



**HAL**  
open science

## Functional Isolation Forest

Guillaume Staerman, Pavlo Mozharovskyi, Stéphan Cléménçon, Florence d'Alché-Buc

► **To cite this version:**

Guillaume Staerman, Pavlo Mozharovskyi, Stéphan Cléménçon, Florence d'Alché-Buc. Functional Isolation Forest. The 11th Asian Conference on Machine Learning 2019, Nov 2019, Nagoya, Japan. hal-02369435

**HAL Id: hal-02369435**

**<https://telecom-paris.hal.science/hal-02369435v1>**

Submitted on 19 Nov 2019

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

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

# Functional Isolation Forest

**Guillaume Staerman**

**Pavlo Mozharovskyi**

**Stephan Cléménçon**

**Florence d’Alché-Buc**

*LTCI, Télécom Paris, Institut Polytechnique de Paris*

GUILLAUME.STAERMAN@TELECOM-PARIS.FR

PAVLO.MOZHAROVSKYI@TELECOM-PARIS.FR

STEPHAN.CLEMENCON@TELECOM-PARIS.FR

FLORENCE.DALCHE@TELECOM-PARIS.FR

**Editors:** Wee Sun Lee and Taiji Suzuki

## Abstract

For the purpose of monitoring the behavior of complex infrastructures (*e.g.* aircrafts, transport or energy networks), high-rate sensors are deployed to capture multivariate data, generally unlabeled, in quasi continuous-time to detect quickly the occurrence of anomalies that may jeopardize the smooth operation of the system of interest. The statistical analysis of such massive data of functional nature raises many challenging methodological questions. The primary goal of this paper is to extend the popular ISOLATION FOREST (IF) approach to Anomaly Detection, originally dedicated to finite dimensional observations, to functional data. The major difficulty lies in the wide variety of topological structures that may equip a space of functions and the great variety of patterns that may characterize abnormal curves. We address the issue of (randomly) splitting the functional space in a flexible manner in order to isolate progressively any trajectory from the others, a key ingredient to the efficiency of the algorithm. Beyond a detailed description of the algorithm, computational complexity and stability issues are investigated at length. From the scoring function measuring the degree of abnormality of an observation provided by the proposed variant of the IF algorithm, a *Functional Statistical Depth* function is defined and discussed, as well as a multivariate functional extension. Numerical experiments provide strong empirical evidence of the accuracy of the extension proposed.

**Keywords:** Anomaly detection, functional data analysis, isolation forest, unsupervised learning

## 1. Introduction

The digital information boom, that goes hand in hand with the recent technological advances in data collection and management (*e.g.* IoT, distributed platforms), offers new perspectives in many areas of human activity (*e.g.* transportation, energy, health, commerce, insurance), and confronts these domains with major scientific challenges for exploiting these observations. The ever growing availability of massive data, often collected in quasi-real time, engendered high expectations, in particular the need of increased automation and computational efficiency, with the goal to design more and more ‘intelligent’ systems. In particular, modern high-rate sensors enabling the continuous observation of the behavior of complex systems pave the way for the design of efficient unsupervised machine-learning approaches to anomaly detection, that may find applications in various domains ranging from fraud surveillance to distributed fleet monitoring through predictive maintenance or

health monitoring of complex systems. However, although many unsupervised learning procedures for anomaly detection (AD in abbreviated form) have been proposed, analyzed and applied in a variety of practical situations (see, *e.g.*, [Chandola et al. \(2009\)](#)), the case of functional data, though of crucial importance in practice (refer to [Ramsay and Silverman \(2005\)](#); [Ferraty and Vieu \(2006\)](#) for an account of Functional Data Analysis) has received much less attention in the literature, the vast majority of methods that are documented in the literature being generally model-based. The main barrier to the design of nonparametric anomaly detection techniques tailored to the functional framework lies in the huge diversity of patterns that may carry the information that is relevant to discriminate between abnormal and normal observations, see [Hubert et al. \(2015\)](#).

It seems indeed far from straightforward to extend machine-learning methods for anomaly detection in the finite-dimensional case such as [Scott and Nowak \(2006\)](#); [Schölkopf et al. \(2001\)](#); [Steinwart et al. \(2005\)](#); [Vert and Vert \(2006\)](#); [Park et al. \(2010\)](#), unless preliminary *filtering techniques* are used. The *filtering approach* consists in projecting the functional data onto an adequate finite dimensional function subspace and using then the coefficients describing the latter to "feed" next some AD algorithm for multivariate data ([Ramsay and Silverman, 2005](#)). The basis functions are either selected through Principal Component Analysis (they correspond in this case to elements of the Karhunen-Loeve basis related to the process under study, supposedly of second order), or else are chosen among a dictionary of "time-frequency atoms" according to their capacity to represent efficiently the data. The representation a priori chosen, which can either enhance artificially certain accessory patterns or else make totally disappear some crucial features, critically determines performance of such an approach, the type of anomalies that can be recovered being essentially shaped by this choice.

The angle embraced in the present article is very different, the goal pursued being to extend the popular ISOLATION FOREST methodology ([Liu et al., 2008, 2012](#)) to the functional setup. This *ensemble learning* algorithm builds a collection of *isolation trees* based on a recursive and randomized tree-structured partitioning procedure. An isolation tree is a binary tree, representing a nested collection of partitions of the finite dimensional feature space, grown iteratively in a top-down fashion, where the cuts are axis perpendicular and random (uniformly, w.r.t. the splitting direction and the splitting value both at the same time). Incidentally, a variant referred to as EXTENDED ISOLATION FOREST ([Hariri et al., 2018](#)), has recently been proposed in the purpose of bias reduction: rather than randomly selecting a perpendicular split, a splitting direction is randomly chosen in the unit ball. An *anomaly score* is assigned to any observation, depending on the length of the path necessary to isolate it from the rest of the data points, the rationale behind this approach being that anomalies should be easier to isolate in a random manner than normal (in the sense of 'non-abnormal') data. Beyond obvious advantages regarding computational cost, scalability (*e.g.* isolation trees can be built from subsamples) and interpretability, the great flexibility offered by Isolation Forest regarding the splitting procedure called recursively makes it appealing when it comes to isolate (multivariate) functions/curves, possibly exhibiting a wide variety of geometrical shapes. It is precisely the goal of this paper to introduce a new generic algorithm, FUNCTIONAL ISOLATION FOREST (FIF) that generalizes (Extended) Isolation Forest to the infinite dimensional context. Avoiding dimensionality reduction steps, this extension is shown to preserve the assets of the original algorithm concerning computational

cost and interpretability. Its efficiency is supported by strong empirical evidence through a variety of numerical results.

