

Manifold Dimension Estimation: An Empirical Study

Bi, Zelong and Pierre Lafaye de Micheaux

School of Mathematics & Statistics, University of New South Wales

September 22, 2025

Abstract

The manifold hypothesis suggests that high-dimensional data often lie on or near a low-dimensional manifold. Estimating the dimension of this manifold is essential for leveraging its structure, yet existing work on dimension estimation is fragmented and lacks systematic evaluation. This article provides a comprehensive survey for both researchers and practitioners. We review often-overlooked theoretical foundations and present eight representative estimators. Through controlled experiments, we analyze how individual factors—such as noise, curvature, and sample size—affect performance. We also compare the estimators on diverse synthetic and real-world datasets, introducing a principled approach to dataset-specific hyperparameter tuning. Our results offer practical guidance and suggest that, for a problem of this generality, simpler methods often perform better.

Keywords: manifold hypothesis, intrinsic dimension, nonlinear dimension reduction

1 Introduction

Scientists and engineers are increasingly training and testing models on high-dimensional datasets (Zhou, 2021). Yet, the success of modern machine learning algorithms appears to defy the curse of dimensionality, which posits that the sample size required to reliably recover underlying structures grows exponentially with the data dimension (Keogh and Mueen, 2011). One possible explanation is the *manifold hypothesis*, which assumes that despite sitting within a high-dimensional space \mathbb{R}^p , the data actually lie on or near a lower d -dimensional manifold (Aghajanyan et al., 2020; Meilă and Zhang, 2024), as illustrated in Figure 1. Under this assumption, at least locally or approximately, the dataset can be described using only d free parameters corresponding to the dimension of the manifold — thus greatly simplifying the problem.

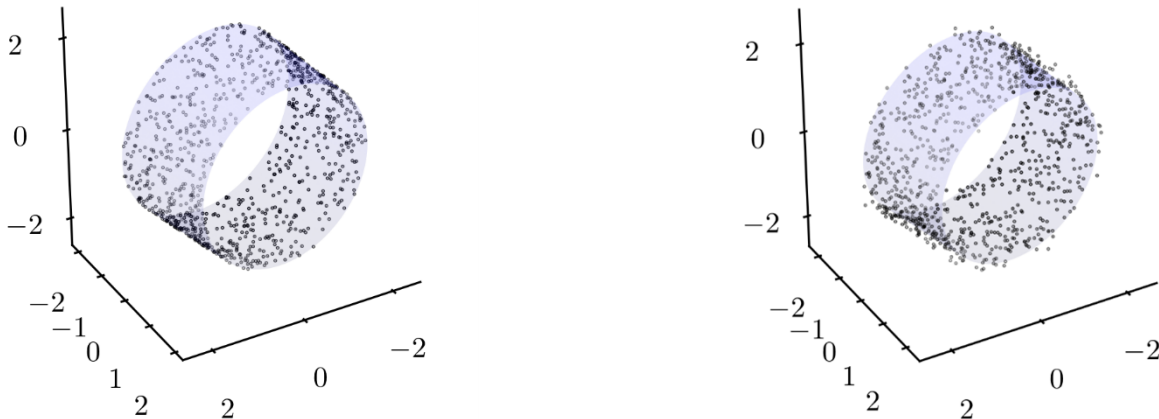


Figure 1: Data points in \mathbb{R}^3 lie on (left) or close to (right) a 2-dimensional manifold (a cylinder).

A natural task that arises from the manifold hypothesis is the *manifold learning problem*, which aims to recover the underlying manifold from the dataset. This task is also referred to as *nonlinear dimension reduction*; see Balasubramanian and Schwartz (2002) and Coifman

and Lafon (2006) for classical methods, and Floryan and Graham (2022) and Alberti et al. (2024) for more recent, data-driven approaches. The most important output of manifold learning is often a low-dimensional representation of the data. This reduced representation is particularly useful because many classical statistical methods do not scale well in high dimensions and can be more effectively applied to the lower-dimensional coordinates. While the integration of manifold learning with traditional statistical analysis remains relatively underexplored, there has been some work in this direction, including regression on manifolds (Aswani et al., 2011; Cheng and Wu, 2013) and factor models for time series (Chan et al., 2017). As high-dimensional data continues to challenge standard statistical techniques, we believe that combining manifold learning with classical methods offers a promising and underutilized alternative.

Among the many properties of the underlying manifold, its dimension is arguably the most important, yet it is often assumed to be known and specified in advance by manifold learning algorithms (Meilă and Zhang, 2024). As a result, researchers have long been interested in the subproblem of *manifold dimension estimation*, and numerous methods have been proposed. But existing work is often fragmented, with researchers approaching the problem from domain-specific perspectives with no unified treatment. Moreover, only limited systematic evaluation of existing techniques have been conducted (such as Campadelli et al. (2015) and Camastra and Staiano (2016), but both are now dated and limited in scope¹), making it difficult for practitioners to make informed choices

This article offers a comprehensive survey of the manifold dimension estimation problem and is organized as follows. Section 2 introduces essential theoretical background in an accessible language, often unfamiliar to practitioners and typically omitted from previ-

