\sigma}(\mathbf{u}) d\mathbb{P}(\mathbf{x}) \right) \\ &= \frac{1}{\alpha\beta} \cdot \xi, \end{aligned} \tag{9}$$

where  $\alpha \doteq \int_{\mathbb{R}^d} m(\mathbf{y}, g) d\mathbb{P}(\mathbf{y}) \in \mathbb{R}$ ,  $\beta \doteq \int_{\mathbb{R}^d} m(\mathbf{v}, h) d\mathbb{P}(\mathbf{v}) \in \mathbb{R}$  and the numerator  $\xi \in \mathcal{H}_{\sigma}$ .

By Tonelli's theorem

$$\begin{aligned}
\xi &= \beta \int_{\mathbb{R}^d} m(\mathbf{x}, g) \Phi_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) - \alpha \int_{\mathbb{R}^d} m(\mathbf{u}, h) \Phi_\sigma(\mathbf{u}) d\mathbb{P}(\mathbf{x}) \\
&= \int_{\mathbb{R}^d} m(\mathbf{x}, g) \Phi_\sigma(\mathbf{x}) \left( \int_{\mathbb{R}^d} m(\mathbf{v}, h) d\mathbb{P}(\mathbf{v}) \right) d\mathbb{P}(\mathbf{x}) \\
&\quad - \int_{\mathbb{R}^d} m(\mathbf{u}, h) \Phi_\sigma(\mathbf{u}) \left( \int_{\mathbb{R}^d} m(\mathbf{y}, g) d\mathbb{P}(\mathbf{y}) \right) d\mathbb{P}(\mathbf{x}) \\
&= \iint_{\mathbb{R}^d \times \mathbb{R}^d} m(\mathbf{x}, g) m(\mathbf{v}, h) \Phi_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{v}) d\mathbb{P}(\mathbf{x}) - \iint_{\mathbb{R}^d \times \mathbb{R}^d} m(\mathbf{u}, h) m(\mathbf{y}, g) \Phi_\sigma(\mathbf{u}) d\mathbb{P}(\mathbf{y}) d\mathbb{P}(\mathbf{u}) \\
&= \iint_{\mathbb{R}^d \times \mathbb{R}^d} \Phi_\sigma(\mathbf{x}) [m(\mathbf{x}, g) m(\mathbf{y}, h) - m(\mathbf{x}, h) m(\mathbf{y}, g)] d\mathbb{P}(\mathbf{x}) d\mathbb{P}(\mathbf{y}). \tag{10}
\end{aligned}$$

Then by adding and subtracting  $m(\mathbf{x}, h)m(\mathbf{y}, h)$  to the term inside, we get

$$\begin{aligned}
m(\mathbf{x}, g)m(\mathbf{y}, h) - m(\mathbf{x}, h)m(\mathbf{y}, g) &= m(\mathbf{y}, h) \{m(\mathbf{x}, g) - m(\mathbf{x}, h)\} \\
&\quad + m(\mathbf{x}, h) \{m(\mathbf{y}, h) - m(\mathbf{y}, g)\}.
\end{aligned}$$

Plugging this back into Eq. (10), we get  $\xi = \xi_1 + \xi_2$  where

$$\begin{aligned}
\xi_1 &= \iint_{\mathbb{R}^d \times \mathbb{R}^d} \Phi_\sigma(\mathbf{x}) \{m(\mathbf{x}, g) - m(\mathbf{x}, h)\} m(\mathbf{y}, h) d\mathbb{P}(\mathbf{y}) d\mathbb{P}(\mathbf{x}) \\
&= \int_{\mathbb{R}^d} m(\mathbf{y}, h) d\mathbb{P}(\mathbf{y}) \int_{\mathbb{R}^d} \Phi_\sigma(\mathbf{x}) \{m(\mathbf{x}, g) - m(\mathbf{x}, h)\} d\mathbb{P}(\mathbf{x}) \\
&= \beta \int_{\mathbb{R}^d} K_\sigma(\cdot, \mathbf{x}) \{m(\mathbf{x}, g) - m(\mathbf{x}, h)\} d\mathbb{P}(\mathbf{x}) \\
&\stackrel{(i)}{=} \beta \cdot [\psi_\sigma * ((m(\cdot, g) - m(\cdot, h)) f(\cdot))],
\end{aligned}$$

where (i) follows from the fact that the kernel  $K_\sigma(\mathbf{x}, \mathbf{y}) = \psi_\sigma(\mathbf{x} - \mathbf{y}) \doteq \sigma^{-d} \psi(\|\mathbf{x} - \mathbf{y}\|_2 / \sigma)$  is translation invariant and  $f$  is the density associated with  $\mathbb{P}$ . Similarly,

$$\begin{aligned}
\xi_2 &= \iint_{\mathbb{R}^d \times \mathbb{R}^d} \Phi_\sigma(\mathbf{x}) m(\mathbf{x}, h) \{m(\mathbf{y}, h) - m(\mathbf{y}, g)\} d\mathbb{P}(\mathbf{x}) d\mathbb{P}(\mathbf{y}) \\
&= \int_{\mathbb{R}^d} [m(\mathbf{y}, h) - m(\mathbf{y}, g)] d\mathbb{P}(\mathbf{y}) \int_{\mathbb{R}^d} \Phi_\sigma(\mathbf{x}) m(\mathbf{x}, h) d\mathbb{P}(\mathbf{x}) \\
&\leq \|m(\cdot, h) - m(\cdot, g)\|_\infty \cdot [\psi_\sigma * (m(\cdot, h) f(\cdot))].
\end{aligned}$$

The upper bound for  $\|\xi_1\|_\infty$  is as follows

$$\begin{aligned}
\|\xi_1\|_\infty &= \beta \|\psi_\sigma * ((m(\cdot, g) - m(\cdot, h)) f(\cdot))\|_\infty \\
&\stackrel{(i)}{\leq} \beta \|\psi_\sigma\|_1 \| (m(\cdot, g) - m(\cdot, h)) f(\cdot) \|_\infty \\
&\stackrel{(ii)}{\leq} \beta \|m(\cdot, g) - m(\cdot, h)\|_\infty \|f\|_\infty,
\end{aligned} \tag{11}$$

where (i) follows from Young's inequality (Hewitt and Ross, 1979, Theorem 20.18) and (ii) follows from the fact that  $\|fg\|_\infty \leq \|f\|_\infty \|g\|_\infty$ . Similarly, for  $\xi_2$  we have

$$\begin{aligned}
\|\xi_2\|_\infty &\leq \|m(\cdot, h) - m(\cdot, g)\|_\infty \|\psi_\sigma * (m(\cdot, h) f(\cdot))\|_\infty \\
&\stackrel{(i)}{\leq} \|m(\cdot, h) - m(\cdot, g)\|_\infty \|\psi_\sigma\|_1 \|m(\cdot, h) f(\cdot)\|_\infty \\
&\stackrel{(ii)}{\leq} \|m(\cdot, h) - m(\cdot, g)\|_\infty \|m(\cdot, h)\|_\infty \|f\|_\infty.
\end{aligned} \tag{12}$$

From the proof of (Vandermeulen and Scott, 2013, Lemma 9, Page 20–22), for  $g, h \in \mathcal{Q}_\sigma$  for fixed constants  $c_1, c_2 > 0$  we have the following two bounds:

$$\alpha, \beta \geq \frac{1}{c_1 \nu_\sigma}, \tag{13}$$

and

$$\|m(\cdot, h) - m(\cdot, g)\|_\infty \leq \|g - h\|_{\mathcal{H}_\sigma} c_2 \nu_\sigma^{-2}, \tag{14}$$

where the last inequality follows from the Lipschitz property of  $\varphi$  and fact that  $\rho$  is strictly convex. Additionally, for  $c_3 = \|\rho'\|_\infty < \infty$  we have

$$\begin{aligned}
m(\mathbf{x}, g) &= \varphi(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) \\
&= \frac{\rho'(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma})}{\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}} \\
&\leq \frac{c_3}{\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}} \\
&\stackrel{(iii)}{\leq} \frac{c_3}{\left| \|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma} - \|g\|_{\mathcal{H}_\sigma} \right|} \\
&= \frac{c_3}{(1 - \delta) \nu_\sigma},
\end{aligned} \tag{15}$$

where (iii) follows from reverse triangle inequality. Plugging the bounds in equations (13), (14) and (15) back into equations (11) and (12) we get,

$$\|\xi_1\|_\infty + \|\xi_2\|_\infty \leq \|f\|_\infty \left( \beta c_2 \nu_\sigma^{-2} \|g - h\|_{\mathcal{H}_\sigma} + \frac{c_2 c_3}{(1 - \delta)} \nu_\sigma^{-3} \|g - h\|_{\mathcal{H}_\sigma} \right).$$

Using this upper bound in Eq. (9) we get

$$\begin{aligned}
\|T_\sigma(g) - T_\sigma(h)\|_\infty &= \left\| \frac{\xi}{\alpha\beta} \right\|_\infty \\
&\leq \frac{\|\xi_1\|_\infty + \|\xi_2\|_\infty}{\alpha\beta} \\
&\stackrel{(iv)}{\leq} \|f\|_\infty \left( \frac{c_1 c_2}{c_1} \nu_\sigma^{-1} \|g - h\|_{\mathcal{H}_\sigma} + \frac{c_2 c_3}{c_1^2 (1 - \delta)} \nu_\sigma^{-1} \|g - h\|_{\mathcal{H}_\sigma} \right) \\
&\stackrel{(v)}{=} C \nu_\sigma^{-1} \|g - h\|_{\mathcal{H}_\sigma} \\
&\stackrel{(vi)}{\leq} C \nu_\sigma^{-1} \|g - h\|_\infty^{\frac{1}{2}},
\end{aligned}$$

where in (iv) we use Eq. (13), in (v) we use the fact that whenever  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$ , we have  $\|f\|_\infty < \infty$  and  $C > 0$  is a constant depending only on  $c_1, c_2, c_3$  and  $\|f\|_\infty$ . Additionally, (vi) holds through an application of Lemma A.1 to  $g - h \in \mathcal{Q}_\sigma \subset \mathcal{D}_\sigma$ . This confirms that  $T_\sigma$  is a contraction mapping. We use this to show that (b) vanishes as  $\sigma \rightarrow 0$ . Since  $f_{\rho,\sigma}, \mathbf{0} \in \mathcal{Q}_\sigma$  and  $f_{\rho,\sigma} - \mathbf{0} \in \mathcal{D}_\sigma$ , we have that

$$\begin{aligned}
\|f_{\rho,\sigma} - f_\sigma\|_\infty &= \|T_\sigma(f_{\rho,\sigma}) - T_\sigma(\mathbf{0})\|_\infty \\
&\leq C \nu_\sigma^{-1} \|f_{\rho,\sigma} - \mathbf{0}\|_\infty^{\frac{1}{2}} \\
&= C \nu_\sigma^{-1} \|f_{\rho,\sigma}\|_\infty^{\frac{1}{2}}.
\end{aligned}$$