The paper is organized as follows. Section 2 recalls the principles under the Isolation Forest algorithm for AD in the multivariate case and introduces the framework we consider for AD based on functional data. In Section 3, the extension to the functional case is presented and its properties are discussed at length. In Section 4, we study the behavior of the new algorithm and compare its performance to alternative methods standing as natural competitors in the functional setup through experiments. In Section 5, extension to multivariate functional data is considered, as well as relation to the data depth function and an application to the supervised classification setting. Eventually, several concluding remarks are collected in Section 6.

## 2. Background and Preliminaries

Here we briefly recall the Isolation Forest algorithm and its advantages (Section 2.1) and next introduce the framework for functional anomaly detection we consider throughout the paper (Section 2.2).

### 2.1. Isolation Forest

As a first go, we describe the Isolation Forest algorithm for AD in the multivariate context in a formalized manner for clarity's sake, as well as the Extended Isolation Forest version, see Liu et al. (2008, 2012) and Hariri et al. (2018) respectively. These two unsupervised algorithms can be viewed as *Ensemble Learning* methods insofar as they build a collection of binary trees and an anomaly scoring function based on the aggregation of the latter. Let  $\mathcal{S}_n = \{x_1, \dots, x_n\}$  be a training sample composed of  $n$  independent realizations of a generic random variable,  $X$ , that takes its value in a finite dimensional Euclidian space,  $\mathbb{R}^d$  say,  $X = (X^{(1)}, \dots, X^{(d)})$ .

An *isolation tree* (*itree* in abbreviated form)  $\mathcal{T}$  of depth  $J \geq 1$  is a proper binary tree that represents a nested sequence of partitions of the feature space  $\mathbb{R}^d$ . The root node corresponds to the whole space  $\mathcal{C}_{0,0} = \mathbb{R}^d$ , while any node of the tree, indexed by the pair  $(j, k)$  where  $j$  denotes the depth of the node with  $0 \leq j < J$  and  $k$ , the node index with  $0 \leq k \leq 2^j - 1$ , is associated to a subset  $\mathcal{C}_{j,k} \subset \mathbb{R}^d$ . A non terminal node  $(j, k)$  has two children, corresponding to disjoint subsets  $\mathcal{C}_{j+1,2k}$  and  $\mathcal{C}_{j+1,2k+1}$  such that  $\mathcal{C}_{j,k} = \mathcal{C}_{j+1,2k} \cup \mathcal{C}_{j+1,2k+1}$ . A node  $(j, k)$  is said to be terminal if it has no children.

Each *itree* is obtained by recursively filtering a subsample of training data of size  $\psi$  in a top-down fashion, by means of the following procedure. The dataset composed of the training observations present at a node  $(j, k)$  is denoted by  $\mathcal{S}_{j,k}$ . At iteration  $k + 2^j$  of the *itree* growing stage, a direction  $m$  in  $\{1, \dots, d\}$ , or equivalently a *split variable*  $X^{(m)}$ , is selected uniformly at random (and independently from the previous draws) as well as a *split value*  $\kappa$  in the interval  $[\min_{x \in \mathcal{S}_{j,k}} x^{(m)}, \max_{x \in \mathcal{S}_{j,k}} x^{(m)}]$  corresponding to the range of the projections of the points in  $\mathcal{S}_{j,k}$  onto the  $m$ -th axis. The children subsets are then defined by  $\mathcal{C}_{j+1,2k} = \mathcal{C}_{j,k} \cap \{x \in \mathbb{R}^d : x^{(m)} \leq \kappa\}$  and  $\mathcal{C}_{j+1,2k+1} = \mathcal{C}_{j,k} \cap \{x \in \mathbb{R}^d : x^{(m)} > \kappa\}$ , the children training datasets being defined as  $\mathcal{S}_{j+1,2k} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k}$  and  $\mathcal{S}_{j+1,2k+1} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k+1}$ .

An *itree*  $\mathcal{T}$  is thus built by iterating this procedure until all training data points are isolated (or the depth limit  $J$  set by the user is attained). A preliminary subsampling stage

can be performed in order to avoid swamping and masking effects, when the size of the dataset is too large. When it isolates any training data point, the *itree* contains exactly  $\psi - 1$  internal nodes and  $\psi$  terminal nodes. An *itree* constructed accordingly to a training subsample allows to assign to each training datapoint  $x_i$  a path length  $h_{\mathcal{T}}(x_i)$ , namely the depth at which it is isolated from the others, *i.e.* the number of edges  $x_i$  traverses from the root node to the terminal node that contains the sole training data  $x_i$ . More generally, it can be used to define an anomaly score for any point  $x \in \mathbb{R}^d$ .

**Anomaly Score prediction.** As the terminal nodes of the *itree*  $\mathcal{T}$  form a partition of the feature space, one may then define the piecewise constant function  $h_{\mathcal{T}} : \mathbb{R}^d \rightarrow \mathbb{N}$  by:  $\forall x \in \mathbb{R}^d$ ,

$$h_{\mathcal{T}}(x) = j \text{ if and only if } x \in \mathcal{C}_{j,k} \text{ and } (j, k) \text{ is a terminal node.}$$

This random path length is viewed as an indication for its degree of abnormality in a natural manner: ideally, the more abnormal the point  $x$ , the higher the probability that the quantity  $h_{\mathcal{T}}(x)$  is small. Hence, the algorithm above can be repeated  $N \geq 1$  times in order to produce a collection of *itrees*  $\mathcal{T}_1, \dots, \mathcal{T}_N$ , referred to as an *iforest*, that defines the scoring function

$$s_n(x) = 2^{-\frac{1}{Nc(\psi)} \sum_{l=1}^N h_{\tau_l}(x)}, \tag{1}$$

where  $c(\psi)$  is the average path length of unsuccessful searches in a binary search tree, see [Liu et al. \(2008\)](#) for further details.

**Extended Isolation Forest.** Observing that the geometry of the abnormal regions of the feature space is not necessarily well-described by perpendicular splits (*i.e.* by unions of hypercubes of the cartesian product  $\mathbb{R}^d$ ), a more flexible variant of the procedure recalled above has been proposed in [Hariri et al. \(2018\)](#), in the purpose of bias reduction. Rather than selecting a direction in  $\{1, \dots, d\}$ , one may choose a direction  $u \in \mathbb{S}_{d-1}$ , denoting by  $\mathbb{S}_{d-1}$  the unit sphere of the euclidian space  $\mathbb{R}^d$ . A node is then cut by choosing randomly and uniformly a threshold value in the range of the projections onto this direction of the training data points lying in the corresponding region. In the case where  $X$ 's ditribution has a density  $f(x)$  w.r.t. a  $\sigma$ -finite measure  $\lambda$  of reference, the goal of anomaly detection can be formulated as the recovery of sublevel sets  $\{x \in \mathbb{R}^d : f(x) \geq q\}$ ,  $q \geq 0$ , (under mild assumptions, they are *minimum volume sets* or *quantile regions*, see [Polonik \(1997\)](#); [Scott and Nowak \(2006\)](#), when measuring the volume by  $\lambda$ ), which may be not accurately approximated by unions of hyperrectangles (in the Gaussian situation for instance, such regions are the complementary sets of ellipsoïds,  $\lambda$  being Lebesgue measure on  $\mathbb{R}^d$ ).