¹During the final stage of writing this article, we became aware of the arXiv preprint by Binnie et al. (2025), which also surveys dimension estimators. However, our work was conducted independently and approaches the problem from different angles, offering unique insights and experimental rigor.

ous surveys. We notably cover manifolds and their constructions, sampling procedures, and the statistical difficulty of the dimension estimation problem. Section 3 details the mathematical construction of eight estimators of the manifold dimension. We chose these specific estimators with care, considering top performers identified in prior work, as well as promising more recent proposals. Section 4 presents our empirical evaluation, including the experimental setup, results, and discussion. We benchmark the performance of these estimators under varying conditions, including among things neighborhood size, curvature, and noise level. While earlier comparative studies provided valuable insights, they did not typically control for these factors individually. We also compare the eight estimators on 18 synthetic data sets and on 3 real-world datasets with different geometry, distribution and noise characteristics. Building on insights from our controlled experiments, we develop an algorithm to tune hyperparameters per dataset for the estimators, addressing a common flaw in prior studies where hyperparameters were fixed across datasets. This leads to more equitable comparisons and more reliable conclusions. Finally, Section 5 concludes with practical recommendations for practitioners and outlines directions for future research. We conduct our study using both `Python` and `R`, and all code and datasets used in our experiments are publicly available.

2 Theoretical Background

In this section, we provide a concise yet comprehensive overview of the theoretical foundations relevant to manifold dimension estimation. We begin with a review of manifolds and their associated structures, then discuss the problem of sampling from manifolds. Finally, we formulate the manifold dimension estimation problem and outline its statistical challenges.

Manifolds—and the mathematical constructions built upon them—can become quite ab-

stract when developed in full generality. Such abstraction, however, is unnecessary for our purposes. Throughout, we adopt the most direct and practical definitions, while avoiding oversimplification. For readers interested in more advanced treatments, we refer to standard references in differential geometry, such as [Lee \(2003\)](#) and [Lee \(2018\)](#).

2.1 Manifold and Dimension

In this article, we view a d -dimensional manifold M simply as a smooth² d -dimensional object, assumed to be sitting inside some higher-dimensional Euclidean space \mathbb{R}^p , called the *ambient space*. In that case, the Euclidean structure of \mathbb{R}^p naturally induces a way to measure distances and angles along M , leading to what is formally called a Riemannian manifold. Intuitively, we can think of M as a geometric object with a flat or curved smooth surface, that when you zoom in closely enough, looks like \mathbb{R}^d . When $d = 1$, a manifold is a curve (e.g., a circle or a helix) and locally resembles a segment of a straight line. When $d = 2$, it appears as a surface (e.g., a sphere or a torus) and locally resembles a plane. See [Figure 2](#) for an illustration of some classical examples. A manifold therefore generalizes the idea of a smooth curve or surface to any dimension d ³.

²Formally, a manifold is smooth if it can be fully described by a collection of (local) coordinate maps (called chart functions), and whenever two chart functions describe the same region (i.e., their domains overlap), the change of variables between their coordinate systems is infinitely differentiable.

³In this work, we restrict attention to manifolds without boundary, excluding objects that have both interior and boundary points (such as closed balls). From the perspective of manifold dimension estimation this distinction is not important, since the probability of sampling boundary points from a continuous distribution is zero. Nevertheless, under the standard definition used here, the boundary of a smooth object cannot itself be modeled as part of the manifold and is therefore assumed to be empty.