Using the triangle inequality  $\|f_{\rho,\sigma}\|_\infty^{\frac{1}{2}} \leq \|f_{\rho,\sigma} - f_\sigma\|_\infty^{\frac{1}{2}} + \|f_\sigma\|_\infty^{\frac{1}{2}}$  we get

$$\begin{aligned}
\|f_{\rho,\sigma} - f_\sigma\|_\infty &\leq C \nu_\sigma^{-1} \left( \|f_{\rho,\sigma} - f_\sigma\|_\infty^{\frac{1}{2}} + \|f_\sigma\|_\infty^{\frac{1}{2}} \right) \\
&= C \nu_\sigma^{-1} \left( \|T_\sigma(f_{\rho,\sigma}) - T_\sigma(\mathbf{0})\|_\infty^{\frac{1}{2}} + \|f_\sigma\|_\infty^{\frac{1}{2}} \right) \\
&\leq C \nu_\sigma^{-1} \left( \left( C \nu_\sigma^{-1} \|f_{\rho,\sigma} - \mathbf{0}\|_\infty^{\frac{1}{2}} \right)^{\frac{1}{2}} + \|f_\sigma\|_\infty^{\frac{1}{2}} \right) \\
&= C^{\frac{3}{2}} \nu_\sigma^{-\frac{3}{2}} \|f_{\rho,\sigma}\|_\infty^{\frac{1}{4}} + C \nu_\sigma^{-1} \|f_\sigma\|_\infty^{\frac{1}{2}}, \tag{16}
\end{aligned}$$

by using the contraction mapping twice. Observe that

$$\|f_{\rho,\sigma}\|_\infty \leq \nu_\sigma \|f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \leq \delta \nu_\sigma^2,$$

where the first inequality follows from Lemma A.1 and the second inequality follows from the fact that  $\|f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \leq \delta \nu_\sigma$  since  $f_{\rho,\sigma} \in \mathcal{Q}_\sigma$ . Furthermore,  $\|f_\sigma\|_\infty = \|\psi_\sigma * f\|_\infty \leq \|\psi_\sigma\|_\infty \|f\|_1 \leq \nu_\sigma$  from Young's inequality. By noting that  $\nu_\sigma = \psi_\sigma(0) = \sigma^{-d} \psi(0)$ , collecting these bounds back into Eq. (16) we get

$$\|f_{\rho,\sigma} - f_\sigma\|_\infty \leq C^{\frac{3}{2}} \delta^{\frac{1}{4}} \nu_\sigma^{-1} + C \nu_\sigma^{-\frac{1}{2}} \sqrt{\psi(0)}.$$

yielding that  $\|f_{\rho,\sigma} - f_\sigma\|_\infty \rightarrow 0$  as  $\sigma \rightarrow 0$ , thereby verifying that (b) vanishes as  $\sigma \rightarrow 0$ . ■

## 6.4 Proof of Theorem 4.3

The proof proceeds in two steps: We first establish the uniform consistency for the robust KDE and then use the bottleneck stability to show consistency of the robust persistence diagrams in  $W_\infty$ . From the stability theorem for persistence diagrams (Cohen-Steiner et al., 2007; Chazal et al., 2016), we have that  $W_\infty(\text{Dgm}(f_{\rho,\sigma}^n), \text{Dgm}(f_{\rho,\sigma})) \leq \|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_\infty$ . Thus, it suffices to show that  $\|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_\infty \xrightarrow{p} 0$  as  $n \rightarrow \infty$ . In order to prove the latter, we adapt the argmax consistency theorem (Van der Vaart, 2000, Theorem 5.7) for minimizers of a risk function.

**Lemma 6.1** (Theorem 5.7, Van der Vaart (2000)). *Given a metric space  $(\mathcal{G}, d)$ , let  $\mathcal{J}_n$  be random functions and  $\mathcal{J}$  be a fixed function of  $g \in \mathcal{G}$  such that for every  $\epsilon > 0$ ,*

$$(1) \inf_{g: d(g, g_0) \geq \epsilon} \mathcal{J}(g) > \mathcal{J}(g_0), \text{ and}$$

$$(2) \sup_{g \in \mathcal{G}} |\mathcal{J}_n(g) - \mathcal{J}(g)| \xrightarrow{p} 0.$$

*Then any sequence  $g_n$  satisfying  $\mathcal{J}_n(g_n) < \mathcal{J}_n(g_0) + O_p(1)$  satisfies  $d(g_n, g_0) \xrightarrow{p} 0$ .*

For  $\mathcal{G} = \mathcal{H}_\sigma \cap \mathcal{D}_\sigma$ , and  $d(f_{\rho,\sigma}^n, f_{\rho,\sigma}) = \|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_\infty$ , in order to establish uniform consistency of the robust KDE, as per Lemma 6.1, we need to verify that conditions (1) and (2) are satisfied.

Condition (1) follows from the strict convexity of  $\mathcal{J}(g)$  in Proposition A.1. Specifically, Kim and Scott (2012) establish that assumptions (A1) – (A3) guarantee the existence and uniqueness of  $f_{\rho,\sigma} = \arg \inf_{g \in \mathcal{G}} \mathcal{J}(g)$ . Then, for any  $g \in \mathcal{G}$  such that  $\|g - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} > \delta$ , we have that  $\mathcal{J}(g) > \mathcal{J}(f_{\rho,\sigma})$ .

We now turn to verifying condition (2). Observe that  $\sup_{g \in \mathcal{G}} |\mathcal{J}_n(g) - \mathcal{J}(g)|$  can be rewritten as the supremum of an empirical process, i.e.,

$$\sup_{g \in \mathcal{G}} |\mathcal{J}_n(g) - \mathcal{J}(g)| = \sup_{\ell_g \in \tilde{\mathcal{F}}} |\mathbb{P}_n \ell_g - \mathbb{P} \ell_g| \doteq \|\mathbb{P}_n - \mathbb{P}\|_{\tilde{\mathcal{F}}},$$

where  $\tilde{\mathcal{F}} = \{\ell_g : g \in \mathcal{G}\}$ , and  $\ell_g(\mathbf{x}) = \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma})$ . Verifying condition (2) reduces to showing that  $\tilde{\mathcal{F}}$  is a Glivenko-Cantelli class.

Define  $\eta(\cdot) = \|\Phi_\sigma(\cdot) - g\|_{\mathcal{H}_\sigma}^2$  and let  $\mathcal{F} = \{\eta_g : g \in \mathcal{G}\}$ . For the continuous map  $\xi : [0, \infty) \rightarrow [0, \infty)$  given by  $\xi(t) = \rho(\sqrt{t})$ , we have that

$$\xi \circ \mathcal{F} = \{\xi(f) : f \in \mathcal{F}\} = \{\xi \circ \eta_g(\cdot) : g \in \mathcal{G}\} = \{\rho(\|\Phi_\sigma(\cdot) - g\|_{\mathcal{H}_\sigma}) : g \in \mathcal{G}\} = \tilde{\mathcal{F}}.$$

By the preservation theorem for Glivenko-Cantelli classes (Van Der Vaart and Wellner, 2000, Theorem 3), it holds that if  $\mathcal{F}$  is a Glivenko-Cantelli class, then  $\tilde{\mathcal{F}}$  is also a Glivenko-Cantelli class. So verifying condition (2) reduces to verifying that  $\mathcal{F}$  is a Glivenko-Cantelli class.

To this end, we first show that  $F(\mathbf{x}_{1:n}) = F(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \sup_{g \in \mathcal{G}} |\mathbb{P}_n \eta_g - \mathbb{P} \eta_g| = \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$  satisfies the self-bounded property for McDiarmid's inequality, i.e.,

$$\sup_{\mathbf{x}_i \neq \mathbf{x}'_i} |F(\mathbf{x}_{1:n}) - F(\mathbf{x}'_{1:n})| \leq \frac{1}{n} \sup_{\mathbf{x}_i, \mathbf{x}'_i} \sup_{g \in \mathcal{G}} \left( \|\Phi_\sigma(\mathbf{x}_i)\|_{\mathcal{H}_\sigma}^2 + \|\Phi_\sigma(\mathbf{x}'_i)\|_{\mathcal{H}_\sigma}^2 + 2|g(\mathbf{x}_i)| + 2|g(\mathbf{x}'_i)| \right).$$

Observe that  $\|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma}^2 = K_\sigma(\mathbf{x}, \mathbf{x}) \leq \|K_\sigma\|_\infty$  and  $|g(\mathbf{x})| \leq \|g\|_\infty < \|K_\sigma\|_\infty$  by Lemma A.1. Thus, we have that

$$\sup_{\mathbf{x}_i \neq \mathbf{x}'_i} |F(\mathbf{x}_{1:n}) - F(\mathbf{x}'_{1:n})| \leq \frac{6 \|K_\sigma\|_\infty}{n}.$$

From (Bartlett and Mendelson, 2002, Theorem 9), we have that with probability greater than  $1 - e^{-\delta}$ ,

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathfrak{R}_n(\mathcal{F}) + \sqrt{\frac{3\delta \|K_\sigma\|_\infty}{n}}, \quad (17)$$

where  $\mathfrak{R}_n(\mathcal{F})$  is the Rademacher complexity of  $\mathcal{F}$  given by,

$$\begin{aligned} \mathfrak{R}_n(\mathcal{F}) &= \mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \|\Phi_\sigma(\mathbf{x}_i) - g\|_{\mathcal{H}_\sigma}^2 \right| \right) \\ &\leq \mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left\{ \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \|\Phi_\sigma(\mathbf{x}_i)\|_{\mathcal{H}_\sigma}^2 \right| + \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \|g\|_{\mathcal{H}_\sigma}^2 \right| + 2 \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(\mathbf{x}_i) \right| \right\} \right) \\ &= \textcircled{1} + \textcircled{2} + \textcircled{3}. \end{aligned}$$

Note that  $\mathbb{E}_\epsilon(f(\epsilon_{1:n}, \mathbf{x}_{1:n})) \doteq \mathbb{E}(f(\epsilon_{1:n}, \mathbf{x}_{1:n}) | \mathbf{x}_{1:n})$  is the conditional expectation of the Rademacher random variables  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ , keeping  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  fixed. First, we have that,

$$\begin{aligned} \textcircled{1} &= \mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \|\Phi_\sigma(\mathbf{x}_i)\|_{\mathcal{H}_\sigma}^2 \right| \right) \stackrel{(i)}{=} \mathbb{E}_\epsilon \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i K_\sigma(\mathbf{x}_i, \mathbf{x}_i) \right| \\ &\stackrel{(ii)}{\leq} \sqrt{\mathbb{E}_\epsilon \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i K_\sigma(\mathbf{x}_i, \mathbf{x}_i) \right|^2} \\ &\leq \sqrt{\mathbb{E}_\epsilon \left( \frac{1}{n^2} \sum_{i,j} \epsilon_i \epsilon_j K_\sigma(\mathbf{x}_i, \mathbf{x}_i) K_\sigma(\mathbf{x}_j, \mathbf{x}_j) \right)} \\ &\stackrel{(iii)}{=} \frac{1}{\sqrt{n}} \|K_\sigma\|_\infty, \end{aligned}$$

where (i) follows from the absence of  $g$  inside the expectation, (ii) follows from Jensen's inequality and (iii) follows from the fact that  $\epsilon_i \perp \epsilon_j$  for  $i \neq j$ . For the second term, we have

$$\begin{aligned} \textcircled{2} &= \mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \|g\|_{\mathcal{H}_\sigma}^2 \right| \right) = \mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \|g\|_{\mathcal{H}_\sigma}^2 \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \right| \right) \\ &\leq \sup_{g \in \mathcal{G}} \|g\|_{\mathcal{H}_\sigma}^2 \sqrt{\mathbb{E}_\epsilon \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \right|^2} \\ &\stackrel{(iv)}{\leq} \frac{1}{\sqrt{n}} \|K_\sigma\|_\infty, \end{aligned}$$



where (iv) follows from the fact that  $\|g\|_{\mathcal{H}_\sigma}^2 \leq \|K_\sigma\|_\infty$ . Lastly, we have

$$\begin{aligned}
\textcircled{3} &= 2\mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(\mathbf{x}_i) \right| \right) \stackrel{(v)}{=} 2\mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \left| \left\langle g, \frac{1}{n} \sum_{i=1}^n \epsilon_i K_\sigma(\cdot, \mathbf{x}_i) \right\rangle_{\mathcal{H}_\sigma} \right| \right) \\
&\stackrel{(vi)}{\leq} 2\mathbb{E}_\epsilon \left( \sup_{g \in \mathcal{G}} \|g\|_{\mathcal{H}_\sigma} \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i K_\sigma(\cdot, \mathbf{x}_i) \right\|_{\mathcal{H}_\sigma} \right) \\
&= 2 \sup_{g \in \mathcal{G}} \|g\|_{\mathcal{H}_\sigma} \mathbb{E}_\epsilon \left( \sqrt{\frac{1}{n^2} \sum_{i,j} \epsilon_i \epsilon_j K_\sigma(\mathbf{x}_i, \mathbf{x}_j)} \right) \\
&\stackrel{(vii)}{\leq} 2 \frac{\|K_\sigma\|_\infty^{\frac{1}{2}}}{n} \sqrt{\mathbb{E}_\epsilon \left( \sum_{i,j} \epsilon_i \epsilon_j K_\sigma(\mathbf{x}_i, \mathbf{x}_j) \right)} \\
&\stackrel{(viii)}{\leq} \frac{2}{\sqrt{n}} \|K_\sigma\|_\infty,
\end{aligned}$$

where (v) follows from the reproducing property, (vi) is obtained from Cauchy-Schwarz inequality, (vii) follows from Jensen's inequality, and (viii) follows from the fact that  $\epsilon_i \perp \epsilon_j$  for  $i \neq j$ . Collecting these three inequalities, we have

$$\mathfrak{R}_n(\mathcal{F}) = \textcircled{1} + \textcircled{2} + \textcircled{3} \leq \frac{4}{\sqrt{n}} \|K_\sigma\|_\infty.$$

