*Review Article*

# Intrinsic Dimension Estimation: Relevant Techniques and a Benchmark Framework

## P. Campadelli,[1] E. Casiraghi,[1] C. Ceruti,[1] and A. Rozza[2]

[1]*Dipartimento di Informatica, Università degli Studi di Milano, Via Comelico 39, 20135 Milano, Italy*
[2]*Research Group, Hyera Software, Via Mattei 2, Coccaglio, 25030 Brescia, Italy*

Correspondence should be addressed to E. Casiraghi; casiraghi@di.unimi.it

When dealing with datasets comprising high-dimensional points, it is usually advantageous to discover some data structure. A fundamental information needed to this aim is the minimum number of parameters required to describe the data while minimizing the information loss. This number, usually called intrinsic dimension, can be interpreted as the dimension of the manifold from which the input data are supposed to be drawn. Due to its usefulness in many theoretical and practical problems, in the last decades the concept of intrinsic dimension has gained considerable attention in the scientific community, motivating the large number of intrinsic dimensionality estimators proposed in the literature. However, the problem is still open since most techniques cannot efficiently deal with datasets drawn from manifolds of high intrinsic dimension and nonlinearly embedded in higher dimensional spaces. This paper surveys some of the most interesting, widespread used, and advanced state-of-the-art methodologies. Unfortunately, since no benchmark database exists in this research field, an objective comparison among different techniques is not possible. Consequently, we suggest a benchmark framework and apply it to comparatively evaluate relevant state-of-the-art estimators.

## 1. Introduction

Since the 1950s, the rapid pace of technological advances allows us to measure and record increasing amounts of data, motivating the urgent need to develop dimensionality reduction systems to be applied on real datasets comprising high-dimensional points.

To this aim, a fundamental information is provided by the *intrinsic dimension* (id) defined by Bennett [1] as the minimum number of parameters needed to generate a data description by maintaining the "intrinsic" structure characterizing the dataset, so that the information loss is minimized.

More recently, a quite intuitive definition employed by several authors in the past has been reported by Bishop in [2], p. 314, where the author writes that "a set in $D$ dimensions is said to have an id equal to $d$ if the data lies entirely within a $d$-dimensional subspace of $\mathfrak{R}^D$."

Though more specific and different id definitions have been proposed in different research fields [3–5], throughout the pattern recognition literature the presently prevailing id definition views a point set as a sample set uniformly drawn from an unknown smooth (or locally smooth) manifold structure, eventually embedded in a higher dimensional space through a nonlinear smooth mapping; in this case, the id to be estimated is the manifold's *topological dimension*.

Due to the importance of id in several theoretical and practical application fields, in the last two decades a great deal of research effort has been devoted to the development of effective id estimators. Though several techniques have been proposed in the literature, the problem is still open for the following main reasons.

At first, it must be highlighted that though Lebesgue's definition of topological dimension [6] (see Section 3.2) is quite clear, in practice its estimation is difficult if only a finite set of points is available. Therefore, id estimation techniques proposed in the literature are either founded on different notions of dimension (e.g., fractal dimensions, Section 3.2.1) approximating the topological one or on various techniques aimed at preserving the characteristics of data-neighborhood distributions, which reflect the topology of the underlying manifold. Besides, the estimated id value markedly changes as the scale used to analyze the input dataset changes [7] (see
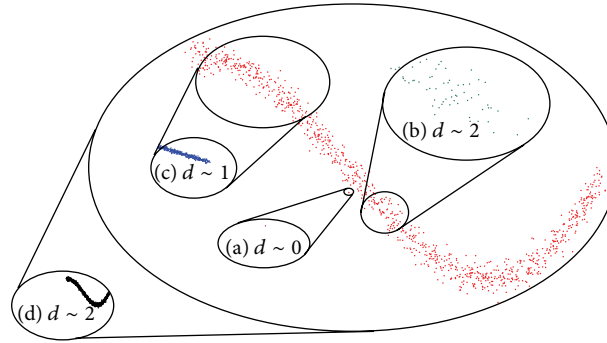
FIGURE 1: At very small scales (a) the dataset seems zero-dimensional; in this example, when the resolution is increased until including all the dataset (b) the id looks larger and seems to equal the embedding space dimension; the same effect happens when it is estimated at noise level (d); the correct id estimate is obtained at an intermediate resolution.

an example in Figure 1), and with the number of available points being practically limited, several methods underestimate id when its value is sufficiently high (namely, id ⩾ 10). Other serious problems arise when the dataset is embedded in higher dimensional spaces through a nonlinear map. Finally, the too high computational complexity of most estimators makes them unpractical when the need is to process datasets comprising huge amounts of high-dimensional data.

In this work, after recalling the application domains of interest, we survey some of the most interesting, widespread used, and advanced id estimators. Unfortunately, since each method has been evaluated on different datasets, it is difficult to compare them by solely analyzing the results reported by the authors. This highlights the need of a benchmark framework, such as the one proposed in this work, to objectively assess and compare different techniques in terms of robustness with respect to parameter settings, high-dimensional datasets, datasets characterized by a high id, and noisy datasets.

The paper is organized as follows: in Section 2 the usefulness of the id knowledge is motivated and interesting id application domains profitably exploiting it are recalled; in Section 3 we survey notable state-of-the-art id estimators, grouping them according to the employed methods; in Section 4 we summarize mostly used experimental settings, we propose a benchmark framework, and we employ it to objectively assess and compare relevant id estimators; in Section 5 conclusions and open research problems are reported.

## 2. Application Domains

In this section we motivate the increasing research interest aimed at the development of automatic id estimators, and we recall different application contexts where the knowledge of the id of the available input datasets is a profitable information.

In the field of pattern recognition, the id is one of the first and fundamental pieces of information required by several dimensionality reduction techniques [8–12], which try to represent the data in a more compact, but still informative, way to reduce the "curse of dimensionality" effects [13].

Furthermore, when using an autoassociative neural network to perform a nonlinear feature extraction, the id value $d$ can suggest a reasonable value for the number of hidden neurons [14]. Indeed, a network with a single hidden layer of neurons with linear activation functions has an error function with a unique global minimum and, at this minimum, the network performs a projection on the subspace spanned by the first $d$ principal components [15] estimated on the dataset (see 8.6.2 of [2]), with $d$ being the number of hidden neurons. Besides, according to statistical learning theory [16], the capacity and generalization capability of a given classifier may depend on the id. More specifically, in the particular case of linear classifiers where the data are drawn from a manifold embedded through an identical map, the Vapnik-Chervonenkis (VC) dimension of the separation hyperplane is $d + 1$ (see [16], pp. 156–158). Since the generalization error depends on the VC dimension, it follows that the generalization capability may depend on the id value $d$. Moreover, in [17] the authors mark that, in order to balance a classifier generalization ability and its empirical error, the complexity of the classification model should also be related to the id of the available dataset. Furthermore, since complex objects can be considered as structures composed by multiple manifolds that must be clustered to be processed separately, the knowledge of the local ids characterizing the considered object is fundamental to obtain a proper clustering [18].

These observations motivate applications employing global or local id estimates to discover some structure within the data. In the following we summarize or simply recall some interesting examples [19–25].

In [19] the authors introduce a fractal dimension estimator, called correlation dimension (CD) estimator (see Section 3.2.1), and show that the id estimate it computes is a reliable approximation of the strange attractor dimension in chaotic dynamical systems.

In the field of gene expression analysis, the work proposed in [20] shows that the id estimate computed by the nearest neighbor estimator (described in [26] and Section 3.2.2) is a lower bound for the number of genes to be used in supervised and unsupervised class separation of cancer and other diseases. This information is important since generally used datasets contain large number of genes and the classification

results strongly depend on the number of genes employed to learn the separation criteria.

In [21], the authors show that `id` estimation methods being derived from the basis theory of fractal dimensions ([7, 19, 27], see Section 3.2.1) can be successfully used to evaluate the model order in signals and time series, which is the number of past samples required to model the time series adequately and is crucial to make reliable predictions. This comparative work employs fractal dimension estimators, since the domain of attraction of nonlinear dynamic systems has a very complex geometric structure, which could be captured by closely related studies on fractal geometry and fractal dimensions.

A noteworthy research work in the field of crystallography [22] employs the fractal CD estimator [19] followed by a correction method [28] that, according to the authors, "is needed because the CD estimator, to give correct estimations of the `id`, requires an unrealistically large number of points." Anyway, the experimental results show that `id` is a useful information to be exploited when analyzing crystal structures. This study not only proves that `id` estimates are especially useful when dealing with practical tasks concerning real data, but also underlines the need to compute reliable estimates on datasets drawn from manifolds characterized by high `id` and embedded in spaces of much greater dimensionality.

The work of Carter et al. [23] is very interesting and notable because it is one of the first considering that the input data might be drawn from a multimanifold structure, where each submanifold has a (possibly) different `id`. To separate the manifolds, the authors compute local `id` estimates, by applying both a fractal dimension estimator (namely, MLE [27]; see Section 3.2.1) and a nearest neighbor-based estimator (described in [29, 30]; see Section 3.2.2) on properly defined data neighborhoods. The authors then show that the computed local `ids` might be helpful for the following interesting applications: (1) "Debiasing global `id` estimates": the negative bias caused both by the limited number of available sample points and by the *curse of dimensionality* is reduced by computing global `id` estimates through a weighted average of the local ones, which assign greater importance to the points away from the boundaries. However the authors themselves note that this method is only applicable for data with a relatively low `id`, since in high dimensions the points lie nearby the boundaries [31]. (2) "Statistical manifold learning": the local `id` estimates are used to reduce the dimension of statistical manifolds [32], that is, manifolds whose points represent a `pdf`. When this step is applied as the first step of document classification applications, and analysis of patients' samples acquired in the field of flow cytometry, it allows us to obtain lower dimensional points showing a good class separation. (3) "Network anomaly detection": considering that the overall complexity of a router network is decreased when few sources account for a disproportionate amount of traffic, a decrease in the `id` of the entire network is searched for. (4) "Clustering": problems of data clustering and image segmentation are dealt with by assuming that different clusters and image patches belong to manifold structures characterized by different complexity (and `ids`).

In [24], to the aim of analyzing gene expression time series, the authors compute `id` estimates by comparing the fractal CD estimator and the nearest neighbor (NN) estimator [26]. The results on both simulated and real data show that NN seems to be more robust than CD with respect to nonlinear embedding and the underlying time-series model.

In the field of geophysical signal processing, hyperspectral images, whose pixels represent spectra generated by the combination of an unknown set of independent contributions, called endmembers, often require estimating the number of endmembers. To this aim, the proposal in [25] is to substitute state-of-the-art algorithms specifically designed to solve this task with `id` estimators. After motivating the idea by describing the relation between the `id` of a dataset and the number of endmembers, the authors choose to experiment two fractal `id` estimators [7, 19] and a nearest neighbor-based one [33]. They obtain the most reliable results with the latest one after opportunely tuning the number of nearest neighbors to be considered.

Finally, other noteworthy examples of research works that profitably exploit `id` and estimate it by usually applying fractal dimension estimators concern financial time series prediction [34], biomedical signal analysis [35–37], analysis of ecological time series [38], radar clutter identification [39], speech analysis [40], data mining and low dimensional representation of (biomedical) time series [41], and plant traits representation [42].

## 3. Intrinsic Dimension Estimators

In this section we survey some of the most notable, recent, and effective state-of-the-art `id` estimators, grouping them according to the main ideas they are based on.

Specifically, in Section 3.1 we describe *projective* `id` *estimators*, which basically process a dataset $\mathbf{P}_N \equiv \{\mathbf{p}_i\}_{i=1}^N \subseteq \Re^D$ to identify a somehow appealing lower dimensional subspace where to project the data; the space dimension of the identified subspace is viewed as the `id` estimate.

More recent projective `id` estimators exploit the assumption of datasets $\mathbf{P}_N \equiv \{\mathbf{p}_i\}_{i=1}^N \subseteq \Re^D$ being uniformly drawn from a smooth (or locally smooth) manifold $\mathscr{M} \subseteq \Re^d$, embedded into a higher $D$-dimensional space through a nonlinear map; this is also the basic assumption of all the other groups of methods that will be referred to as *topological* `id` *estimators* (see Section 3.2) and *graph-based* `id` *estimators* (see Section 3.3).

We note that the taxonomy we are using to group the reviewed methods is different from the one, commonly used by several authors in the past (as an example, see [43]), that viewed methods as *global*, when `id` estimation is performed by considering a dataset as a whole, or *local*, when all the data neighborhoods are analyzed separately and an estimate is computed by combining all the local results. All the recent methods have abandoned the global approach since it is now clear that analyzing a dataset at its biggest scale cannot produce reliable results. They thus estimate the global `id` by somehow combining local `ids`. This way of proceeding comes

from the assumption that the underlying manifold is locally smooth.

*3.1. Projective* id *Estimators.* The first projective id estimators introduced in the literature explicitly compute the mapping that projects input points $\mathbf{P}_N \in \mathfrak{R}^D$ to the subspace $\mathcal{M} \subseteq \mathfrak{R}^d$ minimizing the information loss [27, 43] and therefore view the id as the minimal number of vectors linearly spanning the subspace $\mathcal{M}$. It must be noted that, since these methods were originally designed for exploratory data analysis and dimensionality reduction, they often require the dimensionality of $\mathcal{M}$ (the id to be estimated) to be provided as input parameter. However, their extensions and variants include methodologies to automatically estimate id.

Most of the projective id estimators can be grouped into two main categories: projection techniques based on Multidimensional Scaling (MDS) [44, 45] or its variants, which tend to preserve as much as possible pairwise distances among the data, and projection techniques based on Principal Component Analysis (PCA) [15, 46] and its variants that search for the best projection subspace $\mathcal{M}$ minimizing the projection error.

Some of the best known examples of MDS algorithms are MDSCAL [47–52], Bennett's algorithm [1, 53], Sammon's mapping [54], Curvilinear Component Analysis (CCA) [55], ISOMAP [56], and Local Linear Embedding (LLE) [57]. As shown by experiments reported in [56, 57] ISOMAP and variants of LLE compute the most reliable id estimates. We believe that their better performance is due to the fact that both ISOMAP and LLE have been the first projective methods based on the assumption that the input points are drawn from an underlying manifold, whose curvature might affect the precision of data neighborhoods computed by employing the Euclidean distance. However, as noted in [58, 59], these algorithms have shown that they suffer of all the major drawbacks affecting MDS-based algorithms, which are too much tied by the preservation of the pairwise distance values. Besides, as highlighted in [30], ISOMAP as an id estimator, as well as other spectral based methods like PCA, relies on a specific estimated eigenstructure that may not exist in real data. Regarding LLE, it either requires the id value to be known in advance or may automatically estimate it by analyzing the eigenvalues of the data neighborhoods [60]; however, as outlined in [7, 27], id estimates computed by means of eigenvalue analysis are as unreliable as those computed by most PCA-based approaches. Moreover, in [46] it is noted that methods such as LLE are based on the solution of a sparse eigenvalue problem under the unit covariance constraint; however, due to this imposed constraint, the global shape of the embedded data cannot reflect the underlying manifold.

PCA [15, 46] is one of the most cited and well-known projective id estimators, often used as the first step of several pattern recognition problems, to compute low dimensional representations of the available datasets. When PCA is used for id estimation, the estimate is the number of "most relevant" eigenvectors of the sample covariance matrix, also called principal components (PCs). Due to the promising dimensionality reduction results, several PCA-based approaches, both deterministic and probabilistic, have been published.