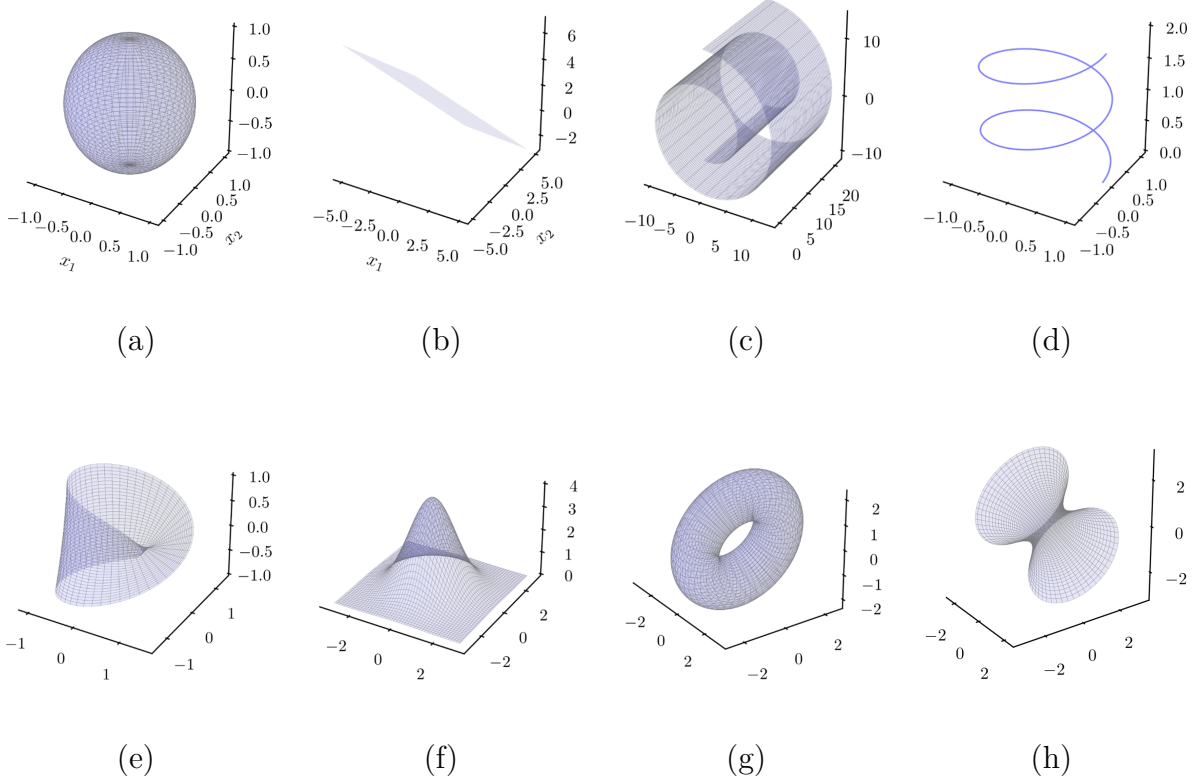
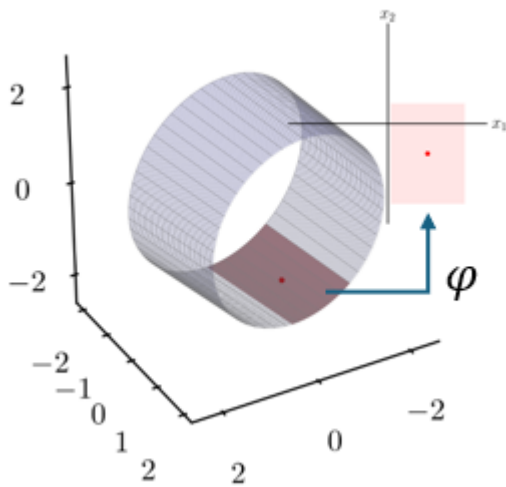


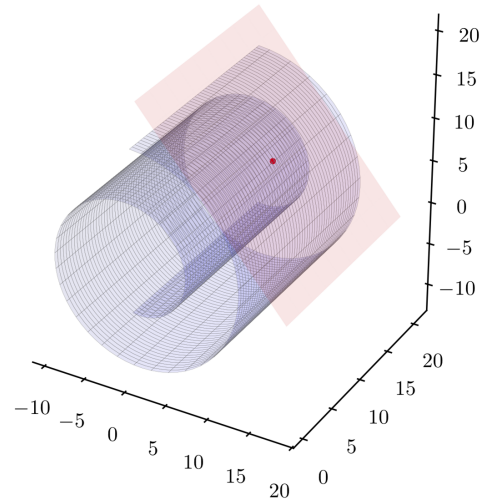
Figure 2: Common manifolds embedded in \mathbb{R}^3 : (a) a 2-sphere; (b) a hyperplane; (c) a Swiss roll with no boundary; (d) a helix with no endpoints; (e) a Möbius strip with no boundary; (f) a Gaussian density surface; (g) a torus; (h) a hyperboloid of one sheet.

Real-world datasets do not necessarily lie on any of the classical manifolds illustrated in Figure 2. The strength of the manifold model is that it can describe much more general geometric structures beyond these familiar examples. What distinguishes a manifold M from an arbitrary subset of \mathbb{R}^p is that it possesses certain local (and global) properties. The most important of these is that M is *locally Euclidean*: every point $\mathbf{x}_0 \in M$ has a neighborhood $U \subseteq M$ which can be described using only d coordinates via a *chart function*⁴, as illustrated in Figure 3a.

⁴A chart function is a map φ that assigns to each point in a neighborhood on the manifold a unique point in \mathbb{R}^d . Chart functions provide a local coordinate system for each “small region” of the manifold, and the compatibility of these maps is what ensures the manifold is well defined.



(a) Chart at one point \mathbf{x}_0 .



(b) Tangent space at one point \mathbf{x}_0 .

Figure 3: Chart and tangent space.

In other words, although M sits inside the higher-dimensional space \mathbb{R}^p (and hence $\mathbf{x}_0 \in \mathbb{R}^p$), it locally (i.e., for all points of M around \mathbf{x}_0) resembles \mathbb{R}^d and can be parameterized by only d free variables. We call this integer d the *intrinsic dimension* of M , and p its *ambient dimension*. In practice, p is known (as this is the dimension of the variables or features in the dataset), and d usually needs to be estimated from the sample.

Beyond the commonly adopted notion of intrinsic dimension d introduced above, several alternative definitions of dimension exist. Some manifold dimension estimators are based on these alternative notions (e.g., [Camastra and Vinciarelli \(2002\)](#), [Bruske and Sommer \(2002\)](#), [Bac and Zinovyev \(2020\)](#)) so we very briefly introduce them below. Notable examples include:

1. *Topological dimension* ([Fedorchuk, 1990](#)), also known as the *Lebesgue covering dimension*, is defined for general topological spaces and remains invariant under homeomorphisms.