Plugging this into Eq. (17), we have with probability greater than  $1 - e^{-\delta}$ ,

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq \frac{8 \|K_\sigma\|_\infty}{\sqrt{n}} + \sqrt{\frac{3\delta \|K_\sigma\|_\infty}{n}},$$

which implies that  $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \rightarrow 0$  as  $n \rightarrow \infty$ , implying that  $\mathcal{F}$  is a Glivenko-Cantelli class. The result, therefore, follows from Lemma 6.1.  $\blacksquare$

## 6.5 Proof of Theorem 4.4

For  $g \in \mathcal{G}$  define the random fluctuation w.r.t.  $f_{\rho,\sigma}$  as

$$\Delta(\mathbf{X}, g) = (\ell_g(\mathbf{X}) - \ell_{f_{\rho,\sigma}}(\mathbf{X})) - (\mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma})).$$

The fluctuation process is an empirical process defined as

$$\begin{aligned}
\Delta_n(g) &= \mathbb{P}_n \Delta(\mathbf{X}, g) = (\mathcal{J}_n(g) - \mathcal{J}_n(f_{\rho,\sigma})) - (\mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma})), \\
&= \mathbb{P}_n (\ell_g - \ell_{f_{\rho,\sigma}}) - \mathbb{P} (\ell_g - \ell_{f_{\rho,\sigma}}).
\end{aligned}$$

We first show that the behaviour of  $\|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}$  is controlled by the tail behaviour of the supremum of the fluctuation process. To this end, for  $\delta > 0$ , let

$$\mathcal{G}_\delta = \left\{ g \in \mathcal{G} : \|g - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} \leq \delta \right\} = B_{\mathcal{H}_\sigma}(f_{\rho,\sigma}, \delta) \cap \mathcal{D}_\sigma.$$

Suppose  $f_{\rho,\sigma}^n$  is such that  $\|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} > \delta$ , then, for sufficiently small  $\lambda \in (0, 1)$  such that  $g = \lambda f_{\rho,\sigma}^n + (1 - \lambda)f_{\rho,\sigma} \in \mathcal{G}_\delta$ , we have that

$$\begin{aligned} \mathcal{J}_n(g) - \mathcal{J}_n(f_{\rho,\sigma}) &\stackrel{(i)}{<} \lambda \mathcal{J}_n(f_{\rho,\sigma}^n) + (1 - \lambda) \mathcal{J}_n(f_{\rho,\sigma}) - \mathcal{J}_n(f_{\rho,\sigma}) \\ &= \lambda \cdot (\mathcal{J}_n(f_{\rho,\sigma}^n) - \mathcal{J}_n(f_{\rho,\sigma})) \stackrel{(ii)}{\leq} 0, \end{aligned} \quad (18)$$

where (i) follows from the strict convexity of  $\mathcal{J}_n$  (Proposition A.1), and (ii) follows from the fact that  $f_{\rho,\sigma}^n = \arg \inf_{g \in \mathcal{G}} \mathcal{J}_n(g)$ . From Proposition A.1, we also know that  $\mathcal{J}$  is strongly convex such that

$$\mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma}) \geq \frac{\mu}{2} \|g - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}^2. \quad (19)$$

Combining equations (18) and (19) we have

$$\begin{aligned} \frac{\mu}{2} \|g - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma}^2 &\leq \mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma}), \\ &= - \left\{ (\mathcal{J}_n(g) - \mathcal{J}_n(f_{\rho,\sigma})) - (\mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma})) \right\} + (\mathcal{J}_n(g) - \mathcal{J}_n(f_{\rho,\sigma})) \\ &\leq -\Delta_n(g) \leq \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)|. \end{aligned}$$

By taking the supremum of the left hand side in the above inequality over all  $g \in \mathcal{G}_\delta$  we have

$$\sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| \geq \frac{\mu}{2} \delta^2 \quad (20)$$

This implies that whenever  $\|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} > \delta$  holds, then the condition in Eq. (20) holds. Therefore,

$$\mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \|f_{\rho,\sigma}^n - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} > \delta \right\} \leq \mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| \geq \frac{\mu}{2} \delta^2 \right\}. \quad (21)$$

We now study the behaviour of the r.h.s. in Eq. (21) using tools from empirical process theory. First, we show that  $F(\mathbf{x}_{1:n}) = F(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)|$  satisfies the self-bounding property.

$$\begin{aligned} \sup_{\mathbf{x}_i \neq \mathbf{x}'_i} |F(\mathbf{x}_{1:n}) - F(\mathbf{x}'_{1:n})| &= \sup_{\mathbf{x}_i \neq \mathbf{x}'_i} \left| \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| - \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| \right|, \\ &\leq \sup_{\mathbf{x}_i \neq \mathbf{x}'_i} \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g) - \Delta'_n(g)|, \\ &= \frac{1}{n} \sup_{\mathbf{x}_i \neq \mathbf{x}'_i} \sup_{g \in \mathcal{G}_\delta} \left| (\ell_g(\mathbf{x}_i) - \ell_{f_{\rho,\sigma}}(\mathbf{x}_i)) - (\ell_g(\mathbf{x}'_i) - \ell_{f_{\rho,\sigma}}(\mathbf{x}'_i)) \right|, \\ &\leq \frac{1}{n} \sup_{\mathbf{x}_i \neq \mathbf{x}'_i} \sup_{g \in \mathcal{G}_\delta} \left| (\ell_g(\mathbf{x}_i) - \ell_{f_{\rho,\sigma}}(\mathbf{x}_i)) \right| + \left| (\ell_g(\mathbf{x}'_i) - \ell_{f_{\rho,\sigma}}(\mathbf{x}'_i)) \right|, \\ &\stackrel{(i)}{\leq} \frac{1}{n} \sup_{g \in \mathcal{G}_\delta} 2M \|g - f_{\rho,\sigma}\|_{\mathcal{H}_\sigma} = \frac{2M\delta}{n}, \end{aligned}$$

where (i) follows from Proposition A.1 that  $\ell_g$  is  $M$ -Lipschitz w.r.t.  $\|\cdot\|_{\mathcal{H}_\sigma}$ . Therefore, from McDiarmid's inequality (Vershynin, 2018, Theorem 2.9.1) we have

$$\mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| > \mathbb{E} \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| + \epsilon \right\} \leq \exp \left( -\frac{n\epsilon^2}{2M^2\delta^2} \right). \quad (22)$$

Next, we find an upper bound for the expected supremum of the fluctuation process. In order to do so, we first show that  $\Delta_n(g)$  has sub-Gaussian increments. For fixed  $g, h \in \mathcal{G}$  we have that  $\mathbb{E}(\Delta(\mathbf{X}, g) - \Delta(\mathbf{X}, h)) = 0$  and

$$\left| \Delta(\mathbf{X}, g) - \Delta(\mathbf{X}, h) \right| \leq \left| \ell_g(\mathbf{X}) - \ell_h(\mathbf{X}) \right| - \left| \mathcal{J}(g) - \mathcal{J}(h) \right| \leq 2M \|g - h\|_{\mathcal{H}_\sigma}.$$

Since  $\left| \Delta(\mathbf{X}, g) - \Delta(\mathbf{X}, h) \right|$  is bounded, it is, therefore, sub-Gaussian and from Vershynin (2018, Example 2.5.8), we have that the sub-Gaussian norm  $\|\Delta(\mathbf{X}, g) - \Delta(\mathbf{X}, h)\|_{\psi_2} \leq 2cM \|g - h\|_{\mathcal{H}_\sigma}$  for  $c > 1/\sqrt{\log 2}$ . Consequently, the fluctuation process has sub-Gaussian increments with respect to the metric  $\|g - h\|_{\mathcal{H}_\sigma}$ , i.e.,

$$\|\Delta_n(g) - \Delta_n(h)\|_{\psi_2} \leq \frac{1}{n} \sqrt{\sum_{i=1}^n \|\Delta(\mathbf{X}_i, g) - \Delta(\mathbf{X}_i, h)\|_{\psi_2}^2} \leq \frac{M}{\sqrt{n}} \|g - h\|_{\mathcal{H}_\sigma}.$$

From the generalized entropy integral (Srebro et al., 2010, Lemma A.3), for a fixed constant  $\gamma > 12/\sqrt{\log 2}$  we have

$$\mathbb{E} \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| \leq \inf_{\alpha > 0} \left\{ 2\alpha + \frac{\gamma M}{\sqrt{n}} \int_\alpha^\delta \sqrt{\log \mathcal{N}(\mathcal{G}_\delta, \|\cdot\|_{\mathcal{H}_\sigma}, \epsilon)} d\epsilon \right\}, \quad (23)$$

where  $\mathcal{N}(\mathcal{G}_\delta, d, \epsilon)$  is the  $\epsilon$ -covering number of the class  $\mathcal{G}_\delta$  with respect to metric  $d$ .

We now turn our attention to finding an upper bound for  $\mathcal{N}(\mathcal{G}_\delta, d, \epsilon)$ . Note that if  $\mathcal{B}_{\mathcal{H}_\sigma}$  is a unit ball in the RKHS, then

$$\begin{aligned} \log \mathcal{N}(\mathcal{G}_\delta, \|\cdot\|_{\mathcal{H}_\sigma}, \epsilon) &= \log \mathcal{N}(\mathcal{B}_{\mathcal{H}_\sigma} \cap \mathcal{D}_\sigma, \|\cdot\|_{\mathcal{H}_\sigma}, \frac{\epsilon}{\delta}) \\ &\stackrel{(i)}{\leq} \log \mathcal{N}(\mathcal{B}_{\mathcal{H}_\sigma} \cap \mathcal{D}_\sigma, \|\cdot\|_\infty, \left(\frac{\epsilon}{\delta}\right)^2) \\ &\leq \log \mathcal{N}(\mathcal{B}_{\mathcal{H}_\sigma}, \|\cdot\|_\infty, \left(\frac{\epsilon}{\delta}\right)^2), \end{aligned}$$

where (i) follows from Lemma A.1 that  $\|g - h\|_{\mathcal{H}_\sigma}^2 \leq \|g - h\|_\infty$ . When the entropy numbers  $e_n(\text{id} : \mathcal{H}_\sigma \rightarrow L_\infty(\mathcal{X}))$  satisfy the assumption, from (Steinwart and Christmann, 2008, Lemma 6.21) we have

$$\log \mathcal{N}(\mathcal{B}_{\mathcal{H}_\sigma}, \|\cdot\|_\infty, \left(\frac{\epsilon}{\delta}\right)^2) \leq \left(\frac{a_\sigma \delta^2}{\epsilon^2}\right)^{2p}.$$

Plugging this into Eq. (23), we have that

$$\mathbb{E} \sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| \leq \inf_{\alpha > 0} \left\{ 2\alpha + \frac{\gamma M a_\sigma \delta^{2p}}{\sqrt{n}} \int_\alpha^\delta \epsilon^{-2p} d\epsilon \right\} = \inf_{\alpha > 0} T(\alpha),$$

where  $T(\alpha)$  is given by

$$T(\alpha) = \begin{cases} 2\alpha + \gamma M \delta \sqrt{\frac{a_\sigma}{n}} \log\left(\frac{\delta}{\alpha}\right) & \text{if } p = \frac{1}{2}, \\ 2\alpha + \frac{\gamma M}{(1-2p)\sqrt{n}} (\delta - \delta^{2p} \alpha^{1-2p}) & \text{if } 0 < p \neq \frac{1}{2} < 1. \end{cases}$$