Among deterministic approaches, we recall the Kernel PCA (KPCA) [61] and the local PCA (LPCA) [62] and its extensions to automatically select the number of PCs [63, 64]. We observe that the work presented in [64] is one of the first works that estimates id by considering an underlying topological structure and therefore applies LPCA on data neighborhoods represented by an Optimally Topology Preserving Map (OTPM) built on clustered data (given an input dataset $\mathbf{P}_N$, its OTPM is usually computed through Topology Representing Networks (TRNs); these are unsupervised neural networks [65] developed to map $\mathbf{P}_N$ to a set of neurons whose learnt connections define proximities in $\mathbf{P}_N$. These proximities correspond to the optimal topology preserving Voronoi tessellation and the corresponding Delaunay triangulation. In other words, TRNs compute connectivity structures that define and perfectly preserve the topology originally present in the data, forming a discrete path-preserving representation of the inner (topological) structure of the topological manifold underlying the dataset $\mathbf{P}_N$). The authors of this method state that their approach is more efficient and less sensitive to noise with respect to the PCA-based approaches. However, they do not show any experimental comparison and, besides, their algorithm employs critical thresholds and a data clustering technique whose result heavily influences the precision of the computed estimate [27].

The usage of a probabilistic approach has been firstly introduced by Tipping and Bishop in [66]. Considering that deterministic methodologies lack an associated probabilistic model for the observed data, their Probabilistic PCA (PPCA) reformulates PCA as the maximum likelihood solution of a specific latent variable model. PPCA and its extensions to both mixture and hierarchical mixture models have been successfully applied to several real problems, but they still provide an id-estimation mechanism depending on critical thresholds. This motivates its subsequent variants [67] and developments, whose examples are Bayesian PCA (BPCA) [68] and two Bayesian model order selection methods introduced in [69, 70]. In [71] the asymptotic consistency of id estimation by (constrained) isotropic version of PPCA is shown with numerical experiments on simulated and real datasets.

While the aforementioned methods have been simply recalled since their id estimation results have shown to be unreliable [7, 27], in the following recent and promising proposals are described with more details.

The Simple Exponential Family PCA (SePCA) [72] has been developed to overcome the assumption of Gaussian-distributed data that makes it difficult to handle all types of practical observations, for example, integers and binary values. SePCA achieves promising results by using exponential family distributions; however, it is highly influenced by critical parameter settings and it is successful only if the data distribution is known, which is often not the case, specially when highly nonlinear manifold structures must be treated.

In [73] the authors propose the Sparse Probability PCA (SPPCA) as a probabilistic version of the Sparse PCA (SPCA) [74]. Precisely, SPCA selects id by forcing the sparsity of the projection matrix that is the matrix containing the PCs. However, based on the consideration that the level of sparsity is not automatically determined by SPCA, SPPCA

employs a Bayesian formulation of SPCA, achieving sparsity by employing a different prior and automatically learning the hyperparameter related to the constraint weight through Evidence Approximation ([75] Section 3.5). The authors' results and also the results of the comparative evaluation proposed in [76] show that this method seems to be less affected by the problems of the aforementioned projective schemes.

An alternative method (MLSVD) [77] applies Singular Value Decomposition (SVD), basically a variant of PCA, locally and in a multiscale fashion to estimate the id characterizing $D$-dimensional datasets drawn from nonlinearly embedded $d$-dimensional manifolds $\mathcal{M}$ corrupted by Gaussian noise. Precisely, exploiting the same ideas of the theoretical PCA-based id estimator presented in [78], the authors note that the best way to avoid the effects of the curvature (induced by the nonlinearity of the embedding) is to apply SVD locally, that is, in hyperspheres $\mathcal{B}(\mathbf{p}, r)$ centered on the data points $\mathbf{p}$ and having radius $r$. However, the choice of $r$ is constrained by the following considerations: (1) $r$ must be big enough to have at least $k \geq d$ neighbors, (2) $r$ must be small enough to ensure that $\mathcal{M} \cap \mathcal{B}$ is linear (or at least smooth), and (3) $r$ must be big enough to ensure that the effects of noise are negligible. When these three constraints are met, the tangent space $T_{\mathcal{M}}^d(\mathbf{p}, r)$, computed by applying SVD on the $k$ neighbors, is a good approximation of the tangent space of $\mathcal{M} \cap \mathcal{B}$ and the number of its relevant eigenvalues corresponds to the (local) id of $\mathcal{M}$. To find a proper value for $r$, the authors propose a multiscale approach that applies SVD on neighborhoods $\mathcal{B}(\mathbf{p}, r_s)$ whose radius varies in a range $r_s \in \{r_L, \ldots, r_H\}$. This allows us to compute $D$ scale-dependent, local singular values $\lambda_1(\mathbf{p}, r_s) \geq \cdots \geq \lambda_D(\mathbf{p}, r_s)$; using a least squares fitting procedure the SVs can be expressed as functions of $r$ whose analysis allows us to identify the range of scales $[r_{\min}, \ldots, r_{\max}]$ not influenced by either noise or curvature. Finally, in the range $r_s = [r_{\min}, \ldots, r_{\max}]$ the squared SVs are analyzed to get the id estimate $\hat{d}$ that maximizes the gap $\Delta(j) = \lambda_j(\mathbf{p}, r_s) - \lambda_{j+1}(\mathbf{p}, r_s)$ for the largest range of $r_s$. The proposed algorithm has been evaluated on unit $d$-dimensional hyperspheres and cubes embedded in $\mathfrak{R}^{100}$ and affected by Gaussian noise. The reported results are very good, while other ten well-known methods [19, 23, 27, 30, 79–81] show that the ids estimated on the same datasets are unreliable also in the absence of noise.

### 3.2. Topological Approaches.

Topological approaches for id estimation consider a manifold $\mathcal{M} \subseteq \mathfrak{R}^d$ embedded in a higher dimensional space $\mathfrak{R}^D$ through a proper (locally) smooth map $\phi : \mathcal{M} \rightarrow \mathfrak{R}^D$ and assume that the given dataset is $\mathbf{P}_N = \{\mathbf{p}_i\}_{i=1}^N = \{\phi(\mathbf{x}_i)\}_{i=1}^N \subset \mathfrak{R}^D$, where $\mathbf{x}_i$ are independent identically distributed (i.i.d.) points drawn from $\mathcal{M}$ through a smooth probability density function (pdf) $f : \mathcal{M} \rightarrow \mathfrak{R}^+$.

Under this assumption the id to be estimated is the manifold's topological dimension, defined either through the firstly proposed Brouwer Large Inductive Dimension [82] or the equivalent Lebesgue Covering Dimension [83]. Since Brouwer's definition has been soon neglected by mathematicians for its difficult proof [83], the commonly adopted topological dimension definition is Lebesgue's Covering Dimension, reported in the following.

*Definition 1* (cover). Given a topological space $\mathcal{X}$, a cover of a set $\mathcal{Y} \subseteq \mathcal{X}$ is a countable collection $\mathcal{C} = \{\mathcal{C}_i\}$ of open sets such that each $\mathcal{C}_i \subset \mathcal{X}$ and $\bigcup_i \mathcal{C}_i \supseteq \mathcal{Y}$.

*Definition 2* (refinement of a cover). A refinement of a cover $\mathcal{C}$ of a set $\mathcal{Y}$ is another cover $\mathcal{C}'$ such that each set in $\mathcal{C}'$ is contained in some sets of $\mathcal{C}$.

*Definition 3* (topological dimension (Lebesgue Covering Dimension)). Given the aforementioned definitions, the topological dimension of the topological space $\mathcal{X}$, also called Lebesgue Covering Dimension, is $d$ if every finite cover of $\mathcal{X}$ admits a refinement $\mathcal{C}'$ such that no subset of $\mathcal{X}$ has more than $d + 1$ intersecting open sets in $\mathcal{C}'$. If no such minimal integer value exists, $\mathcal{X}$ is said to be of infinite topological dimension.

To our knowledge, at the state of the art only two estimators have been explicitly designed to estimate the topological dimension.

One of them, the Tensor Voting Framework (TVF) [84] and its variants [85], relies on the usage of an iterative information diffusion process based on Gestalt principles of perceptual organization [86]. TVF iteratively diffuses local information describing, for each $\mathbf{p}_i \in \mathbf{P}_N$, the tangent space approximating the underlying neighborhood of $\mathcal{M}$. To this aim, the information diffused at each iteration is second-order symmetric positive definite tensors whose eigenvectors span the local tangent space. Practically, during the initialization step a ball tensor $\mathbf{T}_i^0$, which is an identity matrix representing the absence of orientation, is used to initialize a token $p_i$ for each point $\mathbf{p}_i$ as $\{p_i = (\mathbf{p}_i, \mathbf{T}_i^0)\}_{i=1}^N$. During iteration $t$ each token $p_i$ "generates" the set of tensors $\{\mathbf{T}_{i,j}^t\}_{j \neq i}$ that enact as votes cast to neighboring tokens; precisely, $\mathbf{T}_{i,j}^t$ is sent to the $j$th neighbor, and it encodes information related to the local tangent space estimate in $\mathbf{p}_i$ at time $t$ and decays as the curvature and the distance from the $j$th neighbor increase. On the other side, at iteration $t$ each token $p_j$ receives votes that are summed to update the $p_j$'s tensor as $\mathbf{T}_j^{t+1} = \sum_{i \neq j} \mathbf{T}_{i,j}^t$; this essentially refines the estimate of the local tangent space in $\mathbf{p}_j$. Based on the definition of topological dimension provided by Brouwer [82], in [87] it is noted that TVF can be employed to estimate the local ids by identifying the number of most relevant eigenvalues of the computed second-order tensors. Although interesting, this method has a too high computational cost, which makes it unfeasible for spaces of dimension $D \geq 4$.

From the definition of Lebesgue Covering Dimension it can be derived [88] that the topological dimension of any $\mathcal{M} \subseteq \mathfrak{R}^d$ coincides with the affine dimension $d$ of a finite simplicial complex (a simplicial complex in $\mathfrak{R}^d$ has affine dimension $d$ if it is a collection of affine simplexes in $\mathfrak{R}^d$, having at most dimension $d$ or having at most $d + 1$ vertices) covering $\mathcal{M}$. This essentially means that a $d$-dimensional manifold could be approximated by a collection

of $d$-dimensional simplexes (each having at most $d + 1$ vertices); therefore, the topological dimension of $\mathcal{M}$ could be practically estimated by analyzing the number of vertices of the collection of simplexes estimated on $\mathbf{P}_N$. To this aim, in [89] a method is proposed to find the number of relevant positive coefficients that are needed to reconstruct each $\mathbf{p}_i \in \mathbf{P}_N$ from a linear combination of its $k$ neighbors, where $k$ is a parameter to be manually set in the range $d < k \leq D + 1$. This algorithm is based on the fact that neighbors with positive reconstruction coefficients are the vertices of a simplex with dimension equal to the topological dimension of $\mathcal{M}$. Practically, to ensure that $k > d$, its value is set to $D$, the reconstruction coefficients are calculated by means of an optimization problem constrained to be nonnegative, and the coefficients bigger than a user-defined threshold are considered as the relevant ones. The id estimate is then computed by employing two alternative approaches: the first one simply computes the mode of the number of relevant coefficients for each neighborhood and the second one sorts in descending order the coefficients computed for each neighborhood, computes the mean $\bar{\mathbf{c}}$ of the sorted coefficients, and estimates id as the number of relevant values in $\bar{\mathbf{c}}$. Note that, since $k > d$, this method is strongly affected by the curvature of the manifold when the id is big enough. Indeed, the results of the reported experimental evaluation make the authors assert that the method works well only on noisy-free data of low id (id $\leq 6$), under the assumption that the sampling process is uniform and the data points are sufficient.

Though interesting, both approaches have shown to be effective only for manifolds of low curvature as well as low id values.

In the following we survey other id estimators employing two different estimation approaches, which allow us to categorize them. More precisely, in Section 3.2.1 *fractal* id *estimators* are described, which estimate different fractal dimensions since they are good approximations of the topological one; in Section 3.2.2 *nearest neighbors-based* (NN) id *estimators* are recalled, which are often based on the statistical analysis of the distribution of points within small neighborhoods.

*3.2.1. Fractal* id *Estimators.* Since topological dimension cannot be practically estimated, several authors implicitly assume that $\mathcal{M}$ has a somehow fractal structure (see [90] for an exhaustive description of fractal sets) and estimate id by employing fractal dimension estimators, the most relevant of which are surveyed in this section.

Very roughly, since the basis concept of all fractal dimensions is that the volume of a $d$-dimensional ball of radius $r$ scales with its size $s$ as $r^d$ [90, 91], all fractal dimension estimators are based on the idea of counting the number of observations in a neighborhood of radius $r$ to (somehow) estimate the rate of growth of this number. If the estimated growth is $r^d$, then the estimated fractal dimension of the data is considered to be equal to $d$.

Note that all the derived estimators have the fundamental limitation that, in order to get an accurate estimation, the size $N$ of the dataset with id $d$ has to satisfy the inequality proved by Eckmann and Ruelle in [92] for the correlation dimension estimator (CD [19], see below):

$$d < \frac{2}{\log(1/\rho)} * \log N,$$

$$\text{being } \rho = \frac{r}{D} \ll 1, \quad \frac{1}{2}N^2\left(\frac{r}{D}\right)^d \gg 1. \tag{1}$$

This will lead to a large value of $N$, even for a dataset with lower id.

Among fractal dimension estimators, one of the most cited algorithms is presented in [19] and will be referred to as CD in the following. It is an estimator of the correlation dimension ($\dim_{\text{Corr}}$), whose formal definition is as follows.

*Definition 4* (correlation dimension). Given a finite sample set $\mathbf{P}_N$, let

$$C_N(r) = \frac{2}{N(N-1)} \sum_{i=1, i<j}^{N} I\left(r - \|\mathbf{p}_i - \mathbf{p}_j\|\right), \tag{2}$$

where $\|\cdot\|$ is the Euclidean norm and $I(\cdot)$ is the step function used to simulate a closed ball of radius $r$ centered on each $\mathbf{p}_i$ ($I(y) = 0$ if $y < 0$, and $I(y) = 1$ otherwise). Then, for a countable set, the correlation dimension $\dim_{\text{Corr}}$ is defined as

$$\dim_{\text{Corr}} = \lim_{r \to 0} \lim_{N \to \infty} \frac{\log C_N(r)}{\log r}. \tag{3}$$

In practice CD computes an estimate, $\hat{d}$, of $\dim_{\text{Corr}}$ by computing $C_N(r)$ for different $r_i$ and applying least squares to fit a line through the points $(\log r_i; \log C_N(r_i))$. It has to be noted that, to produce correct id estimates, this estimator needs a very large number of data points [22], which is never available for practical applications; however, the computed unreliable estimations can be corrected by the correction method proposed in [28].

The relevance of the CD estimator is shown by its several variants and extensions. An example is the work proposed in [91], where the authors propose a normalized CD estimator for binary data and achieve estimates approximating those computed by CD.

Since CD is heavily influenced by the setting of the scale parameters, in [93] the authors estimate the id by computing the expectation value of $\dim_{\text{Corr}}$ through maximum likelihood estimate of the distribution of distances among points. The estimated $\hat{d}$ is computed as

$$\hat{d} = -\left(\frac{1}{|Q|} \sum_{k=1}^{|Q|} r_k\right)^{-1}, \tag{4}$$

where $Q$ is the set $Q = \{r_k \mid r_k < r\}$ and $r_k$ is the Euclidean distance between two generic data points and $r$ is a real value, called cut-off radius.

To develop an estimator more efficient than CD, in [94] the authors choose a different notion of Fd, namely, the Information Dimension $\dim_I$:

$$\dim_I = -\lim_{\delta \to 0} \frac{\sum_{i=1}^{\mathcal{N}(\delta)} pr_i(\log pr_i)}{\log \delta}, \tag{5}$$

where $\mathcal{N}(\delta)$ is the minimum number of $\delta$-sized hypercubes covering a topological space and $pr_i$ is the probability of finding a point in the $i$th hypercube. Noting that when the scale $\delta$ in (5) is big enough the different coverings used to estimate $\dim_I$ could produce different values for $\mathcal{N}(\delta)$, the authors look for the covering composed by the minimum number $\mathcal{N}_{\min}(\delta)$ of nonempty sets. Similar to the CD algorithm, the id is the average slope of the curve obtained by fitting the points with coordinates $(\log \delta; \sum_{i=1}^{\mathcal{N}_{\min}(\delta)} pr_i \log pr_i)$.