2. *Box-counting dimension* (Falconer, 2013), also called the *Minkowski dimension*, is obtained by analyzing the asymptotic behavior of the number of boxes (sets with equal-length interval sides) required to cover a set. It is widely used in the study of fractals.
3. *Hausdorff dimension* (Falconer, 2013) is defined in terms of the Hausdorff measure (to be briefly introduced next). It applies to metric spaces and plays a central role in fractal geometry as well.

These definitions arise from different contexts and motivations, and are not tailored to manifolds in particular. Nevertheless, for sufficiently “regular” spaces like manifolds, these different notions of dimension coincide.

2.2 Tangent Space and Curvature

We have seen in the previous section that a manifold M is locally Euclidean. This implies that it can be well approximated, near any point $\mathbf{x}_0 \in M$ (i.e., for all points in a small neighborhood $U \subseteq M$ of \mathbf{x}_0), by a flat d -dimensional space (e.g., a plane if $d = 2$) that “touches” the manifold at \mathbf{x}_0 . This d -dimensional space is called the *tangent space* at \mathbf{x}_0 , denoted as $T_{\mathbf{x}_0}M$. For instance, if M is the unit sphere in \mathbb{R}^3 and \mathbf{x}_0 is some point on the sphere, then $T_{\mathbf{x}_0}M$ is the plane tangent to the sphere at \mathbf{x}_0 . See Figure 3b for an illustration of another example with a Swiss roll in \mathbb{R}^3 .

Formally, let $\varphi : U \rightarrow \varphi(U) \subseteq \mathbb{R}^d$ be a chart function (also called a coordinate map since it assigns coordinates to points of U) defined on a small neighborhood U of $\mathbf{x}_0 \in \mathbb{R}^p$. Its inverse $\varphi^{-1} : \varphi(U) \rightarrow U \subseteq \mathbb{R}^p$ parametrizes the neighborhood U in the sense that, given a vector of only d parameters (coordinates) $\mathbf{u} = (u_1, \dots, u_d)^\top \in \varphi(U) \subseteq \mathbb{R}^d$, it returns the corresponding point $\varphi^{-1}(\mathbf{u}) \in \mathbb{R}^p$ of the manifold. The Jacobian matrix of φ^{-1} evaluated

at $\varphi(\mathbf{x}_0) \in \mathbb{R}^d$

$$d\varphi^{-1}|_{\varphi(\mathbf{x}_0)} := \left[\frac{\partial \varphi_i^{-1}}{\partial u_j} \right]_{\varphi(\mathbf{x}_0)} \in \mathbb{R}^{p \times d},$$

has full rank d and defines an injective linear map from \mathbb{R}^d into \mathbb{R}^p . (Throughout, the letter d denotes the differential operator and should not be confused with the letter d , which designates the intrinsic dimension of the manifold M .) Thus the tangent space to M at \mathbf{x}_0 is defined through this Jacobian matrix as

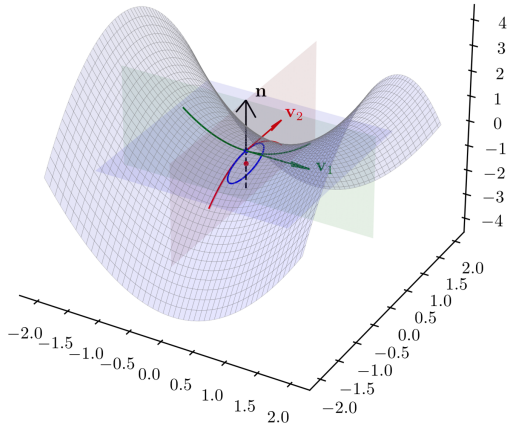
$$T_{\mathbf{x}_0}M = \left\{ \mathbf{x}_0 + d\varphi^{-1}|_{\varphi(\mathbf{x}_0)} \mathbf{u} : \mathbf{u} \in \mathbb{R}^d \right\}.$$

The expression used in the definition is exactly the first-order Taylor expansion of φ^{-1} at $\varphi(\mathbf{x}_0)$, which explains why the tangent space gives the *best linear approximation* to the manifold near that point. The existence of this tangent space is of crucial importance in practice. Indeed, many manifold dimension estimators rely on the fact that small neighborhoods of points on M can be approximated by their tangent spaces, reducing a nonlinear geometric problem to a local linear one.