## 2.2. Functional Data Analysis and Anomaly Detection

A functional random variable  $X$  is a r.v. that takes its values in a space of functions, see, *e.g.*, [Ferraty and Vieu \(2006\)](#). To be more specific, let  $I \subset \mathbb{R}_+$  be a time interval and consider a r.v. taking its values in the Hilbert space  $L_2(I)$  of real valued and square integrable (w.r.t. Lebesgue measure) functions  $x : I \rightarrow \mathbb{R}$ :

$$\begin{aligned} X & : \Omega \longrightarrow L_2(I) \\ \omega & \longmapsto X(\omega) = (X_t(\omega))_{t \in I} . \end{aligned}$$

Without any loss of generality, we restrict ourselves with functions defined on  $[0, 1]$  throughout the paper. In practice, only a finite dimensional marginal  $(X_{t_1}, \dots, X_{t_p})$ ,  $t_1 < \dots < t_p$ ,  $p \geq 1$  and  $(t_1, \dots, t_p) \in [0, 1]^p$  can be observed. However, considering  $(X_{t_1}, \dots, X_{t_p})$  as a discretized curve rather than a simple random vector of dimension  $p$  permits to take into account the dependence structure between the measurements over time, especially when the time points  $t_i$  are not equispaced. To come back to a function from discrete values, interpolation procedures or approximation schemes based on appropriate dictionaries can be used, combined with a preliminary smoothing step when the observations are noisy. From a statistical perspective, the analysis is based on a functional dataset  $\mathcal{S}_n = \{x_1, \dots, x_n\}$  composed of  $n \geq 1$  independent realizations of finite-dimensional marginals of the stochastic process  $X$ , that may be very heterogeneous in the sense that these marginals may correspond to different time points and be of different dimensionality. One may refer to Ramsay and Silverman (2005)'s book for a deep view on Functional Data Analysis (FDA in short). For simplicity, the functional data considered throughout the paper correspond to the observations of independent realizations of  $X$  at the same points.

In this particular context, *functional* anomaly detection aims at detecting the curves that significantly differ from the others among the dataset available. Given the richness of spaces of functions, the major difficulty lies in the huge diversity in the nature of the observed differences, which may not only depend on the locations of the curves. Following in the footsteps of Hubert et al. (2015), one may distinguish between three types of anomalies: *shift* (the observed curve has the same shape as the majority of the sample except that it is shifted away), *amplitude* or *shape* anomalies. All these three types of anomalies can be *isolated/transient* or *persistent*, depending on their duration with respect to that of the observations. One may easily admit that certain types of anomalies are harder to detect than others: for instance, an isolated anomaly in shape compared to an isolated anomaly in amplitude (*i.e.* change point). Although FDA has been the subject of much attention in recent years, very few generic and flexible methods tailored to functional anomaly detection are documented in the machine-learning literature to the best of our knowledge, except for specific types of anomalies (*e.g.* change-points).

In Statistics, although its applications are by no means restricted to AD, the concept of *functional depth* that allows to define a notion of centrality in the path space and a center-outward ordering of the curves of the functional dataset, see, *e.g.*, Cuevas et al. (2007); Claeskens et al. (2014); Hubert et al. (2015), has been used for this purpose. However, since the vast majority of functional depth functions introduced only describe the relative location properties of the sample curves, they generally fail to detect other types of anomalies. Another popular approach, usually referred to as *filtering*, consists in bringing the AD problem to the multivariate case by means of an adequate projection using *Functional Principal Component Analysis* (FPCA) (Ramsay and Silverman, 2005) or a preliminary selected basis of the function space considered (*e.g.* Fourier, wavelets) and apply next an AD algorithm designed for the finite-dimensional setup to the resulting representation. Such methods have obvious drawbacks. In FPCA, estimation of the Kahrnen-Loève basis can be very challenging and lead to loose approximations, jeopardizing next the AD stage, while the *a priori* representation offered by the 'atoms' of a predefined basis or frame may unsuccessfully capture the patterns carrying the relevant information to distinguish abnormal curves from the others. Another approach is based on the notion of *Minimum*

*Volume sets* (MV-sets in shortened version), originally introduced in Einmahl and Mason (1992) and that generalizes the concept of quantile for multivariate distributions and offers a nice nonparametric framework for anomaly detection in finite dimension, see Scott and Nowak (2006)’s work. Given the fact that no analogue of Lebesgue measure on an infinite-dimensional Banach space exists and since, considering a law  $\lambda$  of reference (*e.g.* the Wiener or a Poisson measure) on the function space  $\mathcal{H}$  of interest, the volume  $\lambda(\mathcal{C})$  of a measurable subset  $\mathcal{C} \subset \mathcal{H}$  can be hardly computed in general, it is far from straightforward to extend MV-set estimation to the functional setup.

The angle embraced in this paper is quite different. The direct approach we promote here is free from any preliminary representation stage and can be straightforwardly applied to a functional dataset. Precisely, in the subsequent section, we propose to extend the IF algorithm to the functional data framework, in a very flexible way, so as to deal with a wide variety of anomaly shapes.

### 3. Functional Isolation Forest

We consider the problem of learning a score function  $s : \mathcal{H} \rightarrow \mathbb{R}$  that reflects the degree of anomaly of elements in an infinite dimensional space  $\mathcal{H}$  w.r.t.  $P$ . By  $\mathcal{H}$ , we denote a functional Hilbert space equipped with a scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  such that any  $\mathbf{x} \in \mathcal{H}$  is a real function defined on  $[0, 1]$ . In the following, we describe in detail the proposed FUNCTIONAL ISOLATION FOREST (FIF) algorithm and discuss its properties.

#### 3.1. The FIF algorithm

A Functional Isolation Forest is a collection of *Functional Isolation Trees* (F-*itrees*) built from  $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , a training sample composed of independent realizations of a functional random variable,  $\mathbf{X}$ , that takes its values in  $\mathcal{H}$ . Given a functional observation  $\mathbf{x}$ , the score returned by FIF is a monotone transformation of the empirical mean of the path lengths  $h_{\tau_l}(\mathbf{x})$  computed by the F-*itrees*  $\mathcal{T}_l$ , for  $l = 1, \dots, N$  as defined in Eq. 1 in the multivariate case. While the general construction principle depicted in Section 2.1 remains the same for a F-*itree*, dealing with functional values raises the issue of finding an adequate feature space to represent various properties of a function. A function may be considered as abnormal according to various criteria of location and shape, and the features should permit to measure such properties. Therefore four ingredients have been introduced to handle functional data in a general and flexible way: (i) a set of candidate *Split variables* and (ii) a scalar product both devoted to function representation, (iii) a probability distribution to sample from this set and select a single *Split variable*, (iv) a probability distribution to select a *Split value*. The entire construction procedure of a F-*itree* is described in Figure 1.