This algorithm is compared with the CD estimator, and the experimental tests show that both methods compute the same estimates. However, the achieved computation time is much lower than that of CD.

Considering that CD can severely underestimate the topological dimension if the data distribution on the manifold is nearly nonuniform, in [7] the author proposes the Packing Number (PN), a fractal dimension estimator that approximates the Capacity Dimension ($\dim_{\mathrm{Cap}}$). This choice is motivated by the fact that $\dim_{\mathrm{Cap}}$ does not depend on the data distribution on the manifold and, if both $\dim_{\mathrm{Cap}}$ and the topological dimension exist (which is certainly the case in machine learning applications), the two dimensions agree. To formally define $\dim_{\mathrm{Cap}}$, the $\epsilon$-covering number $\mathcal{N}(\epsilon)$ of a set $\mathcal{S} \subset \mathcal{X}$ must be defined; $\mathcal{N}(\epsilon)$ is the minimum number of open balls $\mathcal{B}(\mathbf{x}_0, \epsilon) = \{\mathbf{x} \in \mathcal{X} : \|\mathbf{x}_0 - \mathbf{x}\| < \epsilon\}$ whose union is a covering of $\mathcal{S}$, where $\|\cdot\|$ is a distance metric. The definition of $\dim_{\mathrm{Cap}}$ of $\mathcal{S} \subset \mathcal{X}$ is based on the observation that the covering number $\mathcal{N}(\epsilon)$ of a $d$-dimensional set is proportional to $\epsilon^{-d}$:

$$\dim_{\mathrm{Cap}} = -\lim_{\epsilon \to 0} \frac{\log \mathcal{N}(\epsilon)}{\log \epsilon}. \quad (6)$$

Since the estimation of $\mathcal{N}(\epsilon)$ is computationally expensive, based on the relation $\mathcal{N}(\epsilon) \leq \mathcal{N}_{\mathrm{Pack}}(\epsilon) \leq \mathcal{N}(\epsilon/2)$, the authors employ the $\epsilon$-Packing Number $\mathcal{N}_{\mathrm{Pack}}(\epsilon)$, defined in [95] as the maximum cardinality of an $\epsilon$-separated set. Employing a greedy algorithm to compute $\mathcal{N}_{\mathrm{Pack}}(\epsilon)$, the estimate, $\hat{d}$, of $\dim_{\mathrm{Cap}}$ is computed as

$$\hat{d}(\epsilon_1, \epsilon_2) = -\frac{\log \mathcal{N}_{\mathrm{Pack}}(\epsilon_1) - \log \mathcal{N}_{\mathrm{Pack}}(\epsilon_2)}{\log \epsilon_1 - \log \epsilon_2}. \quad (7)$$

To estimate $\hat{d}$ a greedy algorithm is used; however, as noted by the author, the dependency of $\hat{d}$ with respect to the order in which the points are visited by the greedy algorithm introduces a high variance. To avoid this problem, the algorithm iterates the procedure $M$ times on random permutations of the data and considers the average as the final id estimate. The comparative evaluation with the CD estimator makes the authors assert that PN "seems more reliable if data contains noise or the distribution on the manifold is not uniform." Unfortunately, also this method is scale-dependent.

To avoid any scale-dependency in [79] the authors propose an estimator (Hein) based on the asymptotes of a smoothed version of (2), obtained by replacing the step

function $I(\cdot)$ with a suitable kernel function. Precisely, they define

$$U(N, h, d) = \frac{2}{N(N-1)} \sum_{1 \leq i < j \leq N} \frac{1}{h^d} K_h\left(\frac{\|\mathbf{p}_i - \mathbf{p}_j\|}{h^2}\right), \quad (8)$$

where $K_h$ is a kernel function with bandwidth $h$ and $d$ is the assumed dimensionality of the manifold from which the points are sampled. Note that, to guarantee the converge of (8), the bandwidth $h$ has to fulfill the constraint $\lim_{N \to \infty}(Nh^d) = \infty$. For this reason the authors formalize $h$ as a function of $N$ and, to achieve scale-independency, propose a method that estimates the id by analyzing the convergence of $U(N, h, d)$ when varying the parameters $N$ and $d$. Precisely, the dataset is subsampled to create sets of different cardinalities $n_{\mathrm{sub}} \in \mathcal{N}_{\mathrm{sub}} = \{N, N/2, N/3, N/4, N/5\}$ and the $D$ curves whose points have coordinates $(U(n_{\mathrm{sub}}, h(n_{\mathrm{sub}}), d), n_{\mathrm{sub}})$ are considered. Employing this information the following id estimator is proposed:

$$\mathrm{Slope}(d) = \max_{n_{\mathrm{sub}} \in \mathcal{N}_{\mathrm{sub}}} \left| \frac{\partial U(n_{\mathrm{sub}}, h(n_{\mathrm{sub}}), d)}{\partial n} \right|,$$

$$\hat{d} = \arg\min_{d \in \{1,\dots,D\}} \mathrm{Slope}(d). \quad (9)$$

This work is notable since the empirical tests are performed on synthetic datasets specifically designed to study the influence of high curvature as well as noise on the proposed estimator. The usefulness of these datasets is confirmed by the fact that they have been also employed to assess several subsequent methods [76, 96].

In [97] the authors present a fractal dimension estimator derived by the analysis of a vector quantizer applied to datasets $\mathbf{P}_N \subseteq \mathfrak{R}^D$. Considering the codebook $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\} \subset \mathfrak{R}^D$ containing $k$ code-vectors $\mathbf{y}_i$, a $k$-point quantizer is defined by a measurable function $Q_k : \mathfrak{R}^D \to \mathcal{Y}$, which brings each data point to one of the code-vectors in $\mathcal{Y}$. This partitions the dataset into $k$ so-called quantizer cells $\mathcal{S}_i = \{\mathbf{p}_i \in \mathbf{P}_N : Q_k(\mathbf{p}_i) = \mathbf{y}_i\}$, where $\log_2(k)$ is called the rate of the quantizer. Being $\mathbf{X}$ a random vector distributed according to a probability distribution $\nu$, the quantization error is $e_r(Q_k \mid \nu) = (E_\nu[\|\mathbf{X} - Q_k(\mathbf{X})\|^r])^{1/r}$, where $r \in [1, \infty)$ and $\|\cdot\|$ is the Euclidean norm in $\mathfrak{R}^D$. Given the set $\mathcal{Q}_k$ of all $D$-dimensional $k$-point quantizers, the performance achieved by an optimal $k$-point quantizer on $\mathbf{X}$ is $e_r^*(Q_k \mid \nu) = \inf_{Q_k \in \mathcal{Q}_k}(e_r(Q_k \mid \nu))$. When the quantizer rate is high, the quantizer cells can be well approximated by $D$-dimensional hyperspheres with radius equal to $\epsilon$ and centered on each code-vector $\mathbf{y}_i \in \mathcal{Y}$. In this case, the regularity of $\nu$ ensures that the probability of such balls is proportional to $\epsilon^{1/d}$, and it can be shown [98] that $e_r^*(Q_k \mid \nu) \approx k^{-1/d}$. This is referred to as the high-rate approximation and motivates the definition of quantization dimension of order $r$:

$$d_r(\nu) = -\lim_{k \to \infty} \frac{\log k}{\log e_r^*(k \mid \nu)}. \quad (10)$$

The theory of high-rate quantization [98] confirms that, for a regular $\nu$ supported on the manifold $\mathcal{M}$, $d_r(\nu)$ exists for each $1 \leq r \leq \infty$ and equals the intrinsic dimension of $\mathcal{M}$. Furthermore, the limit $k \rightarrow \infty$ allows us to motivate the relation between the quantization dimension and the Capacity Dimension. Indeed, according to the theory of high-rate quantization [98, 99], there exists a decreasing sequence $\{\epsilon_k\}$, such that for sufficiently large values of $k$ (i.e., in the high-rate regime that is when $k \rightarrow \infty$) the ratio $-(\log k)/(\log e_r^*(k \mid \nu))$ can be approximated increasingly finely, both from below and from above, by quantities converging to the common value $\dim_{\text{Cap}}$. To practically compute an estimate of the quantization dimension, having fixed the value of $r$, the authors select a range $k_1 \leq k \leq k_2$ of codebook sizes and design a set of quantizers $\{Q_k\}_{k=k_1}^{k_2}$ giving good approximations $\widehat{e}_r(k \mid \nu)$ of $e_r^*(k \mid p)$ over the chosen range of $k$. An id estimate is obtained by fitting the points with coordinates $(\log(k); -\log \widehat{e}_r(k \mid \nu))$ and measuring the average slope over the chosen range $k$. Though the authors mention that their algorithm is less affected by underestimation biases than neighborhood-based methods (see Section 3.2.2), in [23] this statement is confuted with theoretical arguments.

*3.2.2. Nearest Neighbors-Based* id *Estimators.* In this section we consider estimators, referred to as *NN* estimators in the following, that describe data-neighborhoods' distributions as functions of $d$. They usually assume that close points are uniformly drawn from small $d$-dimensional balls (hyperspheres) $\mathscr{B}_d(\mathbf{x}, r)$ having radius (a small radius $r \rightarrow 0 \in \mathfrak{R}^+$ guarantees that samples included into $\mathscr{B}_d(\mathbf{x}, r)$, being less influenced by the curvature induced by the map $\phi$, are approximating well enough the intrinsic structure of the underlying portion of $\mathcal{M}$) $r \rightarrow 0 \in \mathfrak{R}^+$ and being centered on $\mathbf{x} \in \mathcal{M}$.

Practically, given an input dataset $\mathbf{P}_N$, the value of functions $f(d)$ is computed by approximating the sampling process related to $\mathscr{B}_d$ through the $k$-nearest neighbor algorithm (kNN).

Among *NN* id estimators, Trunk's method [100] is often cited as one of the first methods. It formulates the distribution function, $f(d)$, with an ad hoc statistic based on geometric considerations concerning angles; in practice, having fixed a threshold $\gamma$ and a starting value for the parameter $k$, it applies kNN to find the neighbors of each $\mathbf{p}_i \in \mathbf{P}_N$ and calculates the angle $\nu_i$ between the $(k+1)$th-nearest neighbor and the subspace spanned by the $k$-nearest neighbors. Considering a threshold parameter $\gamma$, if $(1/N) \sum_{i=1}^{N} \nu_i \leq \gamma$, then $k$ is considered as the id estimate; otherwise, $k$ is incremented by 1 and the process is repeated. The main limitation of this method is the difficult choice of a proper value for $\gamma$.

The work presented by Pettis et al. [26] is notable since it is one of the first works providing a mathematical motivation for the use of nearest-neighbor distances.

Indeed, for an i.i.d. sample $\mathbf{P}_N \subseteq \mathfrak{R}^D$ drawn from a density distribution $p(\mathbf{x})$ in $\mathfrak{R}^d$, the following approximation holds:

$$\frac{k}{N} \simeq p(\mathbf{x}) V(d) r^d, \tag{11}$$

where $k$ is the number of nearest neighbors to $\mathbf{x}$ within the hypersphere $\mathscr{B}_d(\mathbf{x}, r)$ of radius $r$ and centered on $\mathbf{x}$ and $V(d)$ is the volume of the (unit $d$-dimensional) ball in $\mathfrak{R}^d$.

This means that the proportion of sample points falling in $\mathscr{B}_d(\mathbf{x}, r)$ is roughly approximated by $p(\mathbf{x})$ times the volume of $\mathscr{B}_d(\mathbf{x}, r)$. Since this volume grows as $r^d$, assuming the density $p(x)$ to be a constant, it follows that the number of samples in $\mathscr{B}_d(\mathbf{x}, r)$ is proportional to $r^d$. From the relationship in (11), and assuming that the samples are locally uniformly distributed, the authors derive an id estimator for $d$:

$$\widehat{d} = \frac{\overline{r}_k}{k(\overline{r}_{k+1} - \overline{r}_k)}, \tag{12}$$

where $\overline{r}_k$ is the average of the distances from each sample point to its $k$th nearest neighbors; defining $r_i^{(k)}$ as the distance between $\mathbf{x}_i$ and its $k$th-nearest neighbor, $\overline{r}_k$ is expressed as $\overline{r}_k = (1/N) \sum_{i=1}^{N} r_i^{(k)}$.

Since this algorithm is limited by the choice of a suitable value for parameter $k$, in [63] the authors propose a variant which considers a range of neighborhood sizes $[k_{\min}, k_{\max}]$. However, in the same work the authors themselves show that this technique generally yields an underestimate of the id when its value is high.

Taking into account relation (11), in [101] the number $N_{\mathscr{B}_d}$ of data points in $\mathscr{B}_d(\mathbf{x}, r)$ is described by a polynomial $f(r) = \sum_{s=0}^{d} \beta_s r^s$ of degree $d$. In practice, considering $\mathbf{p}_i, \mathbf{p}_k \in \mathbf{P}_N$, calling $r_{ik} = \|\mathbf{p}_i - \mathbf{p}_k\|$ the interpoint distances, and being $r = \min_{i,k=1}^{N} r_{ik}$, $R$ a parameter adaptively estimated (to estimate $R$ by means of $\mathbf{P}_N$, the radius value corresponding to the first significant peak of the histogram of the $r_{ij}$s is found), a set of $n$ radius values $\mathbf{r} = \{r_j = r + j(R-r)/n\}_{j=1}^{n}$ is selected and used to calculate $n$ pairs $\{(r_j, \widehat{f}(r_j))\}_{j=1}^{n}$, where $\widehat{f}(r_j) = \#[r_{ik} < r_j]_{i,k=1}^{N}$ is the number of interpoint distances strictly lower than $r_j$. To estimate the coefficients $\{\beta_j\}_{j=1}^{D}$, the computed pairs are fit by a least squares fitting procedure that estimates exactly $D+1$ coefficients. Since by hypothesis the degree of $f$ is $d$, the significance test described in [17] is used to estimate the degree $\widehat{d}$ of $\widehat{f}$, which is considered as the id estimate. The comparative evaluation of this algorithm with the well-known Maximum Likelihood Estimator (MLE) [27] and its improved version [102], both described below, has shown that it is more robust than them when dealing with high-dimensional datasets.

MLE [27], one of the most cited estimators, treats the neighbors of each point $\mathbf{p}_i \in \mathbf{P}_N$ as events in a Poisson process and the distance $r^{(j)}(\mathbf{p}_i)$ between the query point $\mathbf{p}_i$ and its $j$th nearest neighbor as the event's arrival time. Since this process depends on $d$, MLE estimates id by maximizing the log-likelihood of the observed process. In practice a local id estimate is computed as

$$\widehat{d}(\mathbf{p}_i, k) = \left( \frac{1}{k} \sum_{j=1}^{k} \log \frac{r^{(k+1)}(\mathbf{p}_i)}{r^{(j)}(\mathbf{p}_i)} \right)^{-1}. \tag{13}$$

Averaging the $\widehat{d}(\mathbf{p}_i, k)$s, the global id estimate is $\widehat{d}(k) = (1/N) \sum_{i=1}^{N} \widehat{d}(\mathbf{p}_i, k)$.

The theoretical stability of the proposed `id` estimator for data living in $C^1$ submanifold of $\mathfrak{R}^D$, $d \le D$, and for data in an affine subspace of $\mathfrak{R}^D$ has been proved, respectively, in [103, 104]. Though the authors' comparative evaluation shows the superior performance of the proposed estimator with respect to the `CD` estimator [19] (see Section 3.2.1) and the `NN` estimator [26], they further improve it by removing its dependency from the parameter $k$; to this end, different values for $k$ are adopted and the computed results are averaged to obtain the final `id` estimate: $\widehat{d} = (1/t) \sum_{k \in \{k_1, \ldots, k_t\}} \widehat{d}(k)$.

Considering that, in practice, `MLE` is highly biased for both large and small values of $k$, a variant of `MLE` is proposed in [102], where the arithmetic mean is substituted with the harmonic average, leading to the following estimator: $\widehat{d}(k) = ((1/N) \sum_{i=1}^N (1/\widehat{d}(\mathbf{p}_i, k)))^{-1}$.