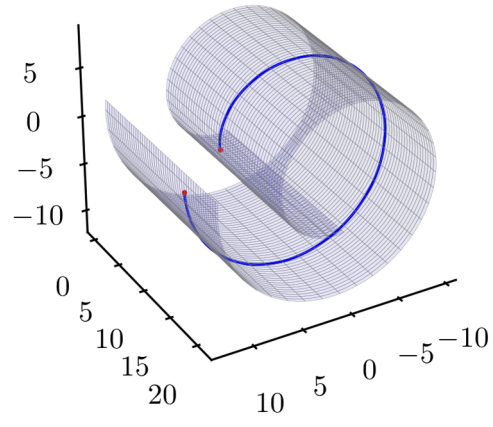
For trivial manifolds like \mathbb{R}^p or its subspaces, the tangent space at any point coincides exactly with the manifold itself; i.e., the linear approximation is perfect everywhere and such manifolds are said to be *flat*. More general manifolds are not flat and can curve non-trivially, the farther we move away from \mathbf{x}_0 , the less accurate the approximation becomes. For curved manifolds, *principal curvatures* (Lee, 2018) can be used to quantify how a manifold bends at various points. This is illustrated in Figure 4a: each direction (e.g., \mathbf{v}_1 or \mathbf{v}_2) in the tangent space $T_{\mathbf{x}_0}M$ (in light blue) can be combined with a normal direction (e.g., \mathbf{n}) to form a plane (in light green or light red) whose intersection with the manifold determines a curve (in dark green or dark red). The principal curvature is then defined as the amount by which this curve bends at \mathbf{x}_0 .⁵ Since there are d linearly independent

⁵The magnitude of the principal curvature is defined as the inverse radius of the osculating circle, which is the unique circle that, when parametrized with unit speed, shares the same position, velocity, and

directions in the tangent space and $p - d$ normal directions, a total of $d(p - d)$ principal curvatures are needed to fully describe the local geometry at \mathbf{x}_0 .



(a) Principal curves (in dark green and red) and osculating circle whose inverse radius gives the curvature of the red curve at \mathbf{x}_0 .



(b) Geodesic (i.e., curve of minimal length) in blue between two points (in red) on a manifold.

Figure 4: Principal curves and their curvature (left) and a geodesic of minimal distance along a manifold (right).

2.3 Geodesics and Distance on a Manifold

By construction, a manifold is also a metric space, i.e., a set together with a notion of distance between its elements. Distances *along the manifold* between two points are measured via curves (of minimal length) that sit on the manifold. These curves are called *geodesics*. They generalize the notion of straight lines in Euclidean spaces. It is important to note that the intrinsic distance between two points on a manifold can differ substantially from their Euclidean distance in the ambient space \mathbb{R}^p ; points that appear close in \mathbb{R}^p may in fact be far apart along the manifold, as illustrated in Figure 4b.

acceleration as the curve at \mathbf{x}_0 .

2.4 Local Graph Representation

Another way to characterize a manifold $M \subseteq \mathbb{R}^p$ is through its *local graph representation*, which expresses M locally as a collection of function graphs. For example, the 2-dimensional unit sphere in \mathbb{R}^3 is defined as the set of triples (x_1, x_2, x_3) satisfying $x_1^2 + x_2^2 + x_3^2 = 1$. In a small neighborhood (say, on the upper hemisphere), this relation can then be used to solved for x_3 , giving

$$x_3 = g(x_1, x_2) := \sqrt{1 - x_1^2 - x_2^2}.$$

Thus, in this neighborhood the sphere coincides with the graph of the function g .

This idea extends to any d -dimensional manifold M . By Proposition 5.16 of [Lee \(2003\)](#), together with the inverse function theorem, for any point $\mathbf{x}_0 \in M$ there exists a neighborhood $U \subseteq M$ in which each $\mathbf{x} \in U$ can be expressed in the new coordinates satisfying

$$\mathbf{x}' = (x'_1, \dots, x'_d, g(x'_1, \dots, x'_d)),$$

with g a smooth function satisfying $g(0) = 0$ and $\nabla g(0) = 0$. In other words, U is precisely the graph of g . The first d new coordinates (x'_1, \dots, x'_d) correspond to directions in the tangent space $T_{\mathbf{x}_0}M$, while the remaining $p - d$ new coordinates are fully determined by g . Setting these last $p - d$ coordinates to zero corresponds to the first-order Taylor approximation to g , which is equivalent to the tangent space approximation for M at \mathbf{x}_0 . The curvature of the manifold at \mathbf{x}_0 is fully encoded in the higher-order terms of the Taylor expansion of g .

2.5 Embedding

The major justification for the manifold hypothesis is the celebrated Nash embedding theorem ([Nash, 1956](#)), which guarantees that any manifold can be embedded into some Euclidean space in a smooth and faithful way, preserving angles and distances. This naturally

raises the question of how small the ambient dimension can be while still faithfully representing the manifold.

For instance, all the 2-dimensional manifolds shown in Figure 2 of Section 2.1 are embedded in \mathbb{R}^3 , where the ambient and intrinsic dimensions are close. We can easily embed them into higher dimensional spaces through linear embeddings with their geometries preserved, but as long as the manifolds remain in some 3-dimensional subspace, the manifold dimension estimation problem on them should remain relatively easy.