**Function representation** To define the set of candidate *Split variables*, a direct extension of the original IF algorithm (Liu et al., 2008) would be to randomly draw an argument value (*e.g.* time), and use functional evaluations at this point to split a node, but this boils down to only rely on instantaneous observations of functional data to capture anomalies, which in practice will be usually interpolated. Drawing a direction on a unit sphere as in Hariri et al. (2018) is no longer possible due to the potentially excessive richness of  $\mathcal{H}$ . To circumvent

**Input:** A subsample  $\{\mathbf{x}_1, \dots, \mathbf{x}_\psi\}$ , a dictionary  $\mathcal{D}$ , a probability measure  $\nu$  and a scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ .

- (a) **Initialization:** The root node indexed by  $(0, 0)$  is associated with the whole input space  $\mathcal{C}_{0,0} = \mathcal{H}$ . The construction starts with the training dataset  $\mathcal{S}_{0,0} = \{\mathbf{x}_1, \dots, \mathbf{x}_\psi\}$  composed of  $n$  i.i.d. realizations of the random variable  $\mathbf{X}$ . Go to (b) with  $(j = 0, k = 0)$ .
- (b) **Stopping criterion:** Test if the node  $(j, k)$  is terminal: a node  $(j, k)$  is declared as *terminal* if the intersection between the current set  $\mathcal{C}_{j,k}$  and the current training set  $\mathcal{S}_{j,k}$  is reduced to a single data point or to a set of predefined cardinal. If the node is terminal, then stop the construction for this node, otherwise go to (c).
- (c) **Children node construction:** A non-terminal node  $(j, k)$  is split in three steps as follows:

1. Choose a *Split variable*  $\mathbf{d}$  according to the probability distribution  $\nu$  on  $\mathcal{D}$ .
2. Choose randomly and uniformly a *Split value*  $\kappa$  in the interval

$$\left[ \min_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{d} \rangle_{\mathcal{H}}, \max_{\mathbf{x} \in \mathcal{S}_{j,k}} \langle \mathbf{x}, \mathbf{d} \rangle_{\mathcal{H}} \right],$$

3. Form the children subsets

$$\begin{aligned} \mathcal{C}_{j+1,2k} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathcal{H} : \langle \mathbf{x}, \mathbf{d} \rangle_{\mathcal{H}} \leq \kappa\}, \\ \mathcal{C}_{j+1,2k+1} &= \mathcal{C}_{j,k} \cap \{\mathbf{x} \in \mathcal{H} : \langle \mathbf{x}, \mathbf{d} \rangle_{\mathcal{H}} > \kappa\}. \end{aligned}$$

as well as the children training datasets

$$\mathcal{S}_{j+1,2k} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k} \text{ and } \mathcal{S}_{j+1,2k+1} = \mathcal{S}_{j,k} \cap \mathcal{C}_{j+1,2k+1}.$$

- (d) **Recursion:** Apply the building procedure starting from (a) to nodes  $(j + 1, 2k)$  and  $(j + 1, 2k + 1)$

**Output:**  $(\mathcal{C}_{(0,0)}, \mathcal{C}_{(1,1)}, \dots)$

Figure 1: Construction procedure of a F-*tree*.

these difficulties, we propose to project the observations on elements of a dictionary  $\mathcal{D} \subset \mathcal{H}$  that is chosen to be rich enough to explore different properties of data and well appropriate to be sampled in a representative manner. More explicitly, given a function  $\mathbf{d} \in \mathcal{D}$ , the projection of a function  $\mathbf{x} \in \mathcal{H}$  on  $\mathcal{D}$ ,  $\langle \mathbf{x}, \mathbf{d} \rangle_{\mathcal{H}}$  defines a feature that partially describes  $\mathbf{x}$ . When considering all the functions of dictionary  $\mathcal{D}$ , one gets a set of candidate *Split variables* that provides a rich representation of function  $\mathbf{X}$ , depending on the nature of the dictionary. Dictionaries have been thoroughly studied in the signal processing community to achieve *sparse coding* of signals, see *e.g.* Mallat and Zhang (1993). They also provide a way to incorporate *a priori* information about the nature of the data, a property very useful in an industrial context in which functional data often come from the observation of a well



known device and thus can benefit from expert knowledge.

**Sampling a *Split variable*** Once a dictionary is chosen, a probability distribution  $\nu$  on  $\mathcal{D}$  is defined to draw a *Split variable*  $\mathbf{d}$ . Note that the choice of the sampling distribution  $\nu$  gives an additional flexibility to orientate the algorithm towards the search for specific properties of the functions.

**Sampling a *Split value*** Given a chosen *Split variable*  $\mathbf{d}$  and a current training dataset  $\mathcal{S}_{j,k}$ , a *Split value* is uniformly drawn in the real interval defined by the smallest and largest values of the projections on  $\mathbf{d}$  when considering the observations present in the node.

**Discussion on the dictionary** The choice of a suited dictionary plays a key role in construction of the FIF anomaly score. The dictionary can consist of deterministic functions, incorporate stochastic elements, contain the observations from  $\mathcal{S}$ , or be a mixture of several mentioned options. In *Computational Harmonic Analysis*, a wide variety of bases or frames, such as wavelets, ridgelets, cosine packets, brushlets and so on, have been developed in the last decades in order to represent efficiently/parsimoniously functions, signals or images exhibiting specific form of singularities (*e.g.* located at isolated points, along hyperplanes) and may provide massive dictionaries. The following ones will be used throughout the article: *mexican hat wavelet dictionary* (MHW), *Brownian motion dictionary* (B), *Brownian bridge dictionary* (BB), *cosine dictionary* (Cos), *uniform indicator dictionary* (UI), *dyadic indicator dictionary* (DI), and the *self-data dictionary* (Self) containing the dataset itself. See Section B and C of the Supplementary Materials for detailed definitions of these dictionaries and further discussion on them, respectively.

**Discussion on the scalar product** Besides the dictionary, the *scalar product* defined on  $\mathcal{H}$  brings some additional flexibility to measure different type of anomaly. While  $L_2$  scalar product allows for detection of *location anomalies*,  $L_2$  scalar product of derivatives (or slopes) would allow to detect anomalies regarding shape. This last type of anomalies can be challenging; *e.g.* Hubert et al. (2015) mention that *shape anomalies* are more difficult to detect, and Mosler and Mozharovskiy (2017) argue that one should consider both location and slope simultaneously for distinguishing complex curves. Beyond these two, a wide diversity of scalar products can be used, involving a variety of  $L_2$ -scalar products related to derivatives of certain orders, like in the definition of Banach spaces such as weighted Sobolev spaces, see Maz'ya (2011).