Though the proposal in [102] seems to achieve more accurate results, it is based on the assumption that neighbors surrounding each $\mathbf{p}_i$ are independent, which is clearly incorrect. To cope with this problem, in [105] an interesting regularized version of `MLE` applies a regularized maximum likelihood technique to distances between neighbors. The comparative evaluation with the aforementioned `MLE` methods [27, 102] makes the authors state that, though the method might be the first to converge to the actual estimate given enough data points, its estimation accuracy is comparable to that achieved by the competing schemes.

In [59, 106] a further improvement of `MLE` is presented; it achieves a better performance by substituting euclidean distances with geodesic ones.

Despite the good results achieved by `MLE`-based approaches, these techniques have shown to be affected by the curvature induced by $\phi$ on the manifold neighborhoods approximated by `kNN`. To reduce this effect, various `id` estimators have been proposed in [76, 107]; here, we review those achieving the most promising experimental results.

In [107] the authors firstly propose a family of `id` estimators (`MiND`$_{\text{ML}*}$), which exploit the `pdf` $g(r; k, d)$ describing the distance $r^{(1)}(\mathbf{x})$ between the center $\mathbf{x}$ of $\mathscr{B}_d(\mathbf{x}, r)$, $\mathbf{x} \in \mathscr{M}$, $r \to 0^+$ and its nearest neighbor. Briefly, formulating $g(r; k, d)$ as a function of the `id` value $d$ ($g(r; k, d) = kdr^{d-1}(1 - r^d)^{k-1}$), the `id` estimator is computed by a maximum likelihood approach.

After noting that this algorithm is still affected by a bias causing underestimations when the dataset dimensionality becomes sufficiently high (i.e., $id \ge 10$), the authors present theoretical considerations which relate the bias to the fact that `id` estimators based on nearest-neighbor distances are often founded on statistics derived under the assumption that the amount of available data is unlimited, which is never the case in practical applications. Based on these considerations, two different estimators, named `MiND`$_{\text{KL}}$ and `IDEA`, are presented.

`MiND`$_{\text{KL}}$ compares the empirical `pdf` of the neighborhood distances computed on the dataset ($g_{\text{Data}}$) with the distribution of the neighborhood distances computed from points uniformly drawn from hyperspheres of known increasing dimensionality ($g_{\text{Sphere}}^d$). The `id` estimate is the dimensionality that minimizes their Kullback-Leibler divergence

$\mathscr{KL}(g_{\text{Data}}, g_{\text{Sphere}}^d)$, which is evaluated by means of the data-driven technique proposed in [108].

`IDEA` relies on the authors' observation that the quantities $1 - r^{(j)}(\mathbf{p}_i)/r^{(k+1)}(\mathbf{p}_i)$ are distributed according to the beta distribution $\beta_{1,d}$ with parameters 1 and $d$, respectively. Therefore, since $\mathbb{E}[\beta_{1,d}] = m = 1/(1 + d)$, a consistent `id` estimator $\widehat{d} \simeq d$ equals

$$\widehat{d} = \frac{\widehat{m}}{1 - \widehat{m}} \simeq d = \frac{m}{1 - m},$$

$$\text{where } \widehat{m} = \frac{1}{Nk} \sum_{i=1}^N \sum_{j=1}^k \frac{r^{(j)}(\mathbf{p}_i)}{r^{(k+1)}(\mathbf{p}_i)} \simeq m. \tag{14}$$

To reduce the effect of the aforementioned bias, `IDEA` finally applies an asymptotic correction step that, inspired by the correction method presented in [28], models the underestimation error by considering both the base algorithm and the given dataset.

Motivated by the promising results achieved by `MiND`$_{\text{KL}}$, in [76] the authors propose its extension, namely, `DANCo`; it further reduces the underestimation effect by combining an estimator employing normalized nearest-neighbor distances with one employing mutual angles. More precisely, `DANCo` compares the statistics estimated on $\mathbf{P}_N$ with those estimated on (uniformly drawn) synthetic datasets of known `id`. The comparisons are performed by two Kullback-Leibler divergences applied to the distribution of normalized nearest-neighbor distances $g(r; k, d)$, having $g(r; k, d) = kdr^{d-1}(1 - r^d)^{k-1}$, and the distribution of pairwise angles $q(\mathbf{x}; \boldsymbol{\nu}, \tau)$, $q(\mathbf{x}; \boldsymbol{\nu}, \tau)$ being the von Mises-Fisher distribution (`VMF`) [109] with parameters $\boldsymbol{\nu}$ and $\tau$.

The `id` estimate $\widehat{d}$ is the one minimizing the sum of the two divergences:

$$\widehat{d} = \underset{d \in \{1, \ldots, D\}}{\arg \min} \; \mathscr{KL}\left(g_{\text{Data}}, g_{\text{Sphere}}^d\right)$$
$$+ \mathscr{KL}\left(q_{\text{Data}}, q_{\text{Sphere}}^d\right). \tag{15}$$

A fast implementation of `DANCo` (Fast-`DANCo`) is also developed. Comparative evaluations show that this algorithm achieves promising results (as shown in [76] and Section 4).

Another work, which is notable because the authors not only prove the consistency in probability of the presented estimators but also derive upper bounds (see (19) below) on the probability of the estimation-error for finite, and large enough, values of $N$, is proposed in [33]. More precisely, the authors introduce two estimators by firstly defining a function $\eta : \mathfrak{R}^D \times \mathfrak{R} \to \mathfrak{R}^+$ slowly varying near the origin (see [33] for a detailed description and motivation of this assumption). The function $\eta$ is then used to express the logarithm of the probability of a point $\mathbf{p}$ of being in the hypersphere $\mathscr{B}_D(\mathbf{p}_i, r)$: $\log(\mathbf{P}(\mathbf{p} \in \mathscr{B}_D(\mathbf{p}_i, r))) = \log(\eta(\mathbf{p}, r)) + d\log(r)$, having $\mathbf{P}(\mathbf{p} \in \mathscr{B}_D(\mathbf{p}_i, r)) = \eta(\mathbf{p}, r)r^d$.

Considering that $\mathbf{P}(\mathbf{p} \in \mathscr{B}(\mathbf{p}_i, r^{(k)}(\mathbf{p}_i))) \approx k/n$ for $N$ big enough, the authors derive the following system of equations:

$$\log\left(\frac{k}{n}\right) \approx \log\left(\eta\left(\mathbf{p}_i, r\right)\right) + \widehat{d}\left(\mathbf{p}_i\right) \log\left(r^{(k)}\left(\mathbf{p}_i\right)\right),$$

$$\log\left(\frac{k}{(2n)}\right) \approx \log\left(\eta\left(\mathbf{p}_i, r\right)\right) \tag{16}$$

$$+ \widehat{d}\left(\mathbf{p}_i\right) \log\left(r^{(\lceil k/2 \rceil)}\left(\mathbf{p}_i\right)\right),$$

and solve it for $\widehat{d}(\mathbf{p}_i)$ to obtain a local id estimate:

$$\widehat{d}\left(\mathbf{p}_i\right) = \frac{\log\left(2\right)}{\log\left(r^{(k)}\left(\mathbf{p}_i\right)/r^{(\lceil k/2 \rceil)}\left(\mathbf{p}_i\right)\right)}. \tag{17}$$

The two proposed estimators are then computed either by averaging ($\widehat{d}_{\text{avg}}$) or by voting ($\widehat{d}_{\text{vote}}$):

$$\widehat{d}_{\text{avg}} = \frac{1}{N}\sum_{i=1}^{N} \widehat{d}\left(\mathbf{p}_i\right),$$

$$\widehat{d}_{\text{vote}} = \underset{d' \in \mathbb{N}^+}{\arg\max} \# \left[\widehat{d}\left(\mathbf{p}_i\right) = d'\right]_{i=1}^{N}, \tag{18}$$

where $\#[\text{cond}]_{i=1}^{N}$ denotes the number of points $\mathbf{p}_i$ for which cond is verified.

Under differentiability assumptions on the function $\eta$ and regularity assumptions on $\mathscr{M}$ the authors prove the consistency in probability of their estimators and provide upper bounds (see (19)) on the probability of the estimation-error for finite, and large enough, values of $N$. However, the derived bounds depend on unknown universal constants $c, c', c'' > 0$:

$$\mathbf{P}\left(\widehat{d}_{\text{avg}} \neq d\right) \leq \exp\left(-\frac{c'N}{\left(Dc^d k\right)^2}\right),$$

$$\mathbf{P}\left(\widehat{d}_{\text{vote}} \neq d\right) \leq \exp\left(-\frac{c''N}{\left(c^d k\right)^2}\right). \tag{19}$$

*3.3. Graph-Based* id *Estimators.* As noted in [96], the work of [110] has cleared up the fact that theories underlying graphs can be applied to solve a variety of statistical problems; indeed, also in the field of id estimation various types of graph structures have been proposed [29, 96, 111, 112] and used for id estimation. Examples are the kNN graph (kNNG) [30], the Minimum Spanning Tree (MST) [113] and its variation, the geodesic MST (GMST) [29]; and the sphere of influence graph (SIG) [114] and its generalization, the $k$−sphere of influence graph (kSIG) [96].

Given a sample set $\mathbf{P}_N = \{\mathbf{p}_i\}_{i=1}^{N}$ a graph built on $\mathbf{P}_N$, usually denoted by $G(\mathbf{P}_N) = (\{\mathbf{p}_i\}_{i=1}^{N}, \{e_{i,j}\}_{i,j \in \{1,\ldots,N\}})$, employs the sample points $\mathbf{p}_i$ as nodes (vertices) of the graph and connects them with weighted arcs (edges) $\{e_{i,j}\}_{i,j \in \{1,\ldots,N\}}$.

A $\text{kNNG}_N(\mathbf{P}_N)$ is built by employing a distance function, which commonly is the Euclidean one, to weight the arcs connecting each $\mathbf{p}_i$ to its kNNs.

A $\text{MST}(\mathbf{P}_N)$ is the spanning tree minimizing the sum of the edge weights. When the weights approximate Geodesic distances [56], a $\text{GMST}_N(\mathbf{P}_N)$ is obtained.

A $\text{SIG}_N(\mathbf{P}_N)$ is defined by connecting nodes $\mathbf{p}_i$ and $\mathbf{p}_j$ iff $\|\mathbf{p}_i - \mathbf{p}_j\| \leq \rho(i) + \rho(j)$, where $\rho(i)$ is the distance between $\mathbf{p}_i$ and its nearest neighbor in $\mathbf{P}_N$. Essentially, two vertices $\mathbf{p}_i$ and $\mathbf{p}_j$ are connected if the corresponding NN hyperspheres intersect. A generalization of $\text{SIG}_N(\mathbf{P}_N)$ is $\text{kSIG}(\mathbf{P}_N)$, where nodes $\mathbf{p}_i$ and $\mathbf{p}_j$ are connected iff $\|\mathbf{p}_i - \mathbf{p}_j\| \leq \rho_k(i) + \rho_k(j)$, $\rho_k(i)$ being the distance between $\mathbf{p}_i$ and its kNN in $\mathbf{P}_N$. This means that the kNN hyperspheres centered on $\mathbf{p}_i$ and $\mathbf{p}_j$ intersect.

In the following we recall interesting id estimators based on $\text{GMST}(\mathbf{P}_N)$, $\text{kNNG}(\mathbf{P}_N)$, and $\text{kSIG}(\mathbf{P}_N)$.

In [29, 30], after defining the length functional $\mathscr{L}(G_N(\mathbf{P}_N)) = \sum |e_{i,j}|^{\gamma}$, $\gamma \in (0,d)$, to build either the $\text{GMST}(\mathbf{P}_N)$ or the $\text{MST}(\mathbf{P}_N)$ of $\text{kNNG}(\mathbf{P}_N)$, graph theories are exploited to estimate both the id of the underlying manifold structure $\mathscr{M}$ and its intrinsic Rènyi $\alpha$-entropy $\mathscr{H}_{\mathscr{M}}$. To this aim, the authors derive the linear model: $\log\mathscr{L}(\text{MST}(\mathbf{P}_N)) = a \log d + b, a = (d-\gamma)/d, b = \log c + \mathscr{H}_{\mathscr{M}}$, $c$ being an unknown constant, and exploit it to define an estimator of both $d$ and $\mathscr{H}_{\mathscr{M}}$. Briefly, a set of cardinalities $\{n_k\}_{k=1}^{K}$ is chosen and, for each $n_k$, the $\text{MST}(\mathbf{P}_{n_k})$ is constructed on the set $\mathbf{P}_{n_k}$, which contains $n_k$ points randomly sampled from $\mathbf{P}_N$, to obtain a set of $K$ pairs $(\log\mathscr{L}(\text{MST}(\mathbf{P}_{n_k})), n_k)$. Fitting them with a least squares procedure the estimates $\widehat{a} \simeq a$ and $\widehat{b} \simeq b$ are computed. Recalling that $a = (d-\gamma)/d$, the id is calculated as $\widehat{d} = \text{round}\{\gamma/(1 - \widehat{a})\} \simeq d$. This process is iterated to produce the final estimate as the average of the obtained results.

The aforementioned kNNG based algorithm [29, 30] is exploited in [112], where the authors consider datasets sampled from a union of disjoint manifolds with possibly different ids. To estimate the local ids, the authors propose a heuristic, which is not described here, to automatically determine the local neighborhoods with similar geometric structures without any prior knowledge on the number of manifolds, their ids, and their sampling distributions.

In [96] the authors present three id estimation approaches, defined as "graph theoretic methods" since the statistics they compute are functions only of graph properties (such as vertex degrees and vertex eccentricities) and do not directly depend on the interpoint distances.

The first statistic, denoted by $S_N^1(\mathbf{P}_N) = \overline{r}_j(\text{kNNG}(\mathbf{P}_N))$ in the following, is based on the reach (the reach $r_{j,i}(\mathbf{p}_i, G)$, in $j$ steps of a node $\mathbf{p}_i \in G$, is the total number of vertices which are connected to $\mathbf{p}_i$ by a path composed of $j$ arcs or less in $G$) of vertices in the $\text{kNNG}(\mathbf{P}_N)$. Considering that the reach of each vertex $\mathbf{p}_i \in \text{kNNG}(\mathbf{P}_N)$ grows as the id increases, in [115] the average reach $\overline{r}_j(\text{kNNG})$ in $j$ steps of vertices in $\text{kNNG}(\mathbf{P}_N)$ is employed: $S_N^1(\mathbf{P}_N) = \overline{r}_j(\text{kNNG}(\mathbf{P}_N)) = (1/N)\sum_{i=1}^{N} r_{j,i}(\mathbf{p}_i, \text{kNNG}(\mathbf{P}_N))$.

The second statistic, denoted by $S_N^2(\mathbf{P}_N) = M_N(\text{MST}(\mathbf{P}_N))$, is computed by considering the degree of vertices in the $\text{MST}(\mathbf{P}_N)$. Recalling that, for datasets $\mathbf{P}_N$ obtained from a continuous distribution on $\mathfrak{R}^d$, the ratio of nodes with a given degree $j$ in $\text{MST}_N(\mathbf{P}_N)$ converges a.s.

to a limit depending only on $j$ and $d$ [116] and that the average degree in a tree is a constant depending only on the number of vertices, the authors empirically observe a dependency between the average degree and the id. This leads to the definition of an id estimator employing the statistic $S_N^2 = M_N(\text{MST}(\mathbf{P}_N)) = (1/N) \sum_{i=1}^N (\deg_{\text{MST}(\mathbf{P}_N)}(\mathbf{p}_i))^2$.

The third statistic, denoted by $S_N^3(\mathbf{P}_N) = U_N^k(\text{kSIG}(\mathbf{P}_N))$, is motivated by studies in the literature [117] showing that the expected number of neighbors shared by a given pair of points depends on the id of the underlying manifold. Accordingly, calling $N_{i,j}$ the number of samples in the intersection of the two kNN hyperspheres centered on $\mathbf{p}_i$ and $\mathbf{p}_j$, intuitions similar to those considered for $\bar{r}_j(\text{kNNG})$ lead to defining of $S_N^3(\mathbf{P}_N) = U_N^k(\text{kSIG}(\mathbf{P}_N)) = (1/n) \sum_{i \leq j} N_{i,j}$.