A more challenging scenario is when a d -dimensional manifold is nontrivially embedded into \mathbb{R}^p . For example, consider the d -dimensional *deformed sphere*, defined through an inverse chart function $\varphi^{-1} : (-\pi, \pi)^d \rightarrow \mathbb{R}^p$ with $p = 2d$ as

$$\varphi^{-1}(\mathbf{u}) = (x_1, \dots, x_{2d}),$$

with

$$x_j = [R + r \cos(2c\pi u_j)] \cos(2\pi u_j) \text{ and } x_{j+d} = [R + r \cos(2c\pi u_j)] \sin(2\pi u_j)$$

for $j = 1, \dots, d$ and $R, r, c > 0$ fixed constants. By definition, a d -dimensional deformed sphere resides in \mathbb{R}^{2d} and is not contained in any lower-dimensional subspaces⁶. Figure 5 shows the deformed sphere⁷ when $d = 1$ for different values of c . The deformed sphere is special in the sense that all its d^2 principal curvatures are nontrivial. Manifolds like the deformed sphere are said to be space-filling or nonlinearly embedded in some literature (Campadelli et al., 2015).

⁶This is evident from the definition, a formal justification is also given in the Appendix.

⁷For $d = 1$ and $c = 0.1$, the curve is not closed when plotted over $u \in (-\pi, \pi)$, since the angular part has period 1 while the radial modulation $\cos(2\pi cu)$ has period $1/c = 10$. The curve closes only after u advances by 10, hence the apparent non-closed “deformed circle” in panel (b).

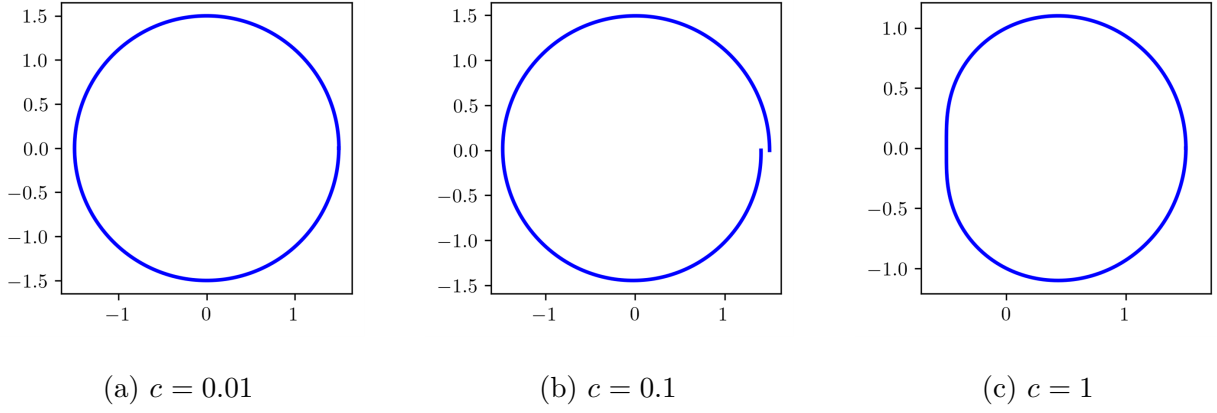


Figure 5: Three deformed spheres ($d = 1$).

2.6 Probability Distributions on Manifolds

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let \mathbf{X} be a random vector defined on it (taking values from \mathbb{R}^p). The manifold hypothesis can be restated by saying that the data distribution is supported entirely on a manifold M , i.e., $\mathbb{P}(\mathbf{X} \in M) = 1$. Moreover, it is often convenient to assume that this distribution admits a density function $f_{\mathbf{X}}$ supported on M .

A few technical challenges must be addressed for a rigorous formulation. The most important is that when $d < p$, the standard Lebesgue measure on \mathbb{R}^p is not appropriate for defining the ‘size’ of the manifold, and one must instead use measures adapted to lower-dimensional sets.

It turns out that the appropriate measure here is the d -dimensional Hausdorff measure, denoted by ν ; see [Diaconis et al. \(2013\)](#) for further details. Assume that the probability distribution $P_{\mathbf{X}}$ on the manifold admits a density function $f_{\mathbf{X}}$ with respect to ν . Given a chart function φ with domain U , one can show that

$$\int_{\varphi(U)} f_{\mathbf{X}}(\varphi^{-1}(\mathbf{u})) J_d \varphi^{-1}(\mathbf{u}) d\mathbf{u} = \int_U f_{\mathbf{X}} d\nu,$$

where

$$J_d \varphi^{-1}(\mathbf{u}) := \sqrt{\det(d\varphi^{-1}(\mathbf{u})^\top d\varphi^{-1}(\mathbf{u}))}.$$

Here $J_d\varphi^{-1}(\mathbf{u})$ plays the role of the *volume distortion factor*: it generalizes the Jacobian determinant from the classical change-of-variables formula, adjusting for how the map φ^{-1} stretches or compresses d -dimensional volumes in the ambient space.

The above result implies that sampling points from $\varphi(U) \subseteq \mathbb{R}^d$ according to the density $f_{\mathbf{X}}(\varphi^{-1}(\mathbf{u})) J_d\varphi^{-1}(\mathbf{u})$, and then mapping them back to M via φ^{-1} , is equivalent in distribution to sampling directly from M with density $f_{\mathbf{X}}$ relative to ν .

The above result, together with the fact that any compact manifold can be covered by finitely many disjoint neighborhoods (up to a set of probability zero), leads to the rejection sampling algorithm below for generating random samples from a prescribed distribution on M . It is worth noting, however, that more efficient sampling methods are available in specific cases. For example, to sample uniformly from the d -dimensional unit sphere in \mathbb{R}^{d+1} , one can draw a vector from $N(0, I_{d+1})$ and normalize it ([Marsaglia, 1972](#)).