### 3.2. Ability of FIF to detect a variety of anomalies

As discussed in Section 2.2, most of state-of-the-art methods have a focus on a certain type of anomalies and are unable to detect various deviations from the normal behavior. The flexibility of the FIF algorithm allows for choosing the scope of the detection by selecting both the scalar product and the dictionary. Nevertheless, by choosing appropriate scalar product and dictionary, FIF is able to detect a great diversity of deviations from normal data. First, to account for both location and shape anomalies, we suggest the following

scalar product that provides a compromise between the both

$$\langle \mathbf{f}, \mathbf{g} \rangle := \alpha \times \frac{\langle \mathbf{f}, \mathbf{g} \rangle_{L_2}}{\|\mathbf{f}\| \|\mathbf{g}\|} + (1 - \alpha) \times \frac{\langle \mathbf{f}', \mathbf{g}' \rangle_{L_2}}{\|\mathbf{f}'\| \|\mathbf{g}'\|}, \quad \alpha \in [0, 1],$$

and illustrate its use right below. Thus, setting  $\alpha = 1$  yields the classical  $L_2$  scalar product,  $\alpha = 0$  corresponds to the  $L_2$  scalar product of derivative, and  $\alpha = 0.5$  is the Sobolev  $W_{1,2}$  scalar product. To illustrate the FIF's ability to detect a wide variety of anomalies at a time, we calculate the FIF anomaly scores with the Sobolev scalar product and the *gaussian wavelets dictionary* for a sample consisting of 105 curves defined as follows (inspired by Cuevas et al. (2007), see Fig. 2):

- 100 curves defined by  $\mathbf{x}(t) = 30(1 - t)^q t^q$  with  $q$  equispaced in  $[1, 1.4]$ ,
- 5 *abnormal* curves composed by one isolated anomaly  $\mathbf{x}_0(t) = 30(1 - t)^{1.2} t^{1.2}$  with a jump in  $t = 0.7$ , one magnitude anomaly  $\mathbf{x}_1(t) = 30(1 - t)^{1.6} t^{1.6}$  and three kind of shape anomalies  $\mathbf{x}_2(t) = 30(1 - t)^{1.2} t^{1.2} + \sin(2\pi t)$ ,  $\mathbf{x}_3(t) = 30(1 - t)^{1.2} t^{1.2}$  noised by  $\varepsilon \sim \mathcal{N}(0, 0.3^2)$  on the interval  $[0.2, 0.8]$  and  $\mathbf{x}_4(t) = 30(1 - t)^{1.2} t^{1.2} + \frac{1}{2} \sin(10\pi t)$ .

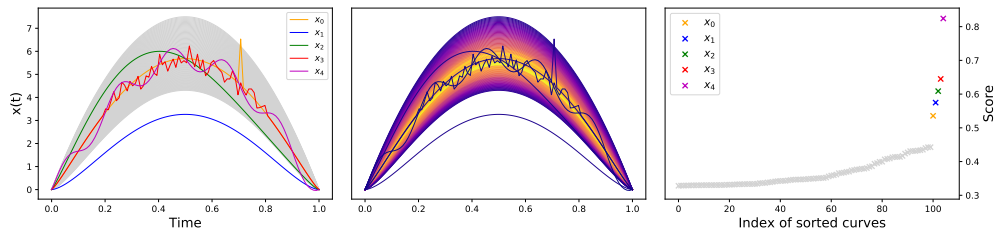


Figure 2: The simulated dataset with the five introduced anomalies (left). The scored dataset (middle), the darker the color, the more the curves are considered anomalies. The sorted anomaly score of the dataset (right).

One can see that the five anomalies, although very different, are all detected by FIF with a significantly different score.

#### 4. Numerical Results

In this section, we provide an empirical study of the proposed algorithm. First, in Section 4.1 we explore the stability and consistency of the score function w.r.t. the probability distribution of a r.v.  $\mathbf{X}$  and the sample size. Furthermore, we examine the influence of proposed dictionaries on the score function and bring performance comparisons with benchmark methods. Second, in Section 4.2, we benchmark the performance of FIF on several real labeled datasets by measuring its ability to recover an "abnormal" class on the test set. In all experiments,  $N$  the number of F-trees is fixed to 100 and the height limit is fixed to  $\lceil \log_2(\psi) \rceil$ .

### 4.1. Impact of the Hyperparameters on Stability

Since functional data are more complex than multivariate data, and the dictionary constitutes an additional source of variance, a question of stability of the FIF anomaly score estimates is of high interest. This issue is even more important because of the absence of theoretical developments due to their challenging nature. The empirical study is conducted on two simulated functional datasets presented in Fig. 3: Dataset (a) is the standard Brownian motion being a classical stochastic process widely used in the literature. Dataset (b) has been used by [Claeskens et al. \(2014\)](#) and has smooth paths. For each dataset, we choose/add four observations for which the FIF anomaly score is computed after training: a normal observation  $\mathbf{x}_0$ , two anomalies  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and a more extreme anomaly  $\mathbf{x}_3$ . We therefore expect the following ranking of the scores:  $s_n(\mathbf{x}_0) < s_n(\mathbf{x}_1) \leq s_n(\mathbf{x}_2) < s_n(\mathbf{x}_3)$ , for both datasets.

Further, we provide an illustration of the empirical convergence of the score. All other parameters being fixed, we increase the number of observations  $n$  when calculating the scores of the four selected observations; the empirical median and the boxplots of the scores computed over 100 random draws of the dataset are shown in Fig. 4. First, one observes score convergence and variance decrease in  $n$ . Further, let us take a closer look at the score tendencies on the example of  $\mathbf{x}_0$  and  $\mathbf{x}_3$ . The score of  $\mathbf{x}_3$  first increases (for dataset (a)) and slightly decreases (for dataset (b)) with growing  $n$  until  $n$  reaches  $\psi = 64$ , which happens because this abnormal observation is isolated quite fast (and thus has short path length) but the  $c(\psi)$  in the denominator of the exponent of (1) increases in  $\psi$ . For  $n > 64$ , the score of  $\mathbf{x}_3$  decreases in  $n$  since  $h_i(\mathbf{x}_3)$  overestimates the real path length of  $\mathbf{x}_3$  for subsamples in which it is absent; frequency of such subsamples grows in  $n$  and equals, *e.g.*, 0.872 for  $n = 500$ . On the other hand, this phenomenon allows to unmask grouped anomalies as mentioned in [Liu et al. \(2008\)](#). The behavior is reciprocal for the typical observation  $\mathbf{x}_0$ . Its FIF anomaly score starts by decreasing in  $n$  since  $\mathbf{x}_0$  tends to belong to the deepest branches of the trees and is always selected while  $\psi < n$ . For larger  $n$ , the path length of  $\mathbf{x}_0$  is underestimated for subsamples where it is absent when growing the tree, which explains slight increase in the score before it stabilizes. A second experiment illustrated in Fig. 5 is conducted to measure the impact of various dictionaries shortly cited in Section 3 and more thoroughly