At the value  $\alpha_0$  where  $T(\alpha_0) = \inf_{\alpha>0} T(\alpha)$ , we have

$$T(\alpha_0) = \begin{cases} \gamma C a_\sigma^{\frac{1}{2}} \cdot \frac{M\delta \log(n)}{\sqrt{n}} & \text{if } p = \frac{1}{2}, \\ \frac{\gamma a_\sigma^p}{(1-2p)} \cdot \frac{M\delta}{\sqrt{n}} - \frac{K p a_\sigma^{\frac{1}{2}}}{(1-2p)} \cdot \frac{M\delta}{n^{1/4p}} & \text{if } 0 < p \neq \frac{1}{2} < 1, \end{cases} \quad (24)$$

for some fixed constant  $C > 3 - \log(9a)$ . Observe that when  $0 < p < \frac{1}{2}$ , the last term of Eq. (24) is negative, and similarly when  $\frac{1}{2} < p < 1$ , the first term is negative. From this, we have that  $T(\alpha_0) \leq M\delta\xi(n, p)$  where

$$\xi(n, p) = \begin{cases} \frac{\gamma a_\sigma^p}{(1-2p)} \cdot \frac{1}{\sqrt{n}} & \text{if } 0 < p < \frac{1}{2}, \\ \gamma C a_\sigma^{\frac{1}{2}} \cdot \frac{\log(n)}{\sqrt{n}} & \text{if } p = \frac{1}{2}, \\ \frac{\gamma p a_\sigma^{\frac{1}{2}}}{2p-1} \cdot \frac{1}{n^{1/4p}} & \text{if } \frac{1}{2} < p < 1. \end{cases}$$

Plugging this into Eq. (22), we have that with probability greater than  $1 - e^{-t}$ ,

$$\sup_{g \in \mathcal{G}_\delta} |\Delta_n(g)| < M\delta\xi(n, p) + M\delta\sqrt{\frac{2t}{n}}. \quad (25)$$

From Eq. (21), this implies that

$$\mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \|f_{\rho, \sigma}^n - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma} > \delta \right\} \leq \mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \sup_{g \in \mathcal{G}_\delta} \Delta_n(g) \geq \frac{\mu\delta^2}{2} \right\}.$$

Thus, in Eq. (25), by letting

$$\frac{\mu\delta^2}{2} = \left( M\delta\xi(n, p) + M\delta\sqrt{\frac{2t}{n}} \right),$$

we have that with probability greater than  $1 - e^{-t}$ ,

$$\|f_{\rho, \sigma}^n - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma} \leq \frac{2M}{\mu} \left( \xi(n, p) + \sqrt{\frac{2t}{n}} \right).$$

Observe that  $\|f_{\rho, \sigma}^n - f_{\rho, \sigma}\|_\infty \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \|f_{\rho, \sigma}^n - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}$ . For  $0 < \alpha < 1$ , by choosing  $\delta_n$  as

$$\delta_n = \frac{2M \|K_\sigma\|_\infty^{\frac{1}{2}}}{\mu} \left( \xi(n, p) + \sqrt{\frac{2 \log(1/\alpha)}{n}} \right),$$

we have that

$$\mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : \|f_{\rho, \sigma}^n - f_{\rho, \sigma}\|_\infty \leq \delta_n \right\} > 1 - \alpha.$$

From the stability of persistence diagrams in Proposition 2.1, this implies that

$$\mathbb{P}^{\otimes n} \left\{ \mathbf{X}_{1:n} : W_\infty(\text{Dgm}(f_{\rho, \sigma}^n), \text{Dgm}(f_{\rho, \sigma})) > \delta_n \right\} \leq \alpha,$$

yielding the desired result. ■

## 7 Conclusion & Discussion

In this paper, we proposed a statistically consistent robust persistent diagram using RKHS-based robust KDE as the filter function. By generalizing the notion of influence function to the space of persistence diagrams, we mathematically and empirically demonstrated the robustness of the proposed method to that of persistence diagrams induced by other filter functions such as KDE. Through numerical experiments, we demonstrated the advantage of using robust persistence diagrams in machine learning applications.

We would like to highlight that most of the theoretical results of this paper crucially hinge on the loss function being convex. As a future direction, we would like to generalize the current results to non-convex loss functions, which potentially yield more robust persistence diagrams. Another important direction we intend to explore is to develop robust persistence diagrams induced by other types of robust density estimators, which enables to understand the power and limitation of the proposed method.

## Acknowledgements

KF is supported in part by JST CREST Grant Number JPMJCR15D3, Japan. SK is partially supported by JSPS KAKENHI Grant Number JP16H02792.

## References

- Henry Adams, Tegan Emerson, Michael Kirby, Rachel Neville, Chris Peterson, Patrick Shipman, Sofya Chepushtanova, Eric Hanson, Francis Motta, and Lori Ziegelmeier. Persistence Images: A stable vector representation of persistent homology. *Journal of Machine Learning Research*, 18(1): 218–252, 2017.
- Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Paul Bendich, Taras Galkowskyi, and John Harer. Improving homology estimates with random walks. *Inverse Problems*, 27(12):124002, 2011.
- Paul Bendich, James S Marron, Ezra Miller, Alex Pieloch, and Sean Skwerer. Persistent homology analysis of brain artery trees. *The Annals of Applied Statistics*, 10(1):198, 2016.
- Gérard Biau, Frédéric Chazal, David Cohen-Steiner, Luc Devroye, and Carlos Rodriguez. A weighted k-Nearest Neighbor density estimate for geometric inference. *Electronic Journal of Statistics*, 5: 204–237, 2011.
- A Braides.  $\Gamma$ -convergence for Beginners. *Oxford Lecture Series in Mathematics and its Applications*, 22, 2002.
- Rickard Brüel-Gabrielsson, Vignesh Ganapathi-Subramanian, Primoz Skraba, and Leonidas J Guibas. Topology-aware surface reconstruction for point clouds. *arXiv preprint arXiv:1811.12543*, 2018.
- Frédéric Chazal and Bertrand Michel. An introduction to topological data analysis: Fundamental and practical aspects for data scientists. *arXiv preprint arXiv:1710.04019*, 2017.

- Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. Geometric inference for probability measures. *Foundations of Computational Mathematics*, 11(6):733–751, 2011.
- Frédéric Chazal, Leonidas J Guibas, Steve Y Oudot, and Primoz Skraba. Persistence-based clustering in Riemannian manifolds. *Journal of the ACM (JACM)*, 60(6):1–38, 2013.
- Frédéric Chazal, Vin De Silva, Marc Glisse, and Steve Oudot. *The Structure and Stability of Persistence Modules*. Springer, 2016.
- Frédéric Chazal, Brittany Fasy, Fabrizio Lecci, Bertrand Michel, Alessandro Rinaldo, Alessandro Rinaldo, and Larry Wasserman. Robust topological inference: Distance to a measure and kernel distance. *Journal of Machine Learning Research*, 18(1):5845–5884, 2017.
- Chao Chen, Xiuyan Ni, Qinxun Bai, and Yusu Wang. A topological regularizer for classifiers via persistent homology. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2573–2582, 2019.
- Yen-Chi Chen. A tutorial on kernel density estimation and recent advances. *Biostatistics & Epidemiology*, 1(1):161–187, 2017.
- David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. *Discrete & Computational Geometry*, 37(1):103–120, 2007.
- Gianni Dal Maso. *An Introduction to  $\Gamma$ -convergence*, volume 8. Springer Science & Business Media, 2012.
- Vincent Divol and Théo Lacombe. Understanding the topology and the geometry of the persistence diagram space via optimal partial transport. *arXiv preprint arXiv:1901.03048*, 2019.
- Herbert Edelsbrunner and John Harer. *Computational Topology: An Introduction*. American Mathematical Society, Providence, RI, 2010.
- Herbert Edelsbrunner, David Letscher, and Afra Zomorodian. Topological persistence and simplification. In *Proceedings 41st Annual Symposium on Foundations of Computer Science*, pages 454–463. IEEE, 2000.
- Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Larry Wasserman, Sivaraman Balakrishnan, and Aarti Singh. Confidence sets for persistence diagrams. *The Annals of Statistics*, 42(6):2301–2339, 2014.
- Marcio Gameiro, Yasuaki Hiraoka, Shunsuke Izumi, Miroslav Kramar, Konstantin Mischaikow, and Vidit Nanda. A topological measurement of protein compressibility. *Japan Journal of Industrial and Applied Mathematics*, 32(1):1–17, 2015.
- Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust Statistics: The Approach Based on Influence Functions*, volume 196. John Wiley & Sons, 2011.
- Allen Hatcher. *Algebraic Topology*. Cambridge University Press, Cambridge, UK, 2002.
- E. Hewitt and K.A. Ross. *Abstract Harmonic Analysis: Volume I*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 1979. ISBN 0387941908.
- Peter J. Huber. *Robust Statistics*. Wiley Series in Probability and Statistics - Applied Probability and Statistics Section Series. Wiley, 2004. ISBN 9780471650720.

- JooSeuk Kim and Clayton D Scott. Robust kernel density estimation. *Journal of Machine Learning Research*, 13(Sep):2529–2565, 2012.
- Longin Jan Latecki, Rolf Lakamper, and T Eckhardt. Shape descriptors for non-rigid shapes with a single closed contour. In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No. PR00662)*, volume 1, pages 424–429. IEEE, 2000.
- Yuriy Mileyko, Sayan Mukherjee, and John Harer. Probability measures on the space of persistence diagrams. *Inverse Problems*, 27(12):124007, 2011.
- Rémi Peyre. Comparison between  $W_2$  distance and  $\dot{H}^{-1}$  norm, and localization of Wasserstein distance. *ESAIM: Control, Optimisation and Calculus of Variations*, 24(4):1489–1501, 2018.
- Jeff M Phillips, Bei Wang, and Yan Zheng. Geometric inference on kernel density estimates. In *31st International Symposium on Computational Geometry (SoCG 2015)*, 2015.
- Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, pages 2199–2207, 2010.
- Bharath Sriperumbudur. On the optimal estimation of probability measures in weak and strong topologies. *Bernoulli*, 22(3):1839–1893, 2016.
- Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer Science & Business Media, 2008.
- Katharine Turner, Yuriy Mileyko, Sayan Mukherjee, and John Harer. Fréchet means for distributions of persistence diagrams. *Discrete & Computational Geometry*, 52(1):44–70, 2014.
- Aad Van Der Vaart and Jon A Wellner. Preservation theorems for Glivenko-Cantelli and uniform Glivenko-Cantelli classes. In *High dimensional probability II*, pages 115–133. Springer, 2000.
- Aad W Van der Vaart. *Asymptotic Statistics*, volume 3. Cambridge University Press, 2000.
- Robert Vandermeulen and Clayton Scott. Consistency of robust kernel density estimators. In *Conference on Learning Theory*, pages 568–591, 2013.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018. doi: 10.1017/9781108231596.
- Cédric Villani. *Topics in Optimal Transportation*. Number 58. American Mathematical Soc., 2003.
- Larry Wasserman. Topological data analysis. *Annual Review of Statistics and Its Application*, 5: 501–532, 2018.
- Xin Xu, Jessi Cisewski-Kehe, Sheridan B Green, and Daisuke Nagai. Finding cosmic voids and filament loops using topological data analysis. *Astronomy and Computing*, 27:34–52, 2019.
- Afra Zomorodian and Gunnar Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.



## A Supplementary Results

In this section, we establish some results which play a key role in the proofs presented in Section 6.

### A.1 Properties of the Risk Functional $\mathcal{J}(g)$

We establish some important properties of the risk functional, given by

$$\mathcal{J}(g) = \int_{\mathbb{R}^d} \ell_g(\mathbf{x}) \, d\mathbb{P}(\mathbf{x}) = \int_{\mathbb{R}^d} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) \, d\mathbb{P}(\mathbf{x}).$$

The following result establishes that some important properties of the robust loss  $\rho$  carry forward to  $\mathcal{J}(g)$ . (i) The Lipschitz property of  $\rho$  is inherited by  $\mathcal{J}(g)$ , (ii) the convexity of  $\rho$  is strengthened to guarantee that  $\mathcal{J}(g)$  is strictly convex, and (iii)  $\mathcal{J}(g)$  is strongly convex with respect to the  $\|\cdot\|_{\mathcal{H}_\sigma}$ -norm around its minimizer.