Based on their theoretical results and empirical tests on synthetically generated datasets characterized by id values $d_j$ in a finite range $\mathbf{F} \subseteq N^+$ (where $\mathbf{F} = \{d_j\}_{d_j=2}^{12}$ in the reported experiments), the authors propose an approximate Bayesian estimator that could indistinctly employ each of the three statistics $S_N^1$, $S_N^2$, and $S_N^3$, denoted by $S_N^*$ in the following. To this aim, they assume that each statistic can be approximated by a Gaussian density $f_{d_j}(\cdot) = \mathcal{N}(\mu(d_j), \sigma^2(d_j))$; to estimate $\mu(d_j)$ and $\sigma^2(d_j)$, for each $d_j \in \mathbf{F}$, $L$ datasets of large size are synthetically generated by random sampling from the uniform distribution on the unit $d_j$-cube. These datasets are then used to estimate the parameters $\tilde{\mu}(d_j) \simeq \mu(d_j)$ and $\tilde{\sigma}^2(d_j) \simeq \sigma^2(d_j)$ that define the approximation $\tilde{f}_{d_j}(\cdot)$, computed on a generic sample set with size $N$ and id $= d_j$, of the Gaussian density $f_{d_j}(\cdot)$ of $S_N^*$.

At this stage, given a new input dataset $\mathbf{P}_N$ having unknown id, the statistic $S_N^*(\mathbf{P}_N) = s_{\mathbf{P}}$ is computed and used to calculate the approximated value $\tilde{f}_{d_j}(s_{\mathbf{P}}) = \mathcal{N}(\tilde{\mu}^2(d_j), \tilde{\sigma}^2(d_j)/N) \simeq f_{d_j}(s_{\mathbf{P}})$. Assuming equal a priori probability for all the $d_j \in \mathbf{F}$, the posterior probability $P[d_j \mid S_N^*]$ is given by

$$P\left[d_j S_N^*\right] = \frac{\tilde{f}_{d_j}(s_{\mathbf{P}})}{\sum_{d_j \in \mathbf{F}} \tilde{f}_{d_j}(s_{\mathbf{P}})}, \quad d_j \in \mathbf{F}, \tag{20}$$

and employed to compute an "a posteriori expected value" of the id:

$$\hat{d} = \text{round} \left\{ \sum_{d_j \in \mathbf{F}} d_j P\left[d_j S_N^*\right] \right\}. \tag{21}$$

The authors evaluate the performance of their methods on synthetic datasets, some of which have been used by similar studies in the literature [79], while the others (challenging ones) are proposed by the authors to have manifolds with nonconstant curvature. The comparison of the achieved results with those obtained by the estimators proposed in [27, 33, 112, 118] has led to the conclusion that none of the methods has a good performance on all the tested datasets. However, graph theoretic approaches would appear to behave better when manifolds of nonconstant curvature are processed.

This interesting comparison strengthens the need of defining a benchmark framework to allow an objective and reproducible comparative evaluation of id estimators. For this reason, in Section 4 we describe our proposal in this direction.

## 4. A Benchmark Proposal

At the present, an objective comparison of different id estimators is not possible due to the lack of a standardized benchmark framework; therefore, in this section, after recalling experimental datasets and evaluation procedures introduced in the literature (see Sections 4.1 and 4.2), we choose some of them to propose a benchmark framework (see Section 4.3) that allows for reproducible and comparable experimental setups. The usefulness of the proposed benchmark is then shown by employing it to compare relevant state-of-the-art id estimators whose code is publicly available (see Section 4.4).

*4.1. Datasets.* The datasets employed in the literature are both synthetically generated datasets and real ones. In the following sections we describe those we choose to use in our benchmark study.

*4.1.1. Synthetic Datasets.* Synthetic datasets are generated by drawing samples from manifolds ($\mathcal{M}$) linearly or nonlinearly embedded in higher dimensional spaces.

The publicly available tool (http://www.mL.uni-saarland.de/code/IntDim/IntDim.htm) proposed by Hein and Audibert in [79] allows us to generate 13 kinds of synthetic datasets by uniformly drawing samples from 13 manifolds of known id; they are schematically described in Table 1, where they are referred to as $\mathcal{M}_*^H$. These manifolds are embedded in higher dimensional spaces through both linear and nonlinear maps and are characterized by different curvatures. We note that manifold $\mathcal{M}_8^H$ is particularly challenging for its high curvature; indeed, when it is used for testing, most relevant id estimators compute pronounced id overestimates (see also the results reported in [107]).

Another interesting dataset [96] is generated by sampling a $d$-dimensional paraboloid, $\mathcal{M}_{Pd}$, nonlinearly embedded in a higher $(3(d + 1))$ dimensional space, according to a multivariate Burr distribution with parameter $\alpha = 1$. Tests on this dataset are particularly challenging since the underlying manifold is characterized by a nonconstant curvature.

To perform tests on datasets generated by employing a smooth nonuniform pdf, we propose the dataset $\mathbf{M}_{\text{beta}}$, obtained as follows: we sample $N$ points in $[0, 1)^{10}$, according to a beta distribution $\beta_{0.5,10}$ with parameters 0.5 and 10, respectively (high skewness), and store them in a matrix $\mathbf{X}_N \in \mathfrak{R}^{N \times 10}$; multiply each point of $\mathbf{X}_N$ ($\mathbf{X}_N(i, j)$) by $\sin(\cos(2\pi \mathbf{X}_N(i, j)))$, thus obtaining a matrix $\mathbf{D}_1 \in \mathfrak{R}^{N \times 10}$; multiply each point of $\mathbf{X}_N$ by $\cos(\sin(2\pi \mathbf{X}_N(i, j)))$, thus obtaining another matrix $\mathbf{D}_2 \in \mathfrak{R}^{N \times 10}$; append $\mathbf{D}_1$ and $\mathbf{D}_2$ to generate a matrix $\mathbf{D}_3 \in \mathfrak{R}^{2500 \times 20}$; append $\mathbf{D}_3$ to its duplicate to finally generate a test dataset containing $N$ points in $\mathfrak{R}^{40}$.
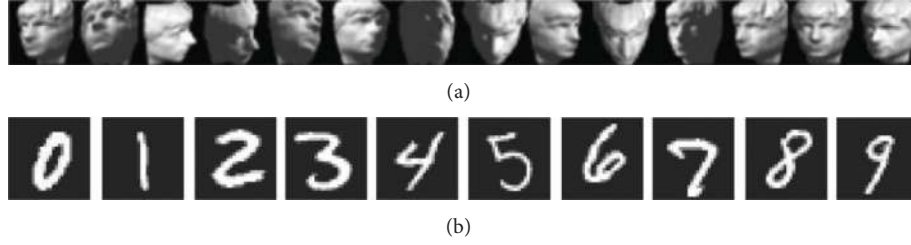
(a)



(b)

FIGURE 2: (a) Samples from ISOMAP face database. (b) Samples from digit "0" to digit "9" in MNIST database.

TABLE 1: The 13 types of synthetic datasets generated with the tool proposed in [79].

| Dataset | Underlying manifold name | Description | $\mathbf{d}$ | $\mathbf{D}$ |
|---|---|---|---|---|
| Synthetic | $\mathscr{M}_1^H$ | $d$-dimensional sphere linearly embedded | $\mathbf{D}-1$ | *User-defined* |
| | $\mathscr{M}_2^H$ | Affine space | 3 | 5 |
| | $\mathscr{M}_3^H$ | Concentrated figure, mistakable with a 3-dimensional one | 4 | 6 |
| | $\mathscr{M}_4^H$ | Nonlinear manifold | 4 | 8 |
| | $\mathscr{M}_5^H$ | 2-dimensional helix | 2 | 3 |
| | $\mathscr{M}_6^H$ | Nonlinear manifold | 6 | 36 |
| | $\mathscr{M}_7^H$ | Swiss-Roll | 2 | 3 |
| | $\mathscr{M}_8^H$ | Nonlinear (highly curved) manifold | 12 | 72 |
| | $\mathscr{M}_9^H$ | Affine space | $\mathbf{D}$ | *User-defined* |
| | $\mathscr{M}_{10}^H$ | $d$-dimensional hypercube | $\mathbf{D}-1$ | *User-defined* |
| | $\mathscr{M}_{11}^H$ | Möebius band 10-times twisted | 2 | 3 |
| | $\mathscr{M}_{12}^H$ | Isotropic multivariate Gaussian | $\mathbf{D}$ | *User-defined* |
| | $\mathscr{M}_{13}^H$ | 1-dimensional helix curve | 1 | *User-defined* |

To further test estimators' performance on nonlinearly embedded manifolds of high id, we propose to generate two datasets, referred to as $\mathbf{M}_{N1}$ and $\mathbf{M}_{N2}$ in the following (a tool to generate the datasets sampled from $d$-dimensional paraboloids, the $\mathbf{M}_{\text{beta}}$ dataset, the $\mathbf{M}_{N1}$ dataset, and the $\mathbf{M}_{N2}$ dataset, is available at http://security.di.unimi.it/~fox721/ dataset_generator.m). Precisely, to generate $\mathbf{M}_{N1}$ we uniformly draw $N$ points in $[0,1]^{18}$, we transform each point by means of $\tan(\mathbf{x}^i\cos(\mathbf{x}^{18-i+1}))$ where $i = 1,\ldots,18$, we obtain points in $\mathfrak{R}^{36}$ by appending each transformed $\mathbf{x}$ to $\arctan(\mathbf{x}^{18-i+1}\sin(\mathbf{x}^i))$, and we duplicate the coordinates of each point to finally generate points in $\mathfrak{R}^{72}$. The id of $\mathbf{M}_{N1}$

is 18, and its points are drawn from a manifold nonlinearly embedded in $\mathfrak{R}^{72}$. To generate $\mathbf{M}_{N2}$ containing $N$ points in $\mathfrak{R}^{96}$, we applied the same procedure on vectors sampled in $[0,1]^{24}$.

*4.1.2. Real Datasets.* Real datasets employed in the literature generally concern problems in the fields of image analysis, signal processing, time series prediction, and biochemistry. Among them, the most known and used datasets are ISOMAP face database [56], MNIST database [119], Isolet dataset [120], $D2$ Santa Fe [121] dataset, and DSVC1 time series [21]. Recently, the Crystal Fingerprint space for the chemical compound silicon dioxide dataset has also been proposed [22].

ISOMAP face database consists in 698 gray-level images of size $64 \times 64$ depicting the face of a sculpture. This dataset has three degrees of freedom: two for the pose and one for the lighting direction (see Figure 2(a)).

MNIST database consists in 70000 gray-level images of size $28 \times 28$ of hand-written digits (see Figure 2(b)). The real id of this database is not actually known, but some works [79, 122] propose similar estimates for the different digits; as an example, the proposed id values for the digit "1" are in the range $\{8,\ldots,11\}$.

Isolet dataset has been generated as follows: 150 subjects spoke the name of each letter of the alphabet twice, thus producing about 52 training examples from each speaker, for a total of 7797 samples. The speakers are grouped into 5 sets of 30 speakers each and are referred to as *isolet*1, *isolet*2, *isolet*3, *isolet*4, and *isolet*5. The real id value characterizing this dataset is not actually known, but a study reported in [123] shows that the correct estimate could be in the range $\{16,\ldots,22\}$.

The version $D2$ of Santa Fe dataset is a time series of 50000 one-dimensional points having nine degrees of freedom (id = 9) and being generated by a simulation of particle motion. In order to estimate the attractor dimension of this time series, it is possible to employ the method of delays described in [124], which generates $D$-dimensional vectors by partitioning the original dataset in blocks containing $D$ consecutive values; as an example, by choosing $D = 50$ a dataset containing 1000 points in $\mathfrak{R}^{50}$ is obtained.

DSVC1 is a time series composed by 5000 samples measured from a hardware realization of Chua's circuit [125]. Employing the method of delays with $D = 20$, a dataset containing 250 points in $\mathfrak{R}^{20}$ is obtained. The id characterizing this dataset is ~2.26 [21].

TABLE 2: Synthetic datasets and real datasets suggested by the benchmark; **N** is the dataset cardinality, **d** is the id, and **D** is the embedding space dimension.

| Dataset | Dataset name | N | d | D |
|---|---|---|---|---|
| Synthetic | $\mathbf{M}_1$ | 2500 | 10 | 11 |
| | $\mathbf{M}_2$ | 2500 | 3 | 5 |
| | $\mathbf{M}_3$ | 2500 | 4 | 6 |
| | $\mathbf{M}_4$ | 2500 | 4 | 8 |
| | $\mathbf{M}_5$ | 2500 | 2 | 3 |
| | $\mathbf{M}_6$ | 2500 | 6 | 36 |
| | $\mathbf{M}_7$ | 2500 | 2 | 3 |
| | $\mathbf{M}_9$ | 2500 | 20 | 20 |
| | $\mathbf{M}_{10a}$ | 2500 | 10 | 11 |
| | $\mathbf{M}_{10b}$ | 2500 | 17 | 18 |
| | $\mathbf{M}_{10c}$ | 2500 | 24 | 25 |
| | $\mathbf{M}_{10d}$ | 2500 | 70 | 71 |
| | $\mathbf{M}_{11}$ | 2500 | 2 | 3 |
| | $\mathbf{M}_{12}$ | 2500 | 20 | 20 |
| | $\mathbf{M}_{13}$ | 2500 | 1 | 13 |
| | $\mathbf{M}_{N1}$ | 2500 | 18 | 72 |
| | $\mathbf{M}_{N2}$ | 2500 | 24 | 96 |
| | $\mathbf{M}_{beta}$ | 2500 | 10 | 40 |
| | $\mathbf{M}_{P3}$ | 2500 | 3 | 12 |
| | $\mathbf{M}_{P6}$ | 2500 | 6 | 21 |
| | $\mathbf{M}_{P9}$ | 2500 | 9 | 30 |
| Real | $\mathbf{M}_{\texttt{DSCV1}}$ | 250 | 2.26 | 20 |
| | $\mathbf{M}_{\texttt{ISOMAP}}$ | 698 | 3.00 | 4096 |
| | $\mathbf{M}_{\texttt{Santa Fe}}$ | 1000 | 9.00 | 50 |
| | $\mathbf{M}_{\texttt{MNIST1}}$ | 70000 | 8.00–11.00 | 784 |
| | $\mathbf{M}_{\texttt{SiO2}}$ | 4738 | 12.00 | 1800 |
| | $\mathbf{M}_{\texttt{Isolet}}$ | 7797 | 16.00–22.00 | 617 |

Crystal Fingerprint spaces, or Crystal Finger spaces, have been recently proposed in crystallography [22] with the aim of representing crystalline structures; these spaces are built starting from the measured distances between atoms in the crystalline structure. The theoretical id of one Crystal Finger space consists in $3N_a + 3$ crystal degrees of freedom, where $N_a$ is the number of atoms in the crystalline unitary cell.

*4.2. Experimental Procedures and Evaluation Measures.* At the state of the art, two approaches have been mainly used to assess id estimators on datasets of known id.

The first one subsamples the test dataset to obtain $T$ subsets of fixed cardinality and computes the percentage of correct estimations. To analyze estimators' behavior with respect to the cardinality of input datasets, this procedure may be repeated by using different cardinality values [29, 30, 79, 122], thus obtaining a distinct performance evaluation measure for each cardinality.

The second approach estimates the id on $T$ permutations of the same dataset and averages the $T$ id estimates to obtain the final one [27, 76, 107, 126]. This value is then compared with the real one to assess the id estimator.

To also test the estimator's robustness with respect to its parameter settings, in [27, 107, 126] the authors apply

a further test, originally proposed by Levina and Bickel in [27]. Precisely, sample sets with different cardinalities are drawn from the standard Gaussian pdf in $\Re^5$ and, for each set, the estimator is applied varying the values of its parameters in fixed ranges; this allows us to analyze the behavior of the id estimate as a function of both the dataset's cardinality and the parameter settings.