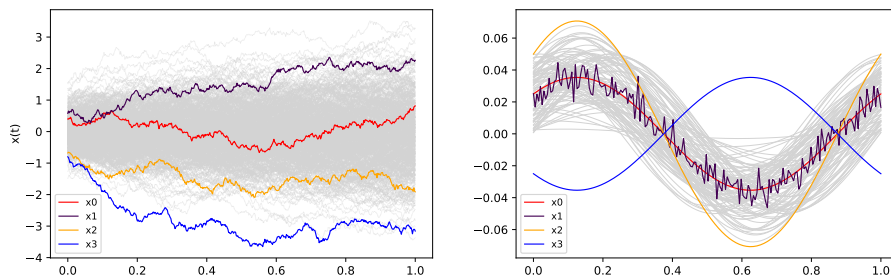


Figure 3: Datasets (a) (left) and (b) (right) containing, respectively, 500 and 200 functional paths with 4 selected observations.

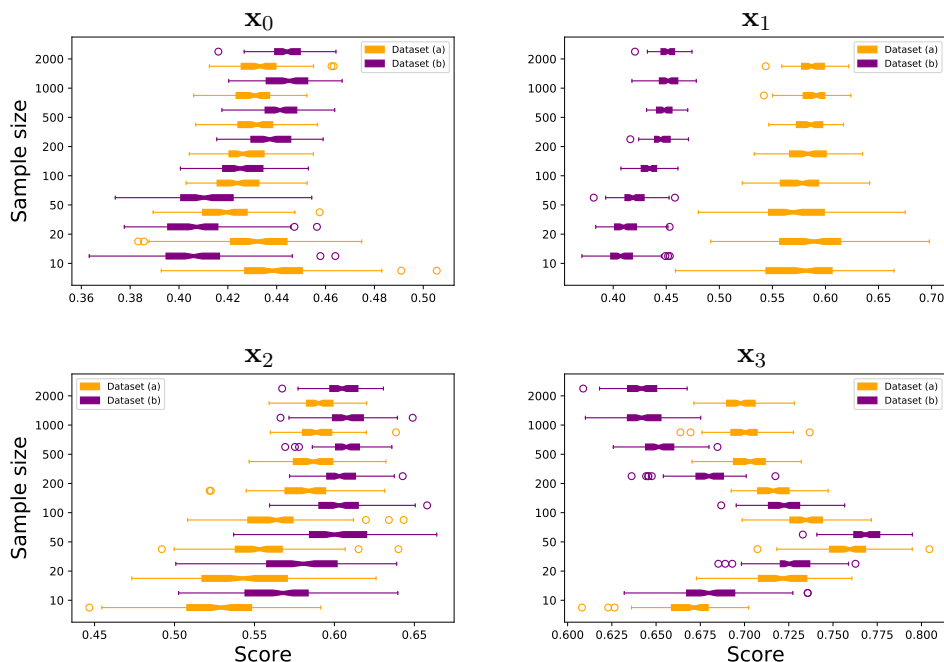


Figure 4: Boxplot (over 100 repetitions) of the FIF score for the observations  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$  for different sample sizes. The orange boxplots represent the dataset (a) while the purple boxplots represent the dataset (b).

described in Section B of the Supplementary Materials;  $L_2$  scalar product is used. One observes that the variance of the score seems to be mostly stable across dictionaries, for both datasets. Thus, random dictionaries like uniform indicator (UI) or Brownian motion (B) do not introduce additional variance into the FIF score. Since we know the expected ranking of the scores, we can observe that FIF relying on the Self, UI, and dyadic indicator (DI) dictionaries fail to make a strong difference between  $\mathbf{x}_0$  and  $\mathbf{x}_1$ . Since  $\mathbf{x}_1$  differs only slightly in the amplitude from the general pattern, these dictionaries seem insufficient to capture this fine dissimilarity: while Self and DI dictionaries simply do not contain enough elements, UI dictionary is too simple to capture this difference (it shares this feature with DI dictionary). For the scalar product  $L_2$  on derivatives (see Fig. 18 in the Supplementary Materials), distinguishing anomalies for the Brownian motion becomes difficult since they differ mainly in location, while for a sine function the scores resemble those with the usual  $L_2$  scalar product. Thus, even though—as seen in Section 3.2—capturing different types of anomalies is one of the general strengths of the FIF algorithm, the dictionary may still have an impact on detection of functional anomalies in particular cases.

More experiments were run regarding the stability of the algorithm, but for sake of space, we describe them in Section C of the Supplementary Materials.

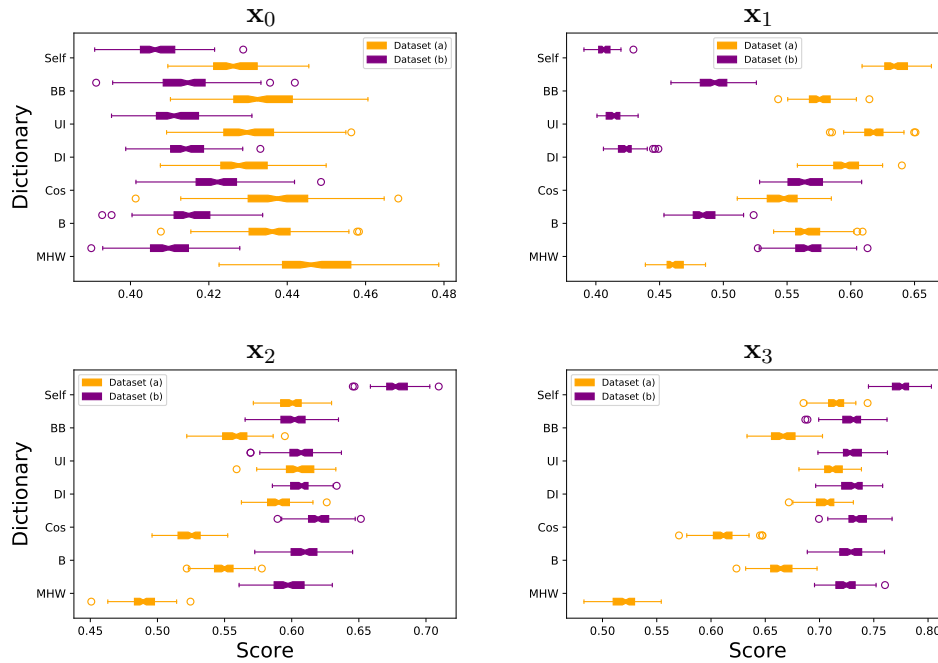


Figure 5: Boxplot (over 100 repetitions) of the FIF score for the observations  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$  for different dictionaries using the  $L_2$  scalar product. The orange boxplots represent the dataset (a) while the purple boxplots represent the dataset (b).