**Proposition A.1** (Convexity and Lipschitz properties of  $\mathcal{J}$ ). *Under assumptions (A1) – (A3),*

- (i) *The risk functionals  $\mathcal{J}(g)$  and  $\mathcal{J}_n(g)$  are  $M$ -Lipschitz w.r.t.  $\|\cdot\|_{\mathcal{H}_\sigma}$ .*
- (ii) *Furthermore, if  $\rho$  is convex,  $\mathcal{J}(g)$  and  $\mathcal{J}_n(g)$  are strictly convex.*
- (iii) *Additionally, under assumption (A4), for  $f_{\rho,\sigma} = \arg \inf_{g \in \mathcal{G}} \mathcal{J}(g)$ , the risk functional satisfies the strong convexity condition*

$$\mathcal{J}(g) - \mathcal{J}(f_{\rho,\sigma}) \geq \frac{\mu}{2} \|f_{\rho,\sigma} - g\|_{\mathcal{H}_\sigma}^2,$$

$$\text{for } \mu = 2 \min \left\{ \varphi \left( 2 \|K_\sigma\|_\infty^{\frac{1}{2}} \right), \rho'' \left( 2 \|K_\sigma\|_\infty^{\frac{1}{2}} \right) \right\}.$$

*Proof.* **Lipschitz property.** Observe that,

$$\begin{aligned} |\ell_{g_1}(\mathbf{x}) - \ell_{g_2}(\mathbf{x})| &= |\rho(\|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma}) - \rho(\|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma})| \\ &\leq M \left| \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma} - \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma} \right| \\ &\leq M \|g_1 - g_2\|_{\mathcal{H}_\sigma}, \end{aligned}$$

where the first inequality follows from the fact that  $\rho$  is  $M$ -Lipschitz and the last inequality follows from reverse triangle inequality. This shows that the loss functions  $\ell_g(\cdot)$  are  $M$ -Lipschitz with respect to  $g$ . For the risk functionals, we have that,

$$\begin{aligned} |\mathcal{J}(g_1) - \mathcal{J}(g_2)| &= \left| \int_{\mathbb{R}^d} (\ell_{g_1}(\mathbf{x}) - \ell_{g_2}(\mathbf{x})) \, d\mathbb{P}(\mathbf{x}) \right| \\ &\leq \int_{\mathbb{R}^d} |\ell_{g_1}(\mathbf{x}) - \ell_{g_2}(\mathbf{x})| \, d\mathbb{P}(\mathbf{x}) \\ &\leq M \|g_1 - g_2\|_{\mathcal{H}_\sigma}, \end{aligned}$$

where the first inequality follows from Jensen's inequality. This verifies that  $\mathcal{J}(g)$  is  $M$ -Lipchitz. The proof for  $\mathcal{J}_n(g)$  is identical.

**Strict Convexity.** We begin by establishing that for translation invariant kernels  $\|\Phi_\sigma(\mathbf{x}) - \cdot\|_{\mathcal{H}_\sigma}$  is strictly convex. Suppose  $g_1, g_2 \in \mathcal{H}_\sigma \cap \mathcal{D}_\sigma$  and  $\lambda \in (0, 1)$ , and let  $g = (1 - \lambda)g_1 + \lambda g_2$ . Then

$$\begin{aligned} \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}^2 &= \|(1 - \lambda)(\Phi_\sigma(\mathbf{x}) - g_1) + \lambda(\Phi_\sigma(\mathbf{x}) - g_2)\|_{\mathcal{H}_\sigma}^2 \\ &= (1 - \lambda)^2 \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma}^2 \\ &\quad + \lambda^2 \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma}^2 + 2\lambda(1 - \lambda) \left\langle \Phi_\sigma(\mathbf{x}) - g_1, \Phi_\sigma(\mathbf{x}) - g_2 \right\rangle_{\mathcal{H}_\sigma}. \end{aligned} \quad (\text{A.1})$$

From Cauchy-Schwarz inequality, we know that

$$\left\langle \Phi_\sigma(\mathbf{x}) - g_1, \Phi_\sigma(\mathbf{x}) - g_2 \right\rangle_{\mathcal{H}_\sigma} \leq \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma} \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma}.$$

In the following, we argue that for translation invariant kernels,

$$\left\langle \Phi_\sigma(\mathbf{x}) - g_1, \Phi_\sigma(\mathbf{x}) - g_2 \right\rangle_{\mathcal{H}_\sigma} < \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma} \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma}, \quad (\text{A.2})$$

for  $g_1 \neq g_2$ . On the contrary, suppose

$$\left\langle \Phi_\sigma(\mathbf{x}) - g_1, \Phi_\sigma(\mathbf{x}) - g_2 \right\rangle_{\mathcal{H}_\sigma} = \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma} \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma}$$

holds. Then this implies that there is a function  $a(\mathbf{x})$ , depending only on  $g_1$  and  $g_2$ , such that  $a(\mathbf{x}) \neq 0$  for  $\mathbf{x} \in \mathbb{R}^d$  and

$$\Phi_\sigma(\mathbf{x}) - g_1 = a(\mathbf{x}) (\Phi_\sigma(\mathbf{x}) - g_2).$$

Rearranging the terms this implies that

$$\Phi_\sigma(\mathbf{x}) = \frac{g_1 - a(\mathbf{x})g_2}{1 - a(\mathbf{x})} = (1 + b(\mathbf{x}))g_1 + b(\mathbf{x})g_2,$$

where  $b(\mathbf{x}) = -a(\mathbf{x})/(1 - a(\mathbf{x}))$  also does not vanish on  $\mathbf{x} \in \mathbb{R}^d$ . For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , from the reproducing property we have

$$\begin{aligned} K_\sigma(\mathbf{x}, \mathbf{y}) &= \left\langle \Phi_\sigma(\mathbf{x}), \Phi_\sigma(\mathbf{y}) \right\rangle_{\mathcal{H}_\sigma} \\ &= \left\langle g_1 + b(\mathbf{x})(g_1 + g_2), g_1 + b(\mathbf{y})(g_1 + g_2) \right\rangle_{\mathcal{H}_\sigma} \\ &= b(\mathbf{x})b(\mathbf{y}) \|g_1 + g_2\|_{\mathcal{H}_\sigma}^2 + (b(\mathbf{x}) + b(\mathbf{y})) \langle g_1, g_1 + g_2 \rangle_{\mathcal{H}_\sigma} + \|g_1\|_{\mathcal{H}_\sigma}^2. \end{aligned}$$

Note that because the kernel is translation invariant, i.e.,  $K_\sigma(\mathbf{x}, \mathbf{x}) = K_\sigma(\mathbf{y}, \mathbf{y}) = \sigma^{-d}\psi(0)$ , this must imply that

$$\begin{aligned} 0 &= (b(\mathbf{x})^2 - b(\mathbf{y})^2) \|g_1 + g_2\|_{\mathcal{H}_\sigma}^2 + 2(b(\mathbf{x}) - b(\mathbf{y})) \langle g_1, g_1 + g_2 \rangle_{\mathcal{H}_\sigma} \\ &= (b(\mathbf{x}) - b(\mathbf{y})) \left( (b(\mathbf{x}) + b(\mathbf{y})) \|g_1 + g_2\|_{\mathcal{H}_\sigma}^2 + 2\langle g_1, g_1 + g_2 \rangle_{\mathcal{H}_\sigma} \right). \end{aligned}$$

Since  $b(\mathbf{x})$  and  $b(\mathbf{y})$  are nonvanishing, the above equation is satisfied only when  $b(\mathbf{x}) = b(\mathbf{y})$ . This implies that  $K_\sigma(\mathbf{x}, \mathbf{y})$  is constant for all  $\mathbf{y}$ , giving us a contradiction. Thus, we have that Eq. (A.2) holds. Plugging this back in Eq. (A.1) we get that for  $\lambda \in (0, 1)$  and  $g = (1 - \lambda)g_1 + \lambda g_2$ ,

$$\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} < (1 - \lambda) \|\Phi_\sigma(\mathbf{x}) - g_1\|_{\mathcal{H}_\sigma} + \lambda \|\Phi_\sigma(\mathbf{x}) - g_2\|_{\mathcal{H}_\sigma}.$$

Since,  $\rho$  is strictly increasing and convex, this implies that

$$\ell_g(\mathbf{x}) < (1 - \lambda)\ell_{g_1}(\mathbf{x}) + \lambda\ell_{g_2}(\mathbf{x}).$$

The map  $\ell_g(\cdot) \mapsto \mathbb{P}\ell_g$  is a linear operator, and  $\ell_g$  is strictly convex in  $g$ , this implies that  $\mathcal{J}(g)$  is also strictly convex in  $g$ . The same holds for  $\mathcal{J}_n(g)$ .

**Strong Convexity around the minimizer.** We now turn our attention to the strong convexity property. For this, we first show that  $\mathcal{J}(g)$  is twice Gâteaux differentiable. Let  $g, h \in \mathcal{G}$ , then the second Gâteaux derivative of the loss  $\ell_g(\mathbf{x}) = \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma})$  at  $g$  in the direction  $h$  is given by,

$$\begin{aligned} \delta^2 \ell(\mathbf{x}, g; h) &= \frac{d^2}{d\alpha^2} \ell(\mathbf{x}, g + \alpha h) \Big|_{\alpha=0} \\ &= \frac{d^2}{d\alpha^2} \rho(\|\Phi_\sigma(\mathbf{x}) - g - \alpha h\|_{\mathcal{H}_\sigma}) \Big|_{\alpha=0} \\ &= \frac{d}{d\alpha} \left[ \varphi(\|\Phi_\sigma(\mathbf{x}) - g - \alpha h\|_{\mathcal{H}_\sigma}) \left( -\langle \Phi_\sigma(\mathbf{x}) - g, h \rangle_{\mathcal{H}_\sigma} + \alpha \|h\|_{\mathcal{H}_\sigma}^2 \right) \right] \Big|_{\alpha=0} \\ &= \varphi(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) \|h\|_{\mathcal{H}_\sigma}^2 + \langle \Phi_\sigma(\mathbf{x}) - g, h \rangle_{\mathcal{H}_\sigma}^2 \frac{\varphi'(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma})}{\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}} \\ &= \varphi(z(\mathbf{x}, g)) \|h\|_{\mathcal{H}_\sigma}^2 + \|h\|_{\mathcal{H}_\sigma}^2 \lambda(\mathbf{x}, g, h) z(\mathbf{x}, g) \varphi'(z(\mathbf{x}, g)), \end{aligned} \quad (\text{A.3})$$

where for a fixed  $g \in \mathcal{G}$ , in the interest of brevity, we define  $z(\mathbf{x}, g) = \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}$  and

$$\lambda(\mathbf{x}, g, h) = \left\langle \frac{\Phi_\sigma(\mathbf{x}) - g}{\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}}, \frac{h}{\|h\|_{\mathcal{H}_\sigma}} \right\rangle_{\mathcal{H}_\sigma}^2 \in [0, 1].$$

Observe that  $z\varphi'(z) = \rho''(z) - \varphi(z)$ , thus Eq. (A.3) becomes

$$\delta^2 \ell(\mathbf{x}, g; h) = \|h\|_{\mathcal{H}_\sigma}^2 ((1 - \lambda(\mathbf{x}, g, h)) \varphi(z(\mathbf{x}, g)) + \lambda(\mathbf{x}, g, h) \rho''(z(\mathbf{x}, g))).$$

From assumption (A4) we have that  $\rho''$  and  $\varphi$  are nonincreasing, and

$$z(\mathbf{x}, g) = \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \leq 2 \|K_\sigma\|_\infty^{\frac{1}{2}}.$$

Thus, we have that

$$\delta^2 \ell(\mathbf{x}, g; h) \geq c \|h\|_{\mathcal{H}_\sigma}^2, \quad (\text{A.4})$$

where

$$c = \min \left\{ \varphi \left( 2 \|K_\sigma\|_\infty^{\frac{1}{2}} \right), \rho'' \left( 2 \|K_\sigma\|_\infty^{\frac{1}{2}} \right) \right\}.$$

We also note that  $\delta^2 \ell(\mathbf{x}, g; h)$  is bounded above. To see this, note that from assumption (A4),  $\rho''$  and  $\varphi$  are bounded and nonincreasing. Consequently, for  $\lambda(\mathbf{x}, g, h) \in (0, 1)$  and

$$C = \max \{ \rho''(0), \varphi(0) \} < \infty,$$

from Eq. (A.3) we have that

$$\delta^2 \ell(\mathbf{x}, g; h) \leq C \|h\|_{\mathcal{H}_\sigma}^2 < \infty.$$