Note that, since id estimators are usually tested on different datasets to evaluate their reliability when confronted by different dataset structures and configurations, in [126] an overall evaluation measure is proposed. This indicator, called Mean Percentage Error (MPE), summarizes all the obtained results in a unique value computed as MPE $= (100/\#\mathbf{M}) \sum_{\mathbf{M}} (|\hat{d}_{\mathbf{M}} - d_{\mathbf{M}}|/d_{\mathbf{M}})$, where $\#\mathbf{M}$ is the number of tested datasets, $\hat{d}_{\mathbf{M}}$ is the id estimated on the dataset $\mathbf{M}$, and $d_{\mathbf{M}}$ is the real id of $\mathbf{M}$. To apply this technique to real datasets whose id belongs to the range $\{d_{\min}, \ldots, d_{\max}\}$, the same authors propose to calculate the associated MPE's term as $\min_{d \in \{d_{\min}, \ldots, d_{\max}\}} (|\hat{d}_{\mathbf{M}} - d|/d_{\mathbf{M}})$, where $d_{\mathbf{M}}$ is the mean of the range.

*4.3. Benchmark.* In this section we propose an evaluation approach which can be used as a standard framework to

assess estimators performance, comparing it to relevant `id` estimators whose code is publicly available. In this benchmark, we suggest to use the following estimators (see Section 3): `Hein`, `MLE`, `kNNG`, `MLSVD`, `BPCA`, `CD`, `MiND`$_{KL}$, and `DANCo` (the source code of the mentioned methods is available at `Hein`: http://www.mL.uni-saarland.de/code.shtml, `MLE`: http://dept.stat.lsa.umich.edu/~elevina/mledim.m, `kNNG`: http://web.eecs.umich.edu/~hero/IntrinsicDim/, `MLSVD`: http://www.math.duke.edu/~mauro/code.html#MSVD, `BPCA`: http://research.microsoft.com/en-us/um/cambridge/projects/infernet/blogs/bayesianpca.aspx, `CD`: http://cseweb.ucsd.edu/lvdmaaten/dr/download.php, `MiND`$_{KL}$, and `DANCo`: http://www.mathworks.it/matlabcentral/fileexchange/40112-intrinsic-dimensionality-estimation-techniques). Note that these estimators cover all the groups described in Section 3, that is, *Projective*, *Fractal*, *Nearest-Neighbors-based*, and *Graph-based* estimators.

The benchmark is composed by following steps:

(1) Test all the considered estimators on both the synthetic and real datasets described below. We highlight that the synthetic datasets whose `id` is a user-defined parameter should be created with sufficiently high `id` values (`id` $\geq 10$).

(2) Comparative evaluation steps are as follows:

    (a) compute the MPE indicator both for synthetic and real datasets,

    (b) compute a ranking test with control methods; to this aim, we suggest the Friedman test with Bonferroni-Dunn post hoc analyses [127],

    (c) perform the tests proposed in [27] to evaluate the robustness, with respect to different cardinalities and parameter settings.

The 21 synthetic datasets used in the benchmark, referred to as $\mathbf{M}_*$ in the following, are listed in Table 2 with their relevant characteristics ($N$, $d$, and $D$). The first 15 datasets are generated with the tool proposed in [79]; they include 4 instances, $\mathbf{M}_{10*}$, of dataset $\mathbf{M}_{10}$, which are drawn from $\mathscr{M}_{10}^H$ after its embedding in $\mathfrak{R}^D$ by setting $D = \{11, 18, 25, 71\}$. Note that we did not include the dataset sampled from $\mathscr{M}_8^H$ (see Table 1) since relevant and recent `id` estimators have similarly produced highly overestimated results when tested on it [107]. Indeed, dealing with highly curved manifolds is still a quite challenging problem in the field.

The last six synthetic datasets are $\mathbf{M}_{N1}$, $\mathbf{M}_{N2}$, $\mathbf{M}_{\text{beta}}$, and 3 instances of dataset $\mathbf{M}_{P*}$, which are sampled from paraboloids $\mathscr{M}_{Pd}$ whose `id` is, respectively, $d = \{3, 6, 9\}$.

To perform multiple tests, 20 instances of each dataset have been generated, and the achieved results have been averaged.

Regarding the real datasets we used the DSVC1 time series [21] ($\mathbf{M}_{\text{DSVC1}}$, `id` $\sim 2.26$), the ISOMAP face database [56] ($\mathbf{M}_{\text{ISOMAP}}$, `id` $= 3$), the `Santa Fe` dataset [121] ($\mathbf{M}_{\text{Santa Fe}}$, `id` $= 9$), the MNIST database [119] ($\mathbf{M}_{\text{MNIST1}}$, `id` $\in \{8, \ldots, 13\}$), the `Isolet` dataset [120] ($\mathbf{M}_{\text{Isolet}}$, `id` $\in \{16, \ldots, 22\}$), and the Crystal Fingerprint space for the chemical compound silicon

TABLE 3: Parameter settings for the different estimators: $k$ represents the number of neighbors, $\gamma$ represents the edge weighting factor for `kNN`, $M$ represents the number of Least Square (LS) runs, $N$ represents the number of resampling trials per LS iteration, $\alpha$ and $\pi$ represent the parameters (shape and rate) of the Gamma prior distributions, which describe the hyperparameters and the observation noise model of `BPCA`, and $\mu$ contains the mean and the precision of the Gaussian prior distribution describing the bias inserted in the inference of `BPCA`.

| Dataset | Method | Parameters |
|---------|--------|------------|
| | MLE | $k_1 = 6, k_2 = 20$ |
| | DANCo | $k = 10$ |
| | kNNG$_1$ | $k_1 = 6, k_2 = 20, \gamma = 1, M = 1, N = 10$ |
| | kNNG$_2$ | $k_1 = 6, k_2 = 20, \gamma = 1, M = 10, N = 1$ |
| Synthetic | BPCA | iters = 2000, $\alpha = (2.0, 2.0)$, $\pi = (2.0, 2.0), \mu = (0.0, 0.01)$ |
| | Hein | *None* |
| | CD | *None* |
| | MLSVD | *None* |
| | MiND$_{KL}$ | $k = 10$ |
| | MLE | $k_1 = 3, k_2 = 8$ |
| | DANCo | $k = 5$ |
| | kNNG$_1$ | $k_1 = 3, k_2 = 8, \gamma = 1, M = 1, N = 10$ |
| | kNNG$_2$ | $k_1 = 3, k_2 = 8, \gamma = 1, M = 10, N = 1$ |
| Real | BPCA | iters = 2000, $\alpha = (2.0, 2.0)$, $\pi = (2.0, 2.0), \mu = (0.0, 0.01)$ |
| | Hein | *None* |
| | CD | *None* |
| | MLSVD | *None* |
| | MiND$_{KL}$ | $k = 5$ |

dioxide $SiO_2$ structure with 3 atoms (this allows us to obtain the $\mathbf{M}_{SiO2}$ dataset containing 4738 points embedded in $\mathfrak{R}^{1800}$ and being characterized by an `id` equal to 12).

To run multiple tests also on $\mathbf{M}_{\text{MNIST1}}$, $\mathbf{M}_{SiO2}$, and $\mathbf{M}_{\text{Isolet}}$, for each of them we generated 5 instances by extracting random subsets containing 2500 points each and we averaged the achieved results.

Table 3 summarizes the parameter values we employed for different estimators. Note that, to relax the dependency of the `kNNG` algorithm from the setting of its parameter $k$, we performed multiple runs with $k_1 \leq k \leq k_2$ and we averaged the achieved results. Furthermore, we tested two versions of the algorithm (referred to as `kNNG`$_1$ and `kNNG`$_2$) obtained by varying the parameters $M$ and $N$.

*4.4. Experimental Results.* Table 4 summarizes the results obtained by the compared estimators on the synthetic datasets, while in Table 5 the results obtained on the real datasets are reported.

Looking at the number of correct estimations computed by each algorithm (highlighted in boldface), we have the following ranking: `MLSVD` is correct on 13 out of 21 synthetic datasets, `DANCo` (correct on 10 out of 21 datasets), `Hein`

TABLE 4: Results achieved on the synthetic datasets. The bottom row reports the MPE achieved by each algorithm; anyhow, for each test dataset the best approximation results are highlighted in boldface.

| Dataset | $d$ | MLE | kNNG$_1$ | kNNG$_2$ | BPCA | Hein | CD | MiND$_{KL}$ | DANCo | MLSVD |
|---|---|---|---|---|---|---|---|---|---|---|
| $M_1$ | 10.00 | 9.10 | 9.16 | 9.89 | 5.45 | 9.45 | 9.12 | 10.30 | 10.09 | **10.00** |
| $M_2$ | 3.00 | 2.88 | 2.95 | 3.03 | **3.00** | **3.00** | 2.88 | **3.00** | **3.00** | **3.00** |
| $M_3$ | 4.00 | 3.83 | 3.75 | 3.82 | **4.00** | **4.00** | 3.23 | **4.00** | **4.00** | 2.08 |
| $M_4$ | 4.00 | 3.95 | 4.05 | 4.76 | 4.25 | **4.00** | 3.88 | 4.15 | **4.00** | 8.00 |
| $M_5$ | 2.00 | 1.97 | 1.96 | 2.06 | **2.00** | **2.00** | 1.98 | **2.00** | **2.00** | **2.00** |
| $M_6$ | 6.00 | 6.39 | 6.46 | 11.24 | 12.00 | **5.95** | 5.91 | 6.50 | 7.00 | 12.00 |
| $M_7$ | 2.00 | 1.96 | 1.97 | 2.09 | **2.00** | **2.00** | 1.93 | 2.07 | **2.00** | 2.35 |
| $M_9$ | 20.00 | 14.64 | 15.25 | 10.59 | 13.55 | 15.50 | 13.75 | 19.15 | 19.71 | **20.00** |
| $M_{10a}$ | 10.00 | 8.26 | 8.62 | 10.21 | 5.20 | 8.90 | 8.09 | 9.85 | 9.86 | **10.00** |
| $M_{10b}$ | 17.00 | 12.87 | 13.69 | 15.38 | 9.46 | 13.85 | 12.30 | 16.25 | 16.62 | **17.00** |
| $M_{10c}$ | 24.00 | 16.96 | 17.67 | 21.42 | 13.3 | 17.95 | 15.58 | 22.55 | 24.28 | **24.00** |
| $M_{10d}$ | 70.00 | 36.49 | 39.67 | 40.31 | 71.00 | 38.69 | 31.4 | 64.38 | 70.52 | **70.00** |
| $M_{11}$ | 2.00 | 2.21 | 1.95 | 2.03 | 1.55 | 2.00 | 2.19 | **2.00** | **2.00** | 1.00 |
| $M_{12}$ | 20.00 | 15.82 | 16.40 | 24.89 | 13.7 | 15.00 | 11.26 | 19.35 | 19.90 | **20.00** |
| $M_{13}$ | 1.00 | **1.00** | 0.97 | 1.07 | 5.70 | **1.00** | 1.14 | **1.00** | **1.00** | **1.00** |
| $M_{N1}$ | 18.00 | 12.25 | 14.26 | 19.8 | 36.00 | 14.10 | 10.40 | 17.76 | 18.76 | **18.00** |
| $M_{N2}$ | 24.00 | 14.72 | 17.62 | 26.87 | 48.00 | 17.76 | 12.43 | 23.76 | 25.76 | **24.00** |
| $M_{beta}$ | 10.00 | 6.36 | 6.45 | 14.77 | 19.7 | 4.00 | 3.05 | 7.00 | 7.00 | **10.00** |
| $M_{P3}$ | 3.00 | 2.89 | 2.93 | 3.12 | 7.00 | 2.00 | 2.43 | **3.00** | **3.00** | 1.00 |
| $M_{P6}$ | 6.00 | 4.96 | 4.98 | 5.82 | 7.00 | 2.66 | 3.58 | 5.04 | **6.00** | 1.00 |
| $M_{P9}$ | 9.00 | 6.35 | 6.89 | 8.04 | 10.95 | 2.85 | 4.55 | 7.00 | **8.00** | 1.00 |
| | MPE | 17.29 | 14.50 | 16.79 | 62.62 | 19.92 | 25.96 | 5.55 | **3.70** | 26.34 |

TABLE 5: Results achieved on the real datasets by the compared approaches. The bottom row reports the MPE achieved by each algorithm; anyhow, for each test dataset the best approximation results are highlighted in boldface (when the real id takes values in a range, we highlighted the results that best approximate the mean value of the range).

| Dataset | id | MLE | kNNG$_1$ | kNNG$_2$ | BPCA | Hein | CD | MiND$_{KL}$ | DANCo | MLSVD |
|---|---|---|---|---|---|---|---|---|---|---|
| $M_{DSCV1}$ | 2.26 | 2.03 | 1.77 | 1.86 | 6.00 | 3.00 | 1.92 | 2.50 | **2.26** | 1.75 |
| $M_{ISOMAP}$ | 3.00 | 4.05 | 3.60 | 4.32 | 4.00 | **3.00** | 3.37 | 3.9 | 4.00 | 1.00 |
| $M_{Santa Fe}$ | 9.00 | 7.16 | 7.28 | 7.43 | 18.00 | 6.00 | 4.39 | 7.60 | **8.19** | 1.00 |
| $M_{MNIST1}$ | 8.00–11.00 | 10.29 | 10.37 | 9.58 | 11.00 | 8.00 | 6.96 | 11.00 | 9.98 | 1.00 |
| $M_{SiO2}$ | 12.00 | 39.28 | 10.24 | 10.36 | 3.00 | 4.80 | 1.05 | 17.20 | **12.60** | 1.00 |
| $M_{Isolet}$ | 16.00–22.00 | 15.78 | 6.50 | 8.32 | 19.00 | 3.00 | 3.65 | 20.00 | 19.00 | 1.00 |
| | MPE | 53.83 | 27.41 | 26.76 | 71.68 | 34.50 | 43.34 | 27.00 | **15.14** | 75.17 |

TABLE 6: Friedman ranking results achieved on all the datasets. The null hypothesis that the algorithms perform comparably is rejected with $p$ value $< 0.00001$.

| Method | Ranking |
|---|---|
| DANCo | 2.40 |
| MiND$_{KL}$ | 3.46 |
| Hein | 4.67 |
| kNNG$_2$ | 5.11 |
| MLSVD | 5.17 |
| kNNG$_1$ | 5.17 |
| MLE | 5.70 |
| CD | 6.63 |
| BPCA | 6.68 |

(correct on 6 out of 21), MiND$_{KL}$ (6 out of 21), BPCA (4 out of 21), and MLE (1 out of 21). It can be further noted that kNNG$_*$,

CD, MLE, and Hein obtain good estimates only for low id manifolds, while they produce underestimated values when processing datasets of high id.

By observing the MPE indicator, which accounts for the precision of the achieved estimates, we obtain a different ranking: DANCo, MiND$_{KL}$, and kNNG$_1$ and kNNG$_2$, MLE, Hein, CD, and MLSVD. This difference is due to the fact that algorithms, such as kNNG$_1$ and kNNG$_2$, MLE, and Hein, most of the times produce results close to the correct value.

Regarding the real datasets, all the algorithms achieve a much worse MPE indicator, and again DANCo is best performing method.