## 4.2. Real Data Benchmarking

To explore the performance of the proposed FIF algorithm, we conduct a comparative study using 13 classification datasets from the UCR repository (Chen et al., 2015). We consider the larger class as normal and some of others as anomalies (see Table 1 for details). When classes are balanced, *i.e.* for 9 datasets out of 13, we keep only part of the anomaly class to reduce its size, always taking the same observations (at the beginning of the table) for a fair comparison. Since the datasets are already split into train/test sets, we use the train part (without labels) to build the FIF and compute the score on the test set. We assess the performance of the algorithm by measuring an Area Under the Receiver Operation Characteristic curve (AUC) on the test set. Both train and test sets are rarely used during learning in unsupervised setting since labels are unavailable when fitting the model. Thus, when fitting the models on unlabeled training data, good performances on the test set show a good generalization power.

**Competitors** FIF is considered with two finite size dictionaries *dyadic indicator*, *the self-data* and the infinite size dictionary *cosines* (with  $\alpha = 1$  and  $\alpha = 0$ ); its parameters are set  $N = 100$ ,  $\psi = \min(256, n)$  and the height limit to  $= \lceil \log_2(\psi) \rceil$ . We contrast the FIF method with three most used multivariate anomaly detection techniques and two functional depths, with default settings. The multivariate methods—*isolation forest* (IF) (Liu et al., 2008), *local outlier factor* (LOF) (Breunig et al., 2000), and *one-class support vector machine*

(OCSVM) (Schölkopf et al., 2001)— are employed after dimension reduction by Functional PCA keeping 20 principal components with largest eigenvalues after a preliminary step of filtering using Haar basis. The depths are the random projection halfspace depth (Cuevas et al., 2007) and the functional Stahel-Donoho outlyingness (Hubert et al., 2015).

**Analysis of the results** Taking into account the complexity of the functional data, as expected there is no method performing generally best. Nevertheless, FIF performs well in most of the cases, giving best results for 10 datasets and second best for 6 datasets. It is worth to mention that the dictionary plays an important role in identifying anomalies, while FIF seems to be rather robust w.r.t. other parameters: The “CinECGTorso” dataset contains anomalies differing in location shift which are captured by the cosine dictionary. Dyadic indicator dictionary allows to detect local anomalies in “TwoLeadECG” and “Yoga” datasets. Self-data dictionary seems suited for Datasets “SonyRobotAI2” and “StarlightCurves” whose challenge is to cope with many different types of anomalies.

	p	training : $n_a/n$	testing : $n_a/n$	normal lab	anomaly lab
Chinatown	24	4 / 14 (29%)	95 / 345	2	1
Coffee	286	5 / 19 (26%)	6 / 19	1	0
ECGFiveDays	136	2 / 16 (12%)	53 / 481	1	2
ECG200	96	31 / 100 (31%)	36 / 100	1	-1
Handoutlines	2709	362 / 1000 (36 %)	133 / 370	1	0
SonyRobotAI1	70	6 / 20 (30 %)	343 / 601	2	1
SonyRobotAI2	65	4 / 20 (20 %)	365 / 953	2	1
StarLightCurves	1024	100 / 673 (15 %)	3482 / 8236	3	1 and 2
TwoLeadECG	82	2 / 14 (14 %)	570 / 1139	1	2
Yoga	426	10 / 173 ( 06 %)	1393 / 3000	2	1
EOGHorizontal	1250	10 / 40 (25 %)	30 / 61	5	6
CinECGTorso	1639	4 / 16 (25 %)	345 / 688	3	4
ECG5000	140	31 / 323 (10 %)	283 / 2910	1	3,4 and 5

Table 1: Datasets considered in performance comparison:  $n$  is the number of instances,  $n_a$  is the number of anomalies.  $p$  is the number of discretization points.

Methods :	DI $_{L_2}$	Cos $_{Sob}$	Cos $_{L_2}$	Self $_{L_2}$	IF	LOF	OCSVM	fHD $_{RP}$	fSDO
Chinatown	<b>0.93</b>	0.82	0.74	0.77	0.69	0.68	0.70	0.76	<b>0.98</b>
Coffee	0.76	<b>0.87</b>	0.73	<b>0.77</b>	0.60	0.51	0.59	0.74	0.67
ECGFiveDays	0.78	0.75	0.81	0.56	0.81	<b>0.89</b>	<b>0.90</b>	0.60	0.81
ECG200	0.86	<b>0.88</b>	<b>0.88</b>	0.87	0.80	0.80	0.79	0.85	0.86
Handoutlines	0.73	<b>0.76</b>	0.73	0.72	0.68	0.61	0.71	0.73	<b>0.76</b>
SonyRobotAI1	<b>0.89</b>	0.80	0.85	0.83	0.79	0.69	0.74	0.83	<b>0.94</b>
SonyRobotAI2	0.77	0.75	0.79	<b>0.92</b>	<b>0.86</b>	0.78	0.80	<b>0.86</b>	0.81
StarLightCurves	0.82	0.81	0.76	<b>0.86</b>	0.76	0.72	0.77	0.77	<b>0.85</b>
TwoLeadECG	<b>0.71</b>	0.61	0.61	0.56	<b>0.71</b>	0.63	<b>0.71</b>	0.65	0.69
Yoga	<b>0.62</b>	0.54	0.60	0.58	0.57	0.52	<b>0.59</b>	0.55	0.55
EOGHorizontal	0.72	<b>0.76</b>	<b>0.81</b>	0.74	0.70	0.69	0.74	0.73	0.75
CinECGTorso	0.70	<b>0.92</b>	<b>0.86</b>	0.43	0.51	0.46	0.41	0.64	0.80
ECG5000	0.93	<b>0.98</b>	<b>0.98</b>	0.95	0.96	0.93	0.95	0.91	0.93

Table 2: AUC of different anomaly detection methods calculated on the test set. Bold numbers correspond to the best result while italics to the second best.

## 5. Extensions of FIF

**Extension to multivariate functions** FIF can be easily extended to the multivariate functional data, *i.e.* when the quantity of interest lies in  $\mathbb{R}^d$  for each moment of time. For this, the coordinate-wise sum of the  $d$  corresponding scalar products is used to project the data onto a chosen dictionary element:  $\langle \mathbf{f}, \mathbf{g} \rangle_{\mathcal{H}^{\otimes d}} := \sum_{i=1}^d \langle f^{(i)}, g^{(i)} \rangle_{\mathcal{H}}$ . The dictionary is then defined in  $(\mathcal{H}([0, 1]))^{\otimes d}$ , *e.g.*, by componentwise application of one or several univariate dictionaries, see Section 3. In the Supplementary Materials we give an illustration of multivariate functional anomaly detection on the MNIST (Lecun et al., 1998) dataset, each digit being seen as a 2D-curve.