The Gâteaux derivative of  $\mathcal{J}(g)$  is, then, given by

$$\begin{aligned}\delta^2 \mathcal{J}(g; h) &= \frac{d^2}{d\alpha^2} \mathcal{J}(g + \alpha h) \Big|_{\alpha=0} = \frac{d^2}{d\alpha^2} \int_{\mathbb{R}^d} \ell(\mathbf{x}, g + \alpha h) d\mathbb{P}(\mathbf{x}) \Big|_{\alpha=0} \\ &= \int_{\mathbb{R}^d} \frac{d^2}{d\alpha^2} \ell(\mathbf{x}, g + \alpha h) d\mathbb{P}(\mathbf{x}) \Big|_{\alpha=0} \\ &= \int_{\mathbb{R}^d} \delta^2 \ell(\mathbf{x}, g; h) d\mathbb{P}(\mathbf{x}).\end{aligned}$$

The exchange of the derivative and integral in the second line follows from the dominated convergence theorem since  $|\delta^2 \ell(\mathbf{x}, g; h)|$  is bounded. This confirms the Gâteaux differentiability of  $\mathcal{J}(g)$ . From Eq. (A.4) we have

$$\delta^2 \mathcal{J}(g; h) = \int_{\mathbb{R}^d} \delta^2 \ell(\mathbf{x}, g; h) d\mathbb{P}(\mathbf{x}) \geq c \|h\|_{\mathcal{H}_\sigma}^2. \quad (\text{A.5})$$

For  $f_{\rho, \sigma} = \arg \inf_{g \in \mathcal{G}} \mathcal{J}(g)$  and  $g \in \mathcal{G}$ , we proceed to show the strong-convexity guarantee. Let  $h = g - f_{\rho, \sigma}$ . From the first-order Taylor approximation for  $\mathcal{J}(g)$  we have,

$$\mathcal{J}(g) = \mathcal{J}(f_{\rho, \sigma}) + \delta \mathcal{J}(f_{\rho, \sigma}, h) + R_2(f_{\rho, \sigma}, h),$$

where the first Gâteaux derivative,  $\delta \mathcal{J}(f_{\rho, \sigma}, h) = 0$  for all  $h$  since  $f_{\rho, \sigma}$  is the unique minimizer of  $\mathcal{J}(g)$  and the remainder term  $R_2(f_{\rho, \sigma}, h)$  is given by

$$\begin{aligned}R_2(f_{\rho, \sigma}, h) &= \frac{1}{2} \int_0^1 (1-t) \delta^2 \mathcal{J}(f_{\rho, \sigma} + th; h) dt \\ &\geq \frac{c}{2} \|h\|_{\mathcal{H}_\sigma}^2 \int_0^1 (1-t) dt = \frac{c}{4} \|h\|_{\mathcal{H}_\sigma}^2,\end{aligned}$$

where the inequality follows from Eq. (A.5). As a result, for any  $g \in \mathcal{G}$  and  $\mu = \frac{c}{2}$  we have that

$$\mathcal{J}(g) - \mathcal{J}(f_{\rho, \sigma}) \geq \frac{\mu}{2} \|g - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma}^2,$$

yielding the desired result. ■

We now turn to examining the behaviour of the risk functional  $\mathcal{J}(g)$  w.r.t. the underlying probability measure  $\mathbb{P}$ . For  $0 \leq \epsilon \leq 1$  and  $\mathbf{x} \in \mathbb{R}^d$ , let  $\mathbb{P}_{\mathbf{x}}^\epsilon = (1 - \epsilon)\mathbb{P} + \epsilon\delta_{\mathbf{x}}$  be a perturbation curve, as defined in Theorem 4.1. The risk functional associated with  $\mathbb{P}_{\mathbf{x}}^\epsilon$  is given by

$$\mathcal{J}_{\epsilon, \mathbf{x}}(g) = \mathbb{P}_{\mathbf{x}}^\epsilon \ell_g = (1 - \epsilon)\mathcal{J}(g) + \epsilon \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}),$$

and  $f_{\rho, \sigma}^{\epsilon, \mathbf{x}} = \inf_{g \in \mathcal{G}} \mathcal{J}_{\epsilon, \mathbf{x}}(g)$  is the minimizer. The convergence of  $f_{\rho, \sigma}^{\epsilon, \mathbf{x}}$  to  $f_{\rho, \sigma}$  can be studied by examining the convergence of  $\mathcal{J}_{\epsilon, \mathbf{x}}$  to  $\mathcal{J}$ . Specifically, under conditions on  $\mathcal{J}$  and  $\mathcal{J}_{\epsilon, \mathbf{x}}$ , it can be shown that  $\|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma} \rightarrow 0$  as  $\epsilon \rightarrow 0$ . The machinery we use here uses the notion of  $\Gamma$ -convergence, which is defined as follows.

**Definition A.1** ( $\Gamma$  convergence). *Given a functional  $F : \mathcal{X} \rightarrow \mathbb{R} \cup \{\pm\infty\}$  and a sequence of functionals  $\{F_n\}_{n \in \mathbb{N}}$ ,  $F_n \xrightarrow{\Gamma} F$  as  $n \rightarrow \infty$  when*

- (i)  $F(\mathbf{x}) \leq \liminf_{n \rightarrow \infty} F_n(\mathbf{x}_n)$  for all  $\mathbf{x} \in \mathcal{X}$  and every  $\{\mathbf{x}_n\}_{n \in \mathbb{N}}$  such that  $d(\mathbf{x}_n, \mathbf{x}) \rightarrow 0$ ;
- (ii) For every  $\mathbf{x} \in \mathcal{X}$ , there exists  $\{\mathbf{x}_n\}_{n \in \mathbb{N}}$ ,  $d(\mathbf{x}_n, \mathbf{x}) \rightarrow 0$  such that  $F(\mathbf{x}) \geq \limsup_{n \rightarrow \infty} F_n(\mathbf{x}_n)$ .

The following result shows that the sequence of functionals  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$   $\Gamma$ -converges to  $\mathcal{J}$ .

**Proposition A.2** ( $\Gamma$ -convergence of  $\mathcal{J}_{\epsilon, \mathbf{x}}$  to  $\mathcal{J}$ ). *Under assumptions (A1)–(A3),*

$$\mathcal{J}_{\epsilon, \mathbf{x}}(g) \xrightarrow{\Gamma} \mathcal{J}(g) \quad \text{as } \epsilon \rightarrow 0.$$

*Proof.* Let  $g \in \mathcal{G}$  and  $\{g_\epsilon\}_{\epsilon > 0}$  be a sequence in  $\mathcal{G}$  such that  $\|g_\epsilon - g\|_{\mathcal{H}_\sigma} \rightarrow 0$  as  $\epsilon \rightarrow 0$ . In order to verify  $\Gamma$ -convergence we first show that the following holds

$$\lim_{\epsilon \rightarrow 0} \left| \mathcal{J}_{\epsilon, \mathbf{x}}(g_\epsilon) - \mathcal{J}(g) \right| = 0.$$

For  $\epsilon > 0$ , using the triangle inequality we have that

$$\begin{aligned} \left| \mathcal{J}_{\epsilon, \mathbf{x}}(g_\epsilon) - \mathcal{J}(g) \right| &\leq \left| \mathcal{J}(g_\epsilon) - \mathcal{J}(g) \right| + \left| \mathcal{J}_{\epsilon, \mathbf{x}}(g_\epsilon) - \mathcal{J}(g_\epsilon) \right| \\ &\stackrel{(i)}{\leq} M \|g_\epsilon - g\|_{\mathcal{H}_\sigma} + \left| \mathcal{J}_{\epsilon, \mathbf{x}}(g_\epsilon) - \mathcal{J}(g_\epsilon) \right| \\ &\stackrel{(ii)}{\leq} M \|g_\epsilon - g\|_{\mathcal{H}_\sigma} + \epsilon \cdot \left| \mathcal{J}(g) - \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) \right|, \end{aligned}$$

where (i) uses the fact that  $\mathcal{J}(g)$  is  $M$ -Lipschitz from Proposition A.1, and (ii) uses the fact that

$$\mathcal{J}_{\epsilon, \mathbf{x}}(g) = (1 - \epsilon)\mathcal{J}(g) + \epsilon\rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}).$$

Since  $\|g_\epsilon - g\|_{\mathcal{H}_\sigma} \rightarrow 0$  as  $\epsilon \rightarrow 0$  we have

$$\lim_{\epsilon \rightarrow 0} \left| \mathcal{J}_{\epsilon, \mathbf{x}}(g_\epsilon) - \mathcal{J}(g) \right| \leq M \lim_{\epsilon \rightarrow 0} \|g_\epsilon - g\|_{\mathcal{H}_\sigma} + \lim_{\epsilon \rightarrow 0} \epsilon \cdot \left| \mathcal{J}(g) - \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) \right| = 0.$$

Since  $\mathcal{J}_{\epsilon, \mathbf{x}}$  and  $\mathcal{J}$  are continuous, using (Dal Maso, 2012, Remark 4.8) it follows that  $\mathcal{J}_{\epsilon, \mathbf{x}}(g) \xrightarrow{\Gamma} \mathcal{J}(g)$ . ■

Now, we examine the coercivity of the sequence  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$ .

**Definition A.2** (Equi-coercivity). *A sequence of functionals  $\{F_n\}_{n \in \mathbb{N}} : \mathcal{X} \rightarrow \mathbb{R} \cup \{\pm\infty\}$  is said to be equi-coercive if for every  $t \in \mathbb{R}$ , there exists a compact set  $K_t \subseteq \mathcal{X}$  such that  $\{\mathbf{x} \in \mathcal{X} : F_n \leq t\} \subseteq K_t$  for every  $n \in \mathbb{N}$ .*

The following result shows that the sequence  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  is equi-coercive.

**Proposition A.3** (Equi-coercivity of  $\mathcal{J}_{\epsilon, \mathbf{x}}$ ). *Under assumptions (A1)–(A3), the sequence of functionals  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  is equi-coercive.*

*Proof.* For  $0 < \epsilon < 1$ ,  $\mathbf{x} \in \mathbb{R}^d$  and  $g \in \mathcal{G}$ , we have that

$$\mathcal{J}_{\epsilon, \mathbf{x}}(g) = (1 - \epsilon)\mathcal{J}(g) + \epsilon\rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}).$$

From (Dal Maso, 2012, Proposition 7.7) in order to show that the sequence of functionals  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  is equi-coercive, it suffices to show that there exists a lower semicontinuous, coercive functional  $F : \mathcal{H}_\sigma \rightarrow \mathbb{R} \cup \{\pm\infty\}$  such that  $F \leq \mathcal{J}_{\epsilon, \mathbf{x}}$  for every  $\epsilon \geq 0$ . To this end consider the functional

$$F(g) = \min \left\{ \mathcal{J}(g), \rho \left( \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \right) \right\}.$$

As  $\mathcal{J}_{\epsilon, \mathbf{x}}$  is a convex combination of  $\mathcal{J}(\cdot)$  and  $\rho \left( \|\Phi_\sigma(\mathbf{x}) - \cdot\|_{\mathcal{H}_\sigma} \right)$ , it implies that  $F \leq \mathcal{J}_{\epsilon, \mathbf{x}}$  for every  $\epsilon \geq 0$ . Additionally, because  $\mathcal{J}(\cdot)$  and  $\rho \left( \|\Phi_\sigma(\mathbf{x}) - \cdot\|_{\mathcal{H}_\sigma} \right)$  are both continuous, it follows that  $F$  is also continuous, and, therefore, lower semicontinuous.