Furthermore, we compute the Friedman ranking test with the Bonferroni-Dunn post hoc analysis as proposed in Section 4.3 to state the quality of the achieved results on both the synthetic and real datasets. Tables 6 and 7 summarize the ranking results.
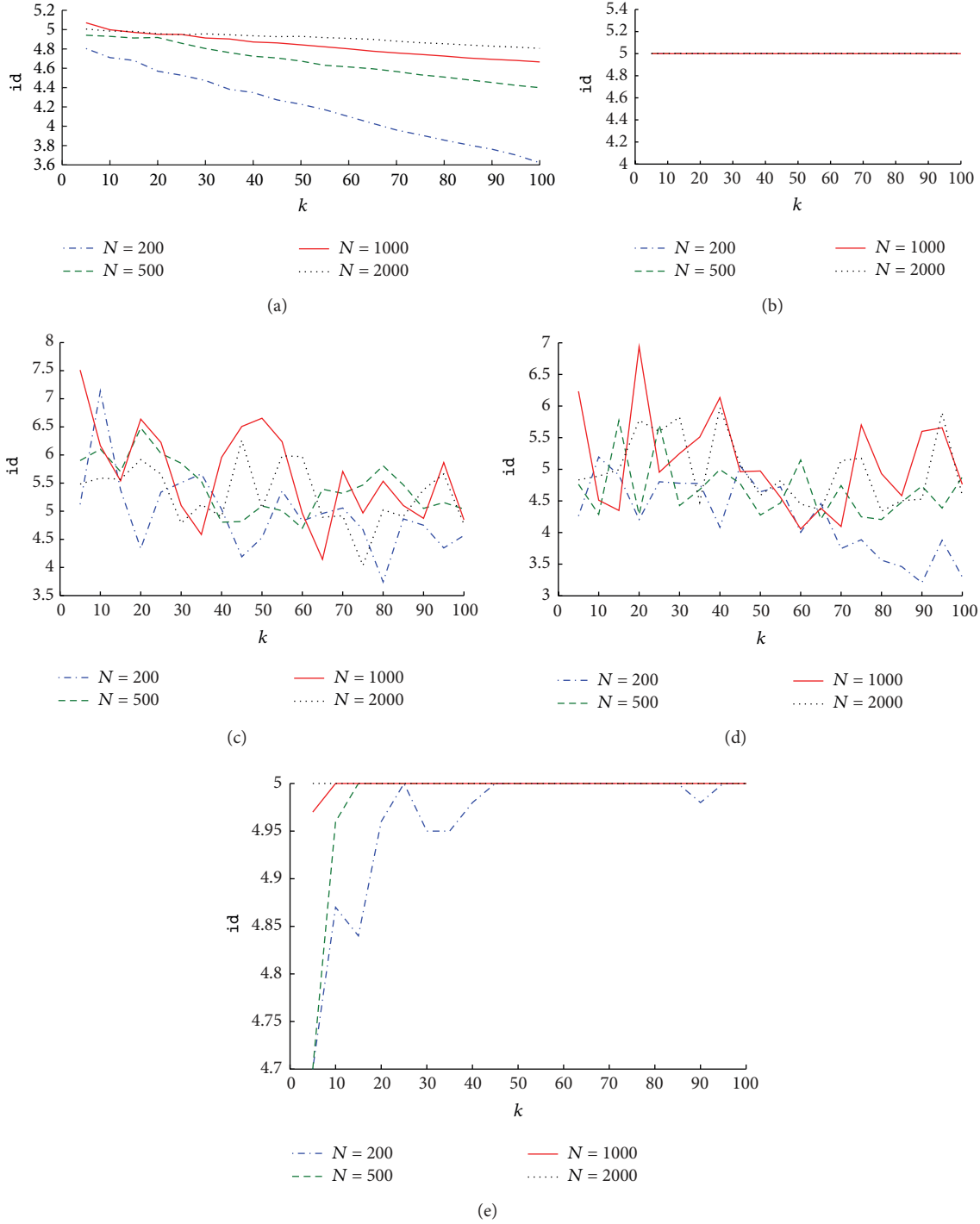
FIGURE 3: Behavior of (a) MLE, (b) DANCo, (c) kNNG$_1$, (d) kNNG$_2$, and (e) MiND$_{KL}$ applied to points drawn from a 5-dimensional standard normal distribution; in this test $N \in \{200, 500, 1000, 2000\}$ and $k \in \{5, \dots, 100\}$.

Finally, we performed the tests proposed in [27] to evaluate the robustness of MiND$_{KL}$, MLE, DANCo, and kNNG$_*$ with respect to the settings of their $k$ parameter. Precisely, these tests employ synthetic datasets subsampled from the standard Gaussian pdf in $\mathfrak{R}^5$ (id = 5). As proposed in Section 4.2, we repeated the tests for datasets with cardinalities

$N \in \{200, 500, 1000, 2000\}$ varying the parameter $k$ in the range $\{5, \dots, 100\}$.

As shown in Figure 3 many of the tested techniques are strongly influenced by the parameter settings; therefore, studying the variability of the algorithms' behavior when changing their parameter settings is of utmost importance.

TABLE 7: Hypothesis testing of significance between techniques. Bonferroni-Dunn's procedure rejects those hypotheses that have a $p$ value $\leq$ 0.0125.

|         | $MiND_{KL}$ | Hein   | $kNNG_1$ | $kNNG_2$ | MLE    | CD     | MLSVD  | BPCA   |
|---------|-------------|--------|----------|----------|--------|--------|--------|--------|
| DANCo   | 0.1567      | 0.0024 | 0.0003   | 0.0002   | 0.0002 | 0.0000 | 0.0000 | 0.0000 |
| $MiND_{KL}$ | * * *   | 0.0801 | 0.0303   | 0.0244   | 0.0055 | 0.0020 | 0.0000 | 0.0000 |
| Hein    | * * *       | * * *  | 0.7528   | 0.6366   | 0.1564 | 0.1474 | 0.0034 | 0.0018 |
| $kNNG_1$ | * * *      | * * *  | * * *    | 0.8557   | 0.3443 | 0.2301 | 0.0164 | 0.0071 |
| $kNNG_2$ | * * *      | * * *  | * * *    | * * *    | 0.9314 | 0.3894 | 0.1113 | 0.0282 |
| MLE     | * * *       | * * *  | * * *    | * * *    | * * *  | 0.3428 | 0.1876 | 0.0307 |
| CD      | * * *       | * * *  | * * *    | * * *    | * * *  | * * *  | 0.7337 | 0.1961 |

## 5. Conclusions and Open Problems

This work presents the base theories of state-of-the-art `id` estimators and surveys the most relevant and recent among them, highlighting their strengths and their drawbacks.

Unfortunately, performing an objective comparative evaluation among the surveyed methods is difficult because, to our knowledge, no benchmark framework exists in this research field; therefore, in Section 4 we propose an evaluation approach that employs both real and synthetic datasets and suggests experiments to evaluate the estimators' robustness with respect to their parameter settings. Note that, the benchmark is designed to evaluate the performance achieved by `id` estimators when both low and high `id` data must be processed; this consideration is due to the fact that, to our knowledge, only few methods [28, 76, 107, 126] have empirically investigated the problem of datasets drawn from manifolds nonlinearly embedded in higher dimensional spaces and characterized by a sufficiently high `id` (i.e., `id` $\geqslant$ 10). However, due to the continuous technological advances, high `id` datasets are becoming more and more common, and the construction of a theoretically well-formed and robust `id` estimator able to deal with high `id` data and limited amount of points remains one of the open research challenges in machine learning. Besides, `id` estimators should be developed by also considering datasets drawn through nonuniform smooth `pdf`s from manifolds $\mathcal{M}$ characterized by a nonconstant curvature; indeed, most of the algorithms are tested by only employing data drawn by means of uniform `pdf`.

We further note that, though the aforementioned problems still need further investigations, most researches in this field are presently focusing on tasks that require to estimate the `id` as the first step. Examples are "multimanifold learning," whose aim is to process datasets drawn from multiple manifolds, each characterized by different `id`, to identify the underlying structures (see [128] for an example); "nonlinear dimensionality reduction"; or "manifold reconstruction," whose aim is to find the mapping that projects the data (embedded in a higher $D$-dimensional space) on the lower $\hat{d}$-dimensional subspace, $\hat{d}$ being the `id` estimated on the input dataset (as examples, see [9, 11, 129]).

## Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

## References

[1] R. S. Bennett, "The intrinsic dimensionality of signal collections," *IEEE Transactions on Information Theory*, vol. 15, no. 5, pp. 517–525, 1969.

[2] C. M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press, Oxford, UK, 1995.

[3] E. Chávez, G. Navarro, R. Baeza-Yates, and J. L. Marroquín, "Searching in metric spaces," *ACM Computing Surveys*, vol. 33, no. 3, pp. 273–321, 2001.

[4] V. Pestov, "An axiomatic approach to intrinsic dimension of a dataset," *Neural Networks*, vol. 21, no. 2-3, pp. 204–213, 2008.

[5] V. Pestov, "Intrinsic dimensionality," *SIGSPATIAL Special*, vol. 2, no. 2, pp. 8–11, 2010.

[6] M. Katetov and P. Simon, "Origins of dimension theory," in *Handbook of the History of General Topology*, vol. 1, 1997.

[7] B. Kégl, "Intrinsic dimension estimation using packing numbers," in *Proceedings of the Neural Information Processing Systems (NIPS '02)*, S. Becker, S. Thrun, and K. Obermayer, Eds., pp. 681–688, MIT Press, 2002.

[8] Z. Zhang and H. Zha, "Adaptive manifold learning," in *Advances in Neural Information Processing Systems*, vol. 17, 2005.

[9] M. Gashler and T. Martinez, "Tangent space guided intelligent neighbor finding," in *Proceedings of the International Joint Conference on Neural Network (IJCNN '11)*, pp. 2617–2624, August 2011.

[10] M. Gashler and T. Martinez, "Robust manifold learning with CycleCut," *Connection Science*, vol. 24, no. 1, pp. 57–69, 2012.

[11] P. Zhang, H. Qiao, and B. Zhang, "An improved local tangent space alignment method for manifold learning," *Pattern Recognition Letters*, vol. 32, no. 2, pp. 181–189, 2011.

[12] N. Verma, "Distance preserving embeddings for general *n*-dimensional manifolds," *Journal of Machine Learning Research*, vol. 14, pp. 2415–2448, 2013.

[13] R. E. Bellman, *Adaptive Control Processes: A Guided Tour*, Princeton University Press, Princeton, NJ, USA, 1961.

[14] M. Kirby, *Geometric Data Analysis: An Empirical Approach to Dimensionality Reduction and the Study of Patterns*, John Wiley & Sons, 2001.

[15] I. T. Jolliffe, *Principal Component Analysis*, Springer Series in Statistics, Springer, New York, NY, USA, 1986.

[16] V. N. Vapnik, *Statistical Learning Theory*, John Wiley & Sons, 1998.

[17] J. H. Friedman, T. Hastie, and R. Tibshirani, *The Elements of Statistical Learning—Data Mining, Inference and Prediction*, Springer, Berlin, Germany, 2009.

[18] P. Campadelli, E. Casiraghi, C. Ceruti, G. Lombardi, and A. Rozza, "Local intrinsic dimensionality based features for

clustering," in *Image Analysis and Processing—ICIAP 2013*, A. Petrosino, Ed., vol. 8156 of *Lecture Notes in Computer Science*, pp. 41–50, Springer, Berlin, Germany, 2013.

[19] P. Grassberger and I. Procaccia, "Measuring the strangeness of strange attractors," *Physica D. Nonlinear Phenomena*, vol. 9, no. 1-2, pp. 189–208, 1983.

[20] H. Lähdesmäki, O. Yli-Harja, W. Zhang, and I. Shmulevich, "Intrinsic dimensionality in gene expression analysis," in *Proceedings of the International Workshop on Genomic Signal Processing and Statistics (GENSIPS '05)*, September 2005.

[21] F. Camastra and M. Filippone, "A comparative evaluation of nonlinear dynamics methods for time series prediction," *Neural Computing and Applications*, vol. 18, no. 8, pp. 1021–1029, 2009.

[22] M. Valle and A. R. Oganov, "Crystal fingerprint space—a novel paradigm for studying crystal-structure sets," *Acta Crystallographica Section A*, vol. 66, no. 5, pp. 507–517, 2010.

[23] K. M. Carter, R. Raich, and A. O. Hero, "On local intrinsic dimension estimation and its applications," *IEEE Transactions on Signal Processing*, vol. 58, no. 2, pp. 650–663, 2010.

[24] J. Lapuyade-Lahorgue and A. Mohammad-Djafari, "Nearest neighbors and correlation dimension for dimensionality estimation. Application to factor analysis of real biological time series data," in *Proceedings of The European Symposium on Artificial Neural Networks (ESANN '11)*, pp. 363–368, Bruges, Belgium, April 2014.

[25] R. Heylen and P. Scheunders, "Hyperspectral intrinsic dimensionality estimation with nearest-neighbor distance ratios," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 6, no. 2, pp. 570–579, 2013.

[26] K. W. Pettis, T. A. Bailey, A. K. Jain, and R. C. Dubes, "An intrinsic dimensionality estimator from near-neighbor information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 1, no. 1, pp. 25–37, 1979.

[27] E. Levina and P. J. Bickel, "Maximum likelihood estimation of intrinsic dimension," in *Proceedings of the NIPS*, vol. 1, pp. 777–784, 2004.

[28] F. Camastra and A. Vinciarelli, "Estimating the intrinsic dimension of data with a fractal-based method," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 10, pp. 1404–1407, 2002.

[29] J. A. Costa and A. O. Hero, "Geodesic entropic graphs for dimension and entropy estimation in manifold learning," *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2210–2221, 2004.

[30] J. A. Costa and A. O. Hero, "Learning intrinsic dimension and entropy of high-dimensional shape spaces," in *Proceedings of the European Signal Processing Conference (EUSIPCO '04)*, pp. 231–252, September 2004.

[31] K. S. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft, "When is 'nearest neighbor' meaningful?" in *Proceedings of the 7th International Conference on Database Theory (ICDT '99)*, pp. 217–235, Springer, London, UK, 1999.

[32] K. M. Carter, R. Raich, W. G. Finn, and A. O. Hero III, "FINE: fisher information nonparametric embedding," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 11, pp. 2093–2098, 2009.

[33] A. M. Farahmand, C. Szepesvári, and J.-Y. Audibert, "Manifold-adaptive dimension estimation," in *Proceedings of the 24th international conference on Machine learning (ICML '07)*, pp. 265–272, June 2007.

[34] J. A. Scheinkman and B. LeBaron, "Nonlinear dynamics and stock returns," *The Journal of Business*, vol. 62, no. 3, pp. 311–337, 1989.

[35] D. R. Chialvo, R. F. Gilmour Jr., and J. Jalife, "Low dimensional chaos in cardiac tissue," *Nature*, vol. 343, no. 6259, pp. 653–657, 1990.

[36] A. Mekler, "Calculation of eeg correlation dimension: large massifs of experimental data," *Computer Methods and Programs in Biomedicine*, vol. 92, no. 1, pp. 154–160, 2008.

[37] G. N. Derry and P. S. Derry, "Age dependence of the menstrual cycle correlation dimension," *Open Journal of Biophysics*, vol. 2, no. 2, pp. 40–45, 2012.

[38] V. Isham, *Statistical Aspects of Chaos: A Review*, Chapman and Hall, London, UK, 1993.

[39] S. Haykin and X. B. Li, "Detection of signals in chaos," *Proceedings of the IEEE*, vol. 83, no. 1, pp. 95–122, 1995.

[40] P. Somervuo, "Speech dimensionality analysis on hypercubical self-organizing maps," *Neural Processing Letters*, vol. 17, no. 2, pp. 125–136, 2003.

[41] B. Hu, T. Rakthanmanon, Y. Hao, S. Evans, S. Lonardi, and E. Keogh, "Towards discovering the intrinsic cardinality and dimensionality of time series using MDL," in *Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence*, vol. 7070 of *Lecture Notes in Computer Science*, pp. 184–197, Springer, Berlin, Germany, 2013.

[42] D. C. Laughlin, "The intrinsic dimensionality of plant traits and its relevance to community assembly," *Journal of Ecology*, vol. 102, no. 1, pp. 186–193, 2014.

[43] F. Camastra, "Data dimensionality estimation methods: a survey," *Pattern Recognition*, vol. 36, no. 12, pp. 2945–2954, 2003.

[44] A. K. Romney, R. N. Shepard, and S. B. Nerlove, *Multidimensionaling Scaling, Volume I: Theory*, Seminar Press, 1972.

[45] A. K. Romney, R. N. Shepard, and S. B. Nerlove, *Multidimensionaling Scaling, Volume II: Applications*, Seminar Press, 1972.