**Connection to data depth** Regarding FIF score as an anomaly ranking yields a connection to the notion of the *statistical depth function* (see Mosler (2013) for an overview), which has been successfully applied in outlier detection (see, *e.g.*, Hubert et al. (2015)). Statistical data depth has been introduced as a measure of centrality (or depth) of an arbitrary observation  $\mathbf{x} \in (\mathcal{H}([0, 1]))^{\otimes d}$  with respect to the data at hand  $\mathcal{S}$ . A data depth measure based on FIF score can be defined for (multivariate) functional data as:  $D_{FIF}(\mathbf{x}; \mathcal{S}) = 1 - s_n(\mathbf{x}; \mathcal{S})$ . Data depth proves to be a useful tool for a low-dimensional data representation called *depth-based map*. Using this property, Li et al. (2012) and Mosler and Mozharovskiy (2017) define a *DD*-plot classifier which consists in applying a multivariate classifier to the depth-based map. Low-dimensional representation is of particular interest for functional data and a *DD*-plot classifier can be defined using the FIF-based data depth. Let  $\mathcal{S}^{trn} = \mathcal{S}^1 \cup \dots \cup \mathcal{S}^q$  be a training set for supervised classification containing  $q$  classes, each subset  $\mathcal{S}^j$  standing for class  $j$ . The depth map is defined as follows:

$$\mathbf{x} \mapsto \phi(\mathbf{x}) = (D_{FIF}(\mathbf{x}; \mathcal{S}^1), \dots, D_{FIF}(\mathbf{x}; \mathcal{S}^q)) \in [0, 1]^q.$$

As an illustration, we apply the depth map to 3 digits (1, 5 and 7, 100 observations per digit for training and 100 testing) of the MNIST dataset after their transformation to two-variate functions using `skimage` python library (see Figure 6). One observes appealing geometrical interpretation (observe, *e.g.*, the location of the abnormally distant—from their corresponding classes—observations) and a clear separation of the classes. To illustrate separability, we apply linear multiclass (one-against-all) SVM in the depth space, which delivers the accuracy of 99% on the test data.

## 6. Conclusion

The Functional Isolation Forest algorithm has been proposed, which is an extension of Isolation Forest to functional data. The combined choice of the dictionary itself, the probability distribution used to pick a *Split variable* and the scalar product used for the projection enables FIF to exhibit a great flexibility in detecting anomalies for a variety of tasks. FIF is extendable to multivariate functional data. When transformed in a data depth definition, FIF can be used for supervised classification via a low-dimensional representation—the depth space. The open-source implementation of the method, along with all reproducing scripts, can be accessed at <https://github.com/Gstaerman/FIF>.

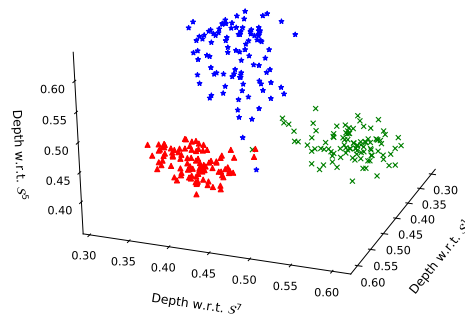


Figure 6: Depth space embedding of the three digits (1, 5 and 7) of the MNIST dataset.

## References

- M.M. Breunig, H.-P. Kriegel, R.T. Ng, and J. Sander. LOF: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, volume 29, pages 93–104. ACM, 2000.
- V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3):15:1–15:58, 2009.
- Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista. The UCR time series classification archive, July 2015. URL [www.cs.ucr.edu/~eamonn/time\\_series\\_data/](http://www.cs.ucr.edu/~eamonn/time_series_data/).
- G. Claeskens, M. Hubert, L. Slaets, and K. Vakili. Multivariate functional halfspace depth. *Journal of American Statistical Association*, 109(505):411–423, 2014.
- A. Cuevas, M. Febrero, and R. Fraiman. Robust estimation and classification for functional data via projection-based depth notions. *Computational Statistics*, 22(3):481–496, 2007.
- J.H.J. Einmahl and D.M. Mason. Generalized quantile processes. *The Annals of Statistics*, 20(2):1062–1078, 1992.
- F. Ferraty and P. Vieu. *Nonparametric Functional Data Analysis*. Springer-Verlag, New York, 2006.
- S. Hariri, M. Carrasco Kind, and R. J. Brunner. Extended isolation forest. *ArXiv e-prints*, 2018. URL <https://arxiv.org/abs/1811.02141>.
- M. Hubert, P.J. Rousseeuw, and P. Segaeert. Multivariate functional outlier detection. *Statistical Methods & Applications*, 24(2):177–202, 2015.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- J. Li, J.A. Cuesta-Albertos, and R.Y. Liu. DD-classifier: Nonparametric classification procedure based on DD-plot. *Journal of the American Statistical Association*, 107(498):737–753, 2012.



- F. T. Liu, K. M. Ting, and Z. Zhou. Isolation forest. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pages 413–422. IEEE Computer Society, 2008.
- F. T. Liu, K. M. Ting, and Z. Zhou. Isolation-based anomaly detection. In *ACM Transactions on Knowledge Discovery from Data (TKDD)*, volume 6, pages 1–39, 2012.
- S.G. Mallat and Z. Zhang. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on signal processing*, 41(12):3397–3415, 1993.
- V. Maz'ya. *Sobolev Spaces: with Applications to Elliptic Partial Differential Equations*. Springer-Verlag, Berlin Heidelberg, 2011.
- K. Mosler. Depth statistics. In *Robustness and Complex Data Structures: Festschrift in Honour of Ursula Gather*, pages 17–34. Springer, Berlin Heidelberg, 2013.
- K. Mosler and P. Mozharovskiy. Fast DD-classification of functional data. *Statistical Papers*, 58(4):1055–1089, 2017.
- C. Park, J.Z. Huang, and Y. Ding. A computable plug-in estimator of minimum volume sets for novelty detection. *Operations Research*, 58(5):1469–1480, 2010.
- W. Polonik. Minimum volume sets and generalized quantile processes. *Stochastic Processes and their Applications*, 69(1):1–24, 1997.
- J.O. Ramsay and B.W. Silverman. *Functional Data Analysis*. Springer-Verlag, New-York, 2005.
- B. Schölkopf, J.C. Platt, J. Shawe-Taylor, A. Smola, and R. Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, 2001.
- C. Scott and R. Nowak. Learning minimum volume sets. *Journal of Machine Learning Research*, 7:665–704, 2006.
- I. Steinwart, D. Hush, and C. Scovel. A classification framework for anomaly detection. *Journal of Machine Learning Research*, 6:211–232, 2005.
- R. Vert and J.-P. Vert. Consistency and convergence rates of one-class SVMs and related algorithms. *Journal of Machine Learning Research*, 7:817–854, 2006.