We now verify that  $F$  is coercive. Since  $\rho$  is strictly increasing we have that

$$\rho \left( \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \right) \rightarrow \infty \quad \text{as} \quad \|g\|_{\mathcal{H}_\sigma} \rightarrow \infty,$$

verifying that  $\rho \left( \|\Phi_\sigma(\mathbf{x}) - \cdot\|_{\mathcal{H}_\sigma} \right)$  is coercive. Next, from the reverse triangle inequality we have that

$$\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \geq \left| \|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma} - \|g\|_{\mathcal{H}_\sigma} \right| = \left| \sqrt{K_\sigma(\mathbf{x}, \mathbf{x})} - \|g\|_{\mathcal{H}_\sigma} \right|.$$

Observe that  $K_\sigma(\mathbf{x}, \mathbf{x}) = \|K_\sigma\|_\infty$ , and because  $\rho$  is strictly increasing we have

$$\rho \left( \left| \|K_\sigma\|_\infty^{\frac{1}{2}} - \|g\|_{\mathcal{H}_\sigma} \right| \right) \leq \rho \left( \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \right).$$

Taking expectations on both sides w.r.t.  $\mathbb{P}$ ,

$$\rho \left( \left| \|K_\sigma\|_\infty^{\frac{1}{2}} - \|g\|_{\mathcal{H}_\sigma} \right| \right) \leq \int_{\mathbb{R}^d} \rho \left( \|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma} \right) d\mathbb{P}(\mathbf{x}) = \mathcal{J}(g).$$

Since

$$\rho \left( \left| \|K_\sigma\|_\infty^{\frac{1}{2}} - \|g\|_{\mathcal{H}_\sigma} \right| \right) \rightarrow \infty \quad \text{as} \quad \|g\|_{\mathcal{H}_\sigma} \rightarrow \infty,$$

it implies that  $\mathcal{J}(g)$  is coercive as well. It follows from this that  $F$  is coercive, and the sequence of functionals  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  is equi-coercive.  $\blacksquare$

Propositions A.2 and A.3 together imply, from the fundamental theorem of  $\Gamma$ -convergence (Braides, 2002), that the sequence of minimizers associated with  $\{\mathcal{J}_{\epsilon, \mathbf{x}}\}$  converge to the minimizer of  $\mathcal{J}$ , i.e.,

$$\|f_{\rho, \sigma}^{\epsilon, \mathbf{x}} - f_{\rho, \sigma}\|_{\mathcal{H}_\sigma} \rightarrow 0 \quad \text{as} \quad \epsilon \rightarrow 0.$$

## A.2 Some Additional Results

Next, we note an important property of the hypothesis class,  $\mathcal{G} = \mathcal{H}_\sigma \cap \mathcal{D}_\sigma$ . The elements of  $\mathcal{G}$  can be shown to have their  $\|\cdot\|_\infty$ -norm related their  $\|\cdot\|_{\mathcal{H}_\sigma}$ -norm.

**Lemma A.1** (Vandermeulen and Scott, 2013, Lemma 6 and Sriperumbudur, 2016, Proposition 5.1). *For every  $g \in \mathcal{H}_\sigma \cap \mathcal{D}_\sigma$ ,*

$$\|g\|_{\mathcal{H}_\sigma}^2 \leq \|g\|_\infty \leq \|K_\sigma\|_\infty^{\frac{1}{2}} \|g\|_{\mathcal{H}_\sigma}.$$

The following result, which is essentially the population analogue of Vandermeulen and Scott (2013, Lemma 7), guarantees that for small enough  $\sigma > 0$ , there exists  $0 < \delta < 1$  such that  $f_{\rho,\sigma}$  is contained in the RKHS ball  $B_{\mathcal{H}_\sigma}(\mathbf{0}, \delta\nu_\sigma)$ , where for brevity we denote  $\nu_\sigma = \|K_\sigma\|_\infty^{1/2}$ . We provide the proof for completeness, however, the proof uses exactly the same ideas from Vandermeulen and Scott (2013). For notational convenience, we also define  $\psi_\sigma(\|\mathbf{x} - \mathbf{y}\|_2) = K_\sigma(\mathbf{x}, \mathbf{y}) = \sigma^{-d}\psi(\|\mathbf{x} - \mathbf{y}\|_2/\sigma)$ .

**Lemma A.2.** *Let  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$  and  $f_{\rho,\sigma}$  be the robust KDE for  $\sigma > 0$ . For sufficiently small  $\sigma > 0$ , there exists  $0 < \delta < 1$  such that  $f_{\rho,\sigma} \in B(\mathbf{0}, \delta\nu_\sigma)$ .*

*Proof.* For  $\mathbb{P} \in \mathcal{M}(\mathbb{R}^d)$ , and  $\mathcal{G} = \mathcal{H}_\sigma \cap \mathcal{D}_\sigma$ , consider the map  $T_\sigma : \mathcal{G} \rightarrow \mathcal{G}$  given by

$$T_\sigma(g) = \int_{\mathbb{R}^d} \frac{\varphi(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma})}{\int_{\mathbb{R}^d} \varphi(\|\Phi_\sigma(\mathbf{y}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{y})} K_\sigma(\cdot, \mathbf{x}) d\mathbb{P}(\mathbf{x}) = \int_{\mathbb{R}^d} K_\sigma(\cdot, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}),$$

for each  $g \in \mathcal{G}$ . Observe that  $w_\sigma \in L_1(\mathbb{P})$  is a non-negative function such that

$$\int_{\mathbb{R}^d} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) = 1. \quad (\text{A.6})$$

Let  $S_\sigma = \text{Im}(T_\sigma) \subset \mathcal{G}$ . It follows from Vandermeulen and Scott (2013, Page 11) that the robust KDE,  $f_{\rho,\sigma} = \arg\inf_{g \in \mathcal{G}} \mathcal{J}(g)$ , is the fixed point of the map  $T_\sigma$  and therefore  $f_{\rho,\sigma} \in S_\sigma$ . For a small  $\epsilon > 0$ , from Vandermeulen and Scott (2013, Lemma 12; Corollary 13) there exist  $r, s > 0$  such that  $\mathbb{P}(B(\mathbf{x}, r)) \leq \epsilon$  and  $\mathbb{P}(B(\mathbf{x}, r+s) \setminus B(\mathbf{x}, r)) \leq \epsilon$  for all  $\mathbf{x} \in \mathbb{R}^d$ . This implies that  $\mathbb{P}(B(\mathbf{x}, r+s)^c) > 1 - 2\epsilon$ . We point out that the constant  $\epsilon$  chosen here is related to  $\sqrt{9/10}$  used by Vandermeulen and Scott (2013) as  $\sqrt{1-\epsilon} = \sqrt{9/10}$ , which, as remarked by the authors, was chosen simply for convenience. Define the sets  $B_\sigma = B_{\mathcal{H}_\sigma}(\mathbf{0}, \nu_\sigma\sqrt{1-\epsilon})$ , and let

$$R_\sigma \doteq S_\sigma \cap B_\sigma^c.$$

In what follows we will show that  $f_{\rho,\sigma}$  does not lie in  $R_\sigma$ . To this end, let  $g = \arg\inf_{h \in R_\sigma} \mathcal{J}(h)$ . It suffices to show that  $\mathcal{J}(g) > \mathcal{J}(\mathbf{0}) > \mathcal{J}(f_{\rho,\sigma})$ . Since  $g \in R_\sigma$ , it must follow that

$$(1-\epsilon)\nu_\sigma^2 < \|g\|_{\mathcal{H}_\sigma}^2 \leq \|g\|_\infty = g(\mathbf{z}), \quad (\text{A.7})$$

for some  $\mathbf{z} \in \mathbb{R}^d$ , where the second inequality follows from Lemma A.1. Since  $g \in S_\sigma$ , there exists a non-negative function  $w_\sigma$  satisfying Eq. (A.6), such that  $g = \int_{\mathbb{R}^d} w_\sigma(\mathbf{x}) K_\sigma(\cdot, \mathbf{x}) d\mathbb{P}(\mathbf{x})$ . Therefore,

$$\begin{aligned} (1-\epsilon)\nu_\sigma^2 &\leq g(\mathbf{z}) = \int_{\mathbb{R}^d} K_\sigma(\mathbf{z}, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \\ &= \int_{B(\mathbf{z}, r)} K_\sigma(\mathbf{z}, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) + \int_{B(\mathbf{z}, r)^c} K_\sigma(\mathbf{z}, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \\ &\stackrel{(i)}{\leq} \nu_\sigma^2 \int_{B(\mathbf{z}, r)} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) + \underbrace{\psi_\sigma(r) \int_{B(\mathbf{z}, r)^c} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x})}_{\leq 1} \\ &\stackrel{(ii)}{\leq} \nu_\sigma^2 \int_{B(\mathbf{z}, r)} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) + \psi_\sigma(r), \end{aligned} \quad (\text{A.8})$$



where (i) follows from the fact that  $\sup_{B(\mathbf{z}, r)^c} K_\sigma(\mathbf{z}, \mathbf{x}) = \psi_\sigma(r)$  and (ii) follows from Eq. (A.6). From Vandermeulen and Scott (2013, Lemma 7), there exists  $\sigma$  small enough such that  $\psi_\sigma(r) < \frac{\epsilon}{2} \nu_\sigma^2$ . Plugging this back in Eq. (A.8) we get

$$\int_{B(\mathbf{z}, r)} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \geq \left(1 - \frac{3\epsilon}{2}\right). \quad (\text{A.9})$$

Additionally,

$$\begin{aligned} \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} g(\mathbf{y}) &= \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} \left( \int_{B(\mathbf{z}, r)} K_\sigma(\mathbf{y}, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) + \int_{B(\mathbf{z}, r)^c} K_\sigma(\mathbf{y}, \mathbf{x}) w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \right) \\ &\leq \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} \sup_{\mathbf{x} \in B(\mathbf{z}, r)} K_\sigma(\mathbf{y}, \mathbf{x}) \int_{B(\mathbf{z}, r)} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \\ &\quad + \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} \sup_{\mathbf{x} \in B(\mathbf{z}, r)} K_\sigma(\mathbf{y}, \mathbf{x}) \int_{B(\mathbf{z}, r)^c} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}) \\ \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} g(\mathbf{y}) &\leq \psi_\sigma(s) + \nu_\sigma^2 \int_{B(\mathbf{z}, r)^c} w_\sigma(\mathbf{x}) d\mathbb{P}(\mathbf{x}). \end{aligned}$$

For a choice of  $\tau > 0$ , there is  $\sigma$  small enough satisfying  $\psi_\sigma(s) \leq \tau$  such that from Eq. (A.9)

$$\sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} g(\mathbf{y}) \leq \tau + \frac{3\epsilon}{2} \nu_\sigma^2. \quad (\text{A.10})$$

Then we have that

$$\begin{aligned} \mathcal{J}(g) &= \int_{\mathbb{R}^d} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}) \\ &= \int_{B(\mathbf{z}, r+s)} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}) + \int_{B(\mathbf{z}, r+s)^c} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}) \\ &\geq \int_{B(\mathbf{z}, r+s)^c} \rho(\|\Phi_\sigma(\mathbf{x}) - g\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}) \\ &= \int_{B(\mathbf{z}, r+s)^c} \rho\left(\sqrt{\nu_\sigma^2 + \|g\|_{\mathcal{H}_\sigma}^2} - 2\langle g, \Phi_\sigma(\mathbf{x}) \rangle_{\mathcal{H}_\sigma}\right) d\mathbb{P}(\mathbf{x}) \\ &\geq \int_{B(\mathbf{z}, r+s)^c} \rho\left(\sqrt{\nu_\sigma^2 + \|g\|_{\mathcal{H}_\sigma}^2} - 2 \sup_{\mathbf{y} \in B(\mathbf{z}, r+s)^c} g(\mathbf{y})\right) d\mathbb{P}(\mathbf{x}). \end{aligned}$$

Plugging in Equations (A.10) and (A.7) we get

$$\mathcal{J}(g) \geq (1 - 2\epsilon) \rho\left(\sqrt{(2 - 4\epsilon)\nu_\sigma^2 - 2\tau}\right).$$

Since  $\rho$  is assumed to be strictly convex, this implies that  $\rho'$  is strictly increasing. Additionally, from (A2) we have that  $\rho'$  is bounded. This implies that, for any  $0 < \alpha < \|\rho'\|_\infty$ , there is  $\beta > 0$  such that  $\rho'(z) > \|\rho'\|_\infty - \alpha$  for all  $z > \beta$ . Using Vandermeulen and Scott (2013, Eq. 11), we have

$$\begin{aligned} \rho\left(\sqrt{(2-4\epsilon)\nu_\sigma^2-2\tau}\right) &= \int_0^{(2-4\epsilon)\nu_\sigma^2-2\tau} \rho'(z)dz \\ &\geq \int_\beta^{(2-4\epsilon)\nu_\sigma^2-2\tau} \rho'(z)dz \\ &\geq \int_\beta^{(2-4\epsilon)\nu_\sigma^2-2\tau} (\|\rho'\|_\infty - \alpha) dz \\ &\geq (\|\rho'\|_\infty - \alpha) \left(\sqrt{(2-4\epsilon)\nu_\sigma^2-2\tau} - \beta\right). \end{aligned}$$