Algorithm 1 Sampling from a density function $f_{\mathbf{X}}$ supported on a manifold M

1: **Input:**

- $f_{\mathbf{X}}$: a density function supported on M
- $\{(U_\alpha, \varphi_\alpha, p_\alpha)\}_{1 \leq \alpha \leq a}$: a partition of M into charts, with chart functions φ_α , and region probabilities p_α , with $\alpha, a \in \mathbb{N}$
- n : sample size

2: **Output:** a random sample of size n from $f_{\mathbf{X}}$

3: For each α , find a constant M_α such that for all $\mathbf{u} \in \varphi_\alpha(U_\alpha)$,

$$f_{\mathbf{X}}(\varphi_\alpha^{-1}(\mathbf{u})) \cdot J_d \varphi_\alpha^{-1}(\mathbf{u}) \leq M_\alpha.$$

4: Initialize $k = 0$

5: **while** $k < n$ **do**

6: Sample $\alpha \in \{1, \dots, a\}$ according to the distribution $\{p_\alpha\}$

7: **repeat**

8: Sample \mathbf{u} uniformly from $\varphi_\alpha(U_\alpha)$

9: Compute acceptance probability:

$$p_{\text{accept}} = \frac{f_{\mathbf{X}}(\varphi_\alpha^{-1}(\mathbf{u})) \cdot J_d \varphi_\alpha^{-1}(\mathbf{u})}{M_\alpha}$$

10: Sample $r \sim \text{Uniform}[0, 1]$

11: **until** $r < p_{\text{accept}}$

12: Set $\mathbf{x}_{k+1} = \varphi_\alpha^{-1}(\mathbf{u})$

13: $k \leftarrow k + 1$

14: **end while**

15: **return** $\mathbf{x}_1, \dots, \mathbf{x}_n$

2.7 Manifold Dimension Estimation

With all the constructions at our disposal, we state the manifold dimension estimation problem as below.

Let $\mathbf{x}_{1:n} := \{\mathbf{x}_k\}_{k=1}^n \subseteq \mathbb{R}^p$ be an observed random sample from a probability distribution with density $f_{\mathbf{X}}$ supported on a d -dimensional manifold M . The problem of manifold dimension estimation is to estimate the intrinsic dimension d by means of an estimator $\hat{d}(\mathbf{x}_{1:n})$.

[Kim et al. \(2019\)](#) analyzed the manifold dimension estimation problem by characterizing its fundamental difficulty in terms of minimax risk. Specifically, consider a random sample of size n in \mathbb{R}^p , drawn from a probability distribution $\mathbb{P} \in \mathcal{P}$, where \mathcal{P} is a class of distributions with bounded densities supported on well-behaved manifolds of intrinsic dimension $d \leq p$. Let \hat{d}_n be an estimator constructed from the sample, and let $d(\mathbb{P})$ denote the intrinsic dimension of the manifold on which \mathbb{P} is supported. The minimax risk is defined as

$$R_n = \inf_{\hat{d}_n} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}} \left[\mathbf{1}_{\{\hat{d}_n \neq d(\mathbb{P})\}} \right],$$

which represents the smallest possible worst-case probability of misestimating the intrinsic dimension.

The authors show that when the manifolds under consideration are compact, complete, and have bounded curvature, the minimax risk satisfies

$$-c_1 n \leq \log R_n \leq -c_2 n,$$

for some constants $c_1, c_2 > 0$. This exponential decay implies that the probability of error decreases at a rate proportional to e^{-cn} , so the dimension estimation problem becomes exponentially easier as the sample size increases.

3 Manifold Dimension Estimators

Over the years, researchers from diverse domains have proposed numerous manifold dimension estimators (e.g., [Verveer and Duin \(1995\)](#); [Kégl \(2002\)](#), [Yang et al. \(2007\)](#); [Carter et al. \(2009\)](#); [Sricharan et al. \(2010\)](#), [Kalantana and Einbecka \(2013\)](#); [Liao et al. \(2014\)](#); [Facco et al. \(2017\)](#); [Lim et al. \(2024\)](#)). Given the sheer volume and variety of these methods, it is neither feasible nor the intention of this article to provide an exhaustive review. Instead, in this section, we present eight representative estimators included in our empirical study. These methods reflect a wide range of design principles, from top performers identified from previous surveys ([Campadelli et al., 2015](#); [Bac et al., 2021](#)) to more recent, promising approaches that have yet to be extensively evaluated. For clarity, we organize them into two broad categories based on their underlying design philosophies.

The majority of manifold dimension estimators belong to the class of *flatness-based estimators*. These methods rely on the assumption that the manifold does not curve substantially and can therefore be locally approximated as flat within sufficiently small neighborhoods. As discussed in the previous section, for any $\mathbf{x}_k \in M$, the tangent space $T_{\mathbf{x}_k} M$ provides the best local linear approximation of the manifold, with deviations attributable to curvature. Flatness-based estimators model \mathbf{x}_k together with its K nearest neighbors as if they were uniformly distributed within a d -dimensional ball of radius $R > 0$, $B_{\mathbf{x}_k}(R)$, embedded in the tangent space centered at \mathbf{x}_k . Under this *flatness assumption*, statistics computed from \mathbf{x}_k and its neighbors can be expressed as functions of the intrinsic dimension d , which motivates a wide range of estimation procedures.