[46] T. Lin and H. Zha, "Riemannian manifold learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 5, pp. 796–809, 2008.

[47] R. N. Shepard, "The analysis of proximities: multidimensional scaling with an unknown distance function. Part I," *Psychometrika*, vol. 27, pp. 125–140, 1962.

[48] R. N. Shepard, "The analysis of proximities: multidimensional scaling with an unknown distance function, part II," *Psychometrika*, vol. 27, pp. 219–246, 1962.

[49] J. B. Kruskal, "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis," *Psychometrika*, vol. 29, pp. 1–27, 1964.

[50] J. B. Kruskal and J. D. Carrol, *Geometrical Models and Badness-of-Fit Functions*, vol. 2, Academic Press, 1969.

[51] R. N. Shepard and J. D. Carroll, *Parametric Representation of Nonlinear Data Structures*, Academic Press, New York, NY, USA, 1969.

[52] J. B. Kruskal, *Linear Transformation of Multivariate Data to Reveal Clustering*, vol. 1, Academic Press, New York, NY, USA, 1972.

[53] C. K. Chen and H. C. Andrews, "Nonlinear intrinsic dimensionality computations," *IEEE Transactions on Computers*, vol. C-23, no. 2, pp. 178–184, 1974.

[54] J. W. J. Sammon, "A nonlinear mapping for data structure analysis," *IEEE Transactions on Computers*, vol. 18, pp. 401–409, 1969.

[55] P. Demartines and J. Hérault, "Curvilinear component analysis: A self-organizing neural network for nonlinear mapping of data sets," *IEEE Transactions on Neural Networks*, vol. 8, no. 1, pp. 148–154, 1997.

[56] J. B. Tenenbaum, V. de Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.

[57] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, 2000.

[58] J. A. Lee and M. Verleysen, *Nonlinear Dimensionality Reduction*, Springer, New York, NY, USA, 2007.

[59] R. Karbauskaite, G. Dzemyda, and E. Mazetis, "Geodesic distances in the maximum likelihood estimator of intrinsic dimensionality," *Nonlinear Analysis: Modelling and Control*, vol. 16, no. 4, pp. 387–402, 2011.

[60] M. Polito and P. Perona, "Grouping and dimensionality reduction by locally linear embedding," *Advances in Neural Information Processing Systems*, vol. 14, pp. 1255–1262, 2001.

[61] B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol. 10, no. 5, pp. 1299–1319, 1998.

[62] K. Fukunaga and D. R. Olsen, "An algorithm for finding intrinsic dimensionality of data," *IEEE Transactions on Computers*, vol. 20, no. 2, pp. 176–183, 1971.

[63] P. J. Verveer and R. P. W. Duin, "An evaluation of intrinsic dimensionality estimators," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 17, no. 1, pp. 81–86, 1995.

[64] J. Brüske and G. Sommer, "Intrinsic dimensionality estimation with optimally topology preserving maps," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 5, pp. 572–575, 1998.

[65] T. Martinetz and K. Schulten, "Topology representing networks," *Neural Networks*, vol. 7, no. 3, pp. 507–522, 1994.

[66] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysis," *Journal of the Royal Statistical Society. Series B. Statistical Methodology*, vol. 61, no. 3, pp. 611–622, 1999.

[67] R. Everson and S. Roberts, "Inferring the eigenvalues of covariance matrices from limited, noisy data," *IEEE Transactions on Signal Processing*, vol. 48, no. 7, pp. 2083–2091, 2000.

[68] C. M. Bishop, "Bayesian PCA," in *Proceedings of the 12th Annual Conference on Neural Information Processing Systems (NIPS '98)*, pp. 382–388, December 1998.

[69] J. J. Rajan and P. J. W. Rayner, "Model order selection for the singular value decomposition and the discrete Karhunen-Loeve transform using a Bayesian approach," *IEE Proceedings—Vision, Image and Signal Processing*, vol. 144, no. 2, pp. 116–123, 1997.

[70] T. P. Minka, "Automatic choice of dimensionality for PCA," Tech. Rep. 514, MIT, 2000.

[71] C. Bouveyron, G. Celeux, and S. Girard, "Intrinsic dimension estimation by maximum likelihood in isotropic probabilistic PCA," *Pattern Recognition Letters*, vol. 32, no. 14, pp. 1706–1713, 2011.

[72] J. Li and D. Tao, "Simple exponential family PCA," in *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS '10)*, pp. 453–460, Sardinia, Italy, May 2010.

[73] Y. Guan and J. G. Dy, "Sparse probabilistic principal component analysis," *Journal of Machine Learning Research*, vol. 5, pp. 185–192, 2009.

[74] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Journal of Computational and Graphical Statistics*, vol. 15, no. 2, pp. 265–286, 2006.

[75] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, NY, USA, 2nd edition, 2006.

[76] C. Ceruti, S. Bassis, A. Rozza, G. Lombardi, E. Casiraghi, and P. Campadelli, "DANCo: an intrinsic dimensionality estimator exploiting angle and norm concentration," *Pattern Recognition*, vol. 47, no. 8, pp. 2569–2581, 2014.

[77] A. V. Little, M. Maggioni, and L. Rosasco, "Multiscale geometric methods for data sets I: multiscale SVD, noise and curvature," MIT-CSAIL-TR 2012-029, 2012.

[78] F. G. Kaslovsky and D. N. Meyer, "Optimal tangent plane recovery from noisy manifold samples," http://xxx.tau.ac.il/abs/1111.4601v2.

[79] M. Hein and J. Y. Audibert, "Intrinsic dimensionality estimation of submanifolds in euclidean space," in *Proceedings of the International Conference on Machine Learning (ICML '05)*, pp. 289–296, 2005.

[80] G. Haro, G. Randall, and G. Sapiro, "Translated poisson mixture model for stratification learning," *International Journal of Computer Vision*, vol. 80, no. 3, pp. 358–374, 2008.

[81] M. Chen, J. Silva, J. Paisley, C. Wang, D. Dunson, and L. Carin, "Compressive sensing on manifolds using a nonparametric mixture of factor analyzers: algorithm and performance bounds," *IEEE Transactions on Signal Processing*, vol. 58, no. 12, pp. 6140–6155, 2010.

[82] L. Brouwer, *Collected Works, Volume I, Philosophy and Foundations of Mathematics and II, Geometry, Analysis, Topology and Mechanics*, North-Holland/American Elsevier, 1976.

[83] I. M. James, *History of Topology*, Mathematics, Elsevier, 1999.

[84] G. Medioni and P. Mordohai, "The tensor voting framework," in *Emerging Topics in Computer Vision*, pp. 191–255, Prentice Hall, 2004.

[85] G. Lombardi, E. Casiraghi, and P. Campadelli, "Curvature estimation and curve inference with tensor voting: a new approach," in *Proceedings of the 10th International Conference on Advanced Concepts for Intelligent Vision Systems (ACIVS '08)*, vol. 5259, pp. 613–624, 2008.

[86] M. Wertheimer, "Untersuchungen zur Lehre von der Gestalt II," in *Psycologische Forshung*, vol. 4, pp. 301–350, 1923, Translation: A Source Book of Gestalt Psychology.

[87] P. Mordohai and G. Medioni, "Dimensionality estimation, manifold learning and function approximation using tensor voting," *Journal of Machine Learning Research*, vol. 11, pp. 411–450, 2010.

[88] J. C. Robinson, *Dimensions, Embeddings, and Attractors*, Cambridge Tracts in Mathematics, Cambridge University Press, 2010.

[89] C.-G. Li, J. Guo, and B. Xiao, "Intrinsic dimensionality estimation within neighborhood convex hull," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 23, no. 1, pp. 31–44, 2009.

[90] K. Falconer, *Fractal Geometry—Mathematical Foundations and Applications*, John Wiley & Sons, 2nd edition, 2003.

[91] N. Tatti, T. Mielikäinen, A. Gionis, and H. Mannila, "What is the dimension of your binary data?" in *Proceedings of the 6th International Conference on Data Mining (ICDM' 06)*, pp. 603–612, December 2006.

[92] J.-P. Eckmann and D. Ruelle, "Fundamental limitations for estimating dimensions and Lyapunov exponents in dynamical

systems," *Physica D. Nonlinear Phenomena*, vol. 56, no. 2-3, pp. 185–187, 1992.

[93] F. Takens, "On the numerical determination of the dimension of an attractor," in *Dynamical Systems and Bifurcations*, B. J. Braaksma, H. W. Broer, and F. Takens, Eds., vol. 1125 of *Lecture Notes in Mathematics*, pp. 99–106, Springer, Berlin, Germany, 1985.

[94] Y. Ashkenazy, "The use of generalized information dimension in measuring fractal dimension of time series," *Physica A: Statistical Mechanics and Its Applications*, vol. 271, no. 3-4, pp. 427–447, 1999.

[95] C. Tricot Jr., "Two definitions of fractional dimension," *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 91, no. 1, pp. 57–74, 1982.

[96] M. R. Brito, A. J. Quiroz, and J. E. Yukich, "Intrinsic dimension identification via graph-theoretic methods," *Journal of Multivariate Analysis*, vol. 116, pp. 263–277, 2013.

[97] M. Raginsky and S. Lazebnik, "Estimation of intrinsic dimensionality using high-rate vector quantization," in *Proceedings of the NIPS*, pp. 1105–1112, 2005.

[98] P. L. Zador, "Asymptotic quantization error of continuous signals and the quantization dimension," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 139–149, 1982.

[99] K. Kumaraswamy, V. Megalooikonomou, and C. Faloutsos, "Fractal dimension and vector quantization," *Information Processing Letters*, vol. 91, no. 3, pp. 107–113, 2004.

[100] G. V. Trunk, "Statistical estimation of the intrinsic dimensionality of a noisy signal collection," *IEEE Transactions on Computers*, vol. 25, no. 2, pp. 165–171, 1976.

[101] M. Fan, H. Qiao, and B. Zhang, "Intrinsic dimension estimation of manifolds by incising balls," *Pattern Recognition*, vol. 42, no. 5, pp. 780–787, 2009.

[102] D. MacKay and Z. Ghahramani, "Comments on maximum likelihood estimation of intrinsic dimension by E. Levina and P. Bickel," 2005, http://www.inference.phy.cam.ac.uk/mackay/dimension/.

[103] M. D. Penrose and J. E. Yukich, "Limit theory for point processes in manifolds," *The Annals of Applied Probability*, vol. 23, no. 6, pp. 2161–2211, 2013.

[104] P. J. Bickel and D. Yan, "Sparsity and the possibility of inference," *Sankhya: The Indian Journal of Statistics*, vol. 70, no. 1, 23 pages, 2008.

[105] M. Das Gupta and T. S. Huang, "Regularized maximum likelihood for intrinsic dimension estimation," in *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI '10)*, P. Grünwald and P. Spirtes, Eds., pp. 220–227, AUAI Press, Catalina Island, Calif, USA, July 2010.

[106] R. Karbauskaite and G. Dzemyda, "Investigation of the maximum likelihood estimator of intrinsic dimensionality," in *Proceedings of the 10th International Conference on Computer Data Analysis and Modeling*, vol. 2, pp. 110–113, 2013.

[107] A. Rozza, G. Lombardi, C. Ceruti, E. Casiraghi, and P. Campadelli, "Novel high intrinsic dimensionality estimators," *Machine Learning*, vol. 89, no. 1-2, pp. 37–65, 2012.

[108] Q. Wang, S. R. Kulkarni, and S. Verdú, "A nearest-neighbor approach to estimating divergence between continuous random vectors," in *Proceedings of the IEEE International Symposium on Information Theory (ISIT '06)*, pp. 242–246, July 2006.

[109] K. V. Mardia, *Statistics of Directional Data*, Academic Press, 1972.

[110] A. J. Quiroz, "Graph-theoretical methods," in *Encyclopedia of Statistical Sciences*, vol. 5, Wiley and Sons, New York, NY, USA, 2006.

[111] A. O. Hero III, B. Ma, O. J. J. Michel, and J. Gorman, "Applications of entropic spanning graphs," *IEEE Signal Processing Magazine*, vol. 19, no. 5, pp. 85–95, 2002.

[112] J. A. Costa, A. Girotra, and A. O. Hero III, "Estimating local intrinsic dimension with k-nearest neighbor graphs," in *Proceedings of the IEEE/SP 13th Workshop on Statistical Signal Processing*, pp. 417–421, July 2005.

[113] J. H. Friedman and L. C. Rafsky, "Graph-theoretic measures of multivariate association and prediction," *Annals of Statistics*, vol. 11, no. 2, pp. 377–391, 1983.

[114] M. D. Penrose and J. E. Yukich, "Central limit theorems for some graphs in computational geometry," *The Annals of Applied Probability*, vol. 11, no. 4, pp. 1005–1041, 2001.

[115] M. R. Brito, A. J. Quiroz, and J. E. Yukich, "Graph-theoretic procedures for dimension identification," *Journal of Multivariate Analysis*, vol. 81, no. 1, pp. 67–84, 2002.

[116] J. M. Steele, L. A. Shepp, and W. F. Eddy, "On the number of leaves of a Euclidean minimal spanning tree," *Journal of Applied Probability*, vol. 24, no. 4, pp. 809–826, 1987.

[117] M. F. Schilling, "Mutual and shared neighbor probabilities: finite- and infinite-dimensional results," *Advances in Applied Probability*, vol. 18, no. 2, pp. 388–405, 1986.

[118] K. Sricharan, R. Raich, and A. O. Hero III, "Optimized intrinsic dimension estimator using nearest neighbor graphs," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '10)*, pp. 5418–5421, Dallas, Tex, USA, March 2010.

[119] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[120] A. Frank and A. Asuncion, *UCI Machine Learning Repository*, UCI, 2010.

[121] F. Pineda and J. Sommerer, "Estimating generalized dimensions and choosing time delays: a fast algorithm," in *Time Series Prediction: Forecasting the Future and Understanding the Past*, pp. 367–385, 1994.

[122] J. A. Costa and A. O. Hero, *Determining Intrinsic Dimension and Entropy of High-Dimensional Shape Spaces*, Birkhäuser, Boston, Mass, USA, 2006.

[123] I. Kivimäki, K. Lagus, I. Nieminen, J. Väyrynen, and T. Honkela, "Using correlation dimension for analysing text data," in *Artificial Neural Networks—ICANN 2010: Proceedings of the 20th International Conference, Thessaloniki, Greece, September 15–18, 2010, Part I*, vol. 6352 of *Lecture Notes in Computer Science*, pp. 368–373, Springer, Berlin, Germany, 2010.

[124] E. Ott, *Chaos in Dynamical Systems*, Cambridge University Press, Cambridge, UK, 1993.

[125] L. O. Chua, M. Komuro, and T. Matsumoto, "The double scroll," *IEEE Transactions on Circuits and Systems*, vol. 32, no. 8, pp. 797–818, 1985.

[126] G. Lombardi, A. Rozza, C. Ceruti, E. Casiraghi, and P. Campadelli, "Minimum neighbor distance estimators of intrinsic dimension," in *Machine Learning and Knowledge Discovery in Databases: Proceedings of the European Conference, ECML PKDD 2011, Athens, Greece, September 5–9, 2011, Part II*, vol. 6912 of *Lecture Notes in Computer Science*, pp. 374–389, Springer, Berlin, Germany, 2011.

[127] J. Jaccard, M. A. Becker, and G. Wood, "Pairwise multiple comparison procedures: a review," *Psychological Bulletin*, vol. 96, no. 3, pp. 589–596, 1984.

[128] D. Gong, X. Zhao, and G. Medioni, "Robust multiple manifolds structure learning," in *Proceedings of the 29th International Conference on Machine Learning (ICML' 12)*, pp. 321–328, July 2012.

[129] J. Wei, H. Peng, Y.-S. Lin, Z.-M. Huang, and J.-B. Wang, "Adaptive neighborhood selection for manifold learning," in *Proceedings of the International Conference on Machine Learning and Cybernetics (ICMLC '08)*, vol. 1, pp. 380–384, IEEE, Kunming, China, July 2008.