Without loss of generality, we can assume  $\|\rho'\|_\infty = 1$ . Choosing  $\alpha$ ,  $\tau$  and  $\sigma$  small enough we obtain

$$\mathcal{J}(g) \geq \nu_\sigma.$$

Now we note that

$$\begin{aligned} \mathcal{J}(\mathbf{0}) &= \int_{\mathbb{R}^d} \rho(\|\Phi_\sigma(\mathbf{x})\|_{\mathcal{H}_\sigma}) d\mathbb{P}(\mathbf{x}) \\ &= \rho(\nu_\sigma) \\ &= \rho(0) + \int_0^{\nu_\sigma} \rho'(z)dz \\ &\leq \rho(0) + \|\rho'\|_\infty \int_0^{\nu_\sigma} dz = \nu_\sigma. \end{aligned}$$

Thus, we obtain that  $\mathcal{J}(g) > \mathcal{J}(\mathbf{0})$ . We have  $g = \arg \inf_{h \in R_\sigma} \mathcal{J}(h)$  and  $f_{\rho,\sigma} = \arg \inf_{h \in \mathcal{G}} \mathcal{J}(h)$ , and, additionally we know that  $f_{\rho,\sigma} \neq \mathbf{0}$ . It follows that since  $\mathcal{J}(f_{\rho,\sigma}) \leq \mathcal{J}(\mathbf{0}) < \mathcal{J}(g)$ , then  $f_{\rho,\sigma} \notin R_\sigma$  as  $\sigma \rightarrow 0$ . Taking  $\delta = \sqrt{1-\epsilon}$ , we get the desired result.  $\blacksquare$

## B Background on Persistent Homology

Given a set of points  $\mathbb{X}_n = \{\mathbf{x}_1 \dots \mathbf{x}_n\}$  in a metric space  $(\mathcal{X}, d)$  their topology is encoded in a geometric object called a simplicial complex  $\mathcal{K} \subseteq 2^{\mathbb{X}_n}$ .

**Definition B.1.** (Hatcher, 2002). *A simplicial complex  $\mathcal{K}$  is a collection of simplices  $\langle \sigma \rangle$  i.e. points, lines, triangles, tetrahedra and its higher dimensional analogues, such that*

1.  $\forall \tau \preceq \sigma, \sigma \in \mathcal{K}$  we have  $\tau \in \mathcal{K}$ ;
2.  $\forall \sigma, \tau \in \mathcal{K}$ , we have that  $\sigma \cap \tau \preceq \sigma, \tau$  or  $\sigma \cap \tau = \emptyset$ .

For a given spatial resolution  $r > 0$ , the simplicial complex for  $\mathbb{X}_n$ , given by  $\mathcal{K}(\mathbb{X}_n, r)$ , can be constructed in multiple ways. For example, the Vietoris-Rips complex is the simplicial complex

$$\mathcal{K}_r = \{\sigma \subseteq \mathbb{X}_n : \bigcap_{\mathbf{x} \in \sigma} B(\mathbf{x}, r) \neq \emptyset\},$$

and the Čech complex is given by

$$\mathcal{K}_r = \{\sigma \subseteq \mathbb{X}_n : \max_{\mathbf{x}_i, \mathbf{x}_j \in \sigma} d(\mathbf{x}_i, \mathbf{x}_j) \leq r\}.$$

More generally, if  $\mathcal{K}$  is a simplicial complex constructed using an approximation of the space  $\mathcal{X}$  (e.g., triangulation, surface mesh, grid, etc.), and  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  a filter function,  $\phi$  induces the map  $\phi : \mathcal{K} \rightarrow \mathbb{R}$ . Then,  $\mathcal{K}_r = \phi^{-1}([0, r])$  encodes the information in the sublevel set of  $\phi$  at resolution  $r$ . Similarly,  $\mathcal{K}^r$  encodes the information in the superlevel sets at resolution  $r$ .

For  $0 \leq k \leq d$ , the  $k^{th}$ -homology (Hatcher, 2002) of a simplicial complex  $\mathcal{K}$ , given by  $H_k(\mathcal{K})$  is an algebraic object encoding its topology as a vector-space (over a fixed field). Using the Nerve lemma,  $H_k(\mathcal{K}(\mathbb{X}_n, r))$  is isomorphic to the homology of its union of  $r$ -balls,  $H_k(\bigcup_{i=1}^n B_r(\mathbf{x}_i))$ . The ordered sequence  $\{\mathcal{K}(\mathbb{X}_n, r)\}_{r>0}$  forms a *filtration*, encoding the evolution of topological features over a spectrum of resolutions. For  $0 < r < s$ , the simplicial complex  $\mathcal{K}(\mathbb{X}_n, r)$  is a *sub-simplicial complex* of  $\mathcal{K}(\mathbb{X}_n, s)$ . Their homology groups are associated with the inclusion maps

$$\iota_r^s : H_k(\mathcal{K}(\mathbb{X}_n, r)) \hookrightarrow H_k(\mathcal{K}(\mathbb{X}_n, s)),$$

which in turn carry information on the number of non-trivial  $k$ -cycles. As the resolution  $r$  varies, the evolution of the topology is captured in the filtration. Roughly speaking, new cycles (e.g., connected components, loops, voids and higher order analogues) can appear or existing cycles can merge. Formally, a new  $k$ -cycle  $\sigma_k$  with homology class  $[\alpha_k]$  is *born* at  $b \in \mathbb{R}$  if  $[\alpha_k] \notin \text{Im}(\iota_{b-\epsilon, b}^k)$  for all  $\epsilon > 0$  and  $[\alpha_k] \in \text{Im}(\iota_{b, b+\delta}^k)$  for some  $\delta > 0$ . The same  $k$ -cycle born at  $b$  dies at  $d > b$  if  $\iota_{b, d-\delta}^k([\alpha_k]) \notin \text{Im}(\iota_{b-\epsilon, d-\delta}^k)$  and  $\iota_{b, d}^k([\alpha_k]) \in \text{Im}(\iota_{b-\epsilon, d}^k)$  for all  $\epsilon > 0$  and  $0 < \delta < d - b$ . Persistent homology,  $PH_*(\phi)$ , is an algebraic module which tracks the persistence pairs  $(b, d)$  of births  $b$  and deaths  $d$  across the entire filtration. By collecting all persistence pairs  $(b, d)$ , the persistent homology is represented as a persistence diagram

$$\text{Dgm}(\mathcal{K}(\mathbb{X}_n)) \doteq \{(b, d) \in \mathbb{R}^2 : 0 \leq b < d \leq \infty\}.$$

The persistence diagram is a multiset of points on the space  $\Omega = \{(x, y) : 0 \leq x < y \leq \infty\}$ , such that each point  $(x, y)$  in the persistence diagram corresponds to a distinct topological feature which existed in  $\mathcal{K}(\mathbb{X}_n, r)$  for  $x \leq r < y$ . Given a persistence diagram  $\mathbf{D}$  and  $1 \leq p \leq \infty$  the *degree- $p$  total persistence* of  $\mathbf{D}$  is given by

$$\text{pers}_p(\mathbf{D}) = \left( \sum_{(b, d) \in \mathbf{D}} |d - b|^p \right)^{\frac{1}{p}}.$$

The space of persistence diagrams, given by  $\mathcal{D}_p = \{\mathbf{D} : \text{pers}_p(\mathbf{D}) < \infty\}$ , is endowed with the family of  $p$ -Wasserstein metrics  $W_p$ . Given two persistence diagrams  $\mathbf{D}_1, \mathbf{D}_2 \in \mathcal{D}_p$ , the  $p$ -Wasserstein distance is given by

$$W_p(\mathbf{D}_1, \mathbf{D}_2) \doteq \left( \inf_{\gamma \in \Gamma} \sum_{\mathbf{z} \in \mathbf{D}_1 \cup \Delta} \|\mathbf{z} - \gamma(\mathbf{z})\|_\infty^p \right)^{\frac{1}{p}}, \quad (\text{B.1})$$

where  $\Gamma = \{\gamma : \mathbf{D}_1 \cup \Delta \rightarrow \mathbf{D}_2 \cup \Delta\}$  is the set of all bijections from  $\mathbf{D}_1$  to  $\mathbf{D}_2$  including the diagonal  $\Delta = \{(x, y) \in \mathbb{R}^2 : 0 \leq x = y \leq \infty\}$  with infinite multiplicity.

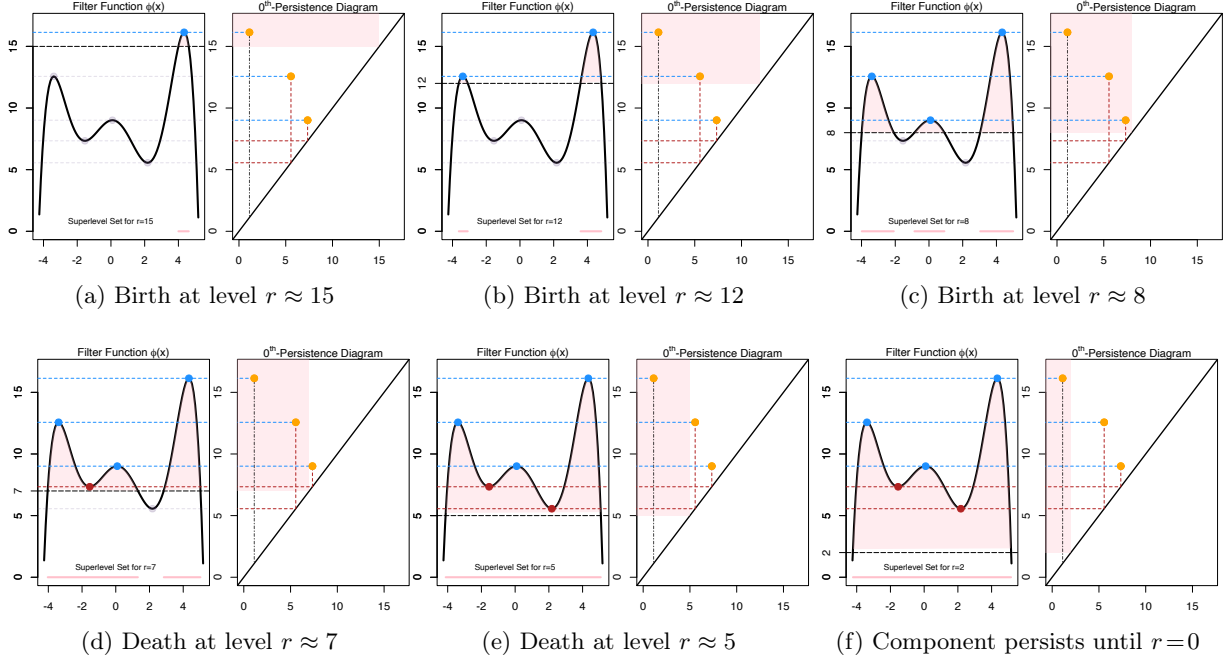


Figure 9: An example for the superlevel filtration of  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ . (a) As the superlevel set enters  $r \approx 15$ , the first connected component is born, corresponding to the blue dot on the highest peak of  $\phi$ . The superlevel set for  $r = 15$  is depicted in pink below. This is recorded as a birth in the corresponding orange dot enclosed in the pink shaded region of the persistence diagram. (b) As the  $r$  enters  $r \approx 12$ , another connected component is born. This is recorded as the second orange dot in the shaded region of the persistence diagram. (c) Again, at  $r \approx 8$ , a third connected component is born at the lowest peak of  $\phi$ . The three connected components in the superlevel set are shaded in pink below the function. The persistence diagram has three orange dots corresponding to these three connected components. (d) As  $r$  enters the first valley of  $\phi$ , depicted by the red dot, two connected components merge (i.e., one of the existing connected components die). By convention, the most recent persistent feature is merged into the older one, i.e., the connected component from (c) merges into the one from (b), and thus, it dies at this resolution. In the persistence diagram, this is noted by the fact that the orange dot born in (c) dies at resolution  $r \approx 7$ . At this stage, there are only two orange dots in the pink shaded region of the persistence diagram, corresponding to the two pink connected components in the superlevel set of  $\phi$ . (e) When  $r$  enters the second valley of  $\phi$ , the connected component from (b) merges into the connected component from (a), and form a single connected component. The orange dot in the persistence diagram records the death of this feature. (f) The single connected component persists from then on, and eventually dies at  $r = 0$ .