Our first category therefore comprises six representative flatness-based estimators: **Local PCA**, **MADA**, **MLE**, **DanCo**, **TLE**, and **TwoNN**. This list is not exhaustive: several older but well-known estimators also fall into this category, including **CC** ([Grassberger and Procaccia, 1983](#); [Kalantana and Einbecka, 2013](#)), **kNN** ([Costa and Hero, 2004a,b](#)),

and **Hein** (Hein and Audibert, 2005), along with a number of lesser-known approaches (Bennett, 1969; Trunk, 1968; Johnsson et al., 2014; Thordsen and Schubert, 2022). Many of these have already been reviewed in earlier surveys (Verveer and Duin, 1995; Campadelli et al., 2015; Camastra and Staiano, 2016).

Estimators that do not rely on the flatness assumption are classified into the second category. In this study, we focus on two recent, promising methods: **CA-PCA** and **Wasserstein**. Other well-known methods in this category include **ISOMAP** (Balasubramanian and Schwartz, 2002), which seek a low-dimensional embedding into the Euclidean space that preserves geometric properties such as geodesic distances, and **Diffusion Maps** (Coifman and Lafon, 2006), which construct a transition matrix that captures the manifold’s intrinsic geometry. These two algorithms, however, are primarily designed for the broader task of recovering the full manifold structure rather than solely estimating its dimension, though they can be readily adapted for this purpose.

Throughout this section, let $\{\mathbf{x}_k\}_{k=1}^n \subseteq \mathbb{R}^p$ denote an observed sample drawn from a density function $f := f_{\mathbf{X}}$ supported on an underlying d -dimensional manifold M . For each observation \mathbf{x}_k , we denote $(\mathbf{x}_k^1, \dots, \mathbf{x}_k^K)$ the ordered (from closest to furthest) K -tuple of its nearest neighbors in Euclidean distance. We denote $\|\mathbf{x}\|$ the Euclidean norm of a vector $\mathbf{x} \in \mathbb{R}^p$.

3.1 Flatness-based Estimators

3.1.1 The Local PCA estimator

Local PCA (Fukunaga and Olsen, 1971; Fan et al., 2010) applies principal component analysis (PCA) to each local neighborhood. For the neighborhood of \mathbf{x}_k , let $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$

denote the eigenvalues of the sample covariance matrix

$$\hat{\Sigma}_k = \frac{1}{K} \sum_{\ell=0}^K (\mathbf{x}_k^\ell - \bar{\mathbf{x}}_k)(\mathbf{x}_k^\ell - \bar{\mathbf{x}}_k)^\top,$$

where $\mathbf{x}_k^0 = \mathbf{x}_k$ and

$$\bar{\mathbf{x}}_k = \frac{1}{K+1} \sum_{\ell=0}^K \mathbf{x}_k^\ell.$$

Under the flatness assumption, the smallest $p - d$ eigenvalues of the population covariance matrix Σ_k vanish. Accordingly, we expect the empirical covariance $\hat{\Sigma}_k$ to approximate this structure, with $\hat{\lambda}_{d+1}, \dots, \hat{\lambda}_p$ close to zero. In practice, each neighborhood yields an estimate \hat{d}_k (expected to be close to d), defined as the number of eigenvalues judged to be significantly nonzero, and the corresponding eigenvectors span an “estimate” of the tangent space $T_{\mathbf{x}_k} M$. A global estimate is then obtained by averaging over all N neighborhoods:

$$\hat{d} = \frac{1}{N} \sum_{k=1}^N \hat{d}_k.$$

This method is straightforward, computationally fast, and widely used in practice. However, many studies on manifold dimension estimators do not compare against it. A commonly cited reason is the need to select a threshold for determining which eigenvalues are significant, which introduces some arbitrariness. A typical rule, proposed in (Fukunaga and Olsen, 1971) is to take the index of the smallest eigenvalue which exceeds 5% of the largest eigenvalue, i.e.,

$$\hat{d}_k = \operatorname{argmax}_{1 \leq q \leq p} \hat{\lambda}_q > 0.05 \hat{\lambda}_1.$$

Our experiments show that with this simple rule, **Local PCA** often outperforms most alternative estimators.

Several methods have been proposed to avoid the subjective threshold choice. For example, Bouveyron et al. (2011) impose an isotropic variance structure within each neighborhood, which yields a consistent maximum likelihood estimator of the local dimension. More

recently, [Lim et al. \(2024\)](#) proposed an alternative criterion:

$$\hat{d}_k = \operatorname{argmin}_{1 \leq q \leq p} \sum_{j=1}^p \left(\frac{\hat{\lambda}_j}{R^2} - \frac{1}{j+2} 1_{\{j \leq q\}} \right)^2,$$

motivated by the fact that if \mathbf{x}_k and its neighbors are uniformly distributed in $B_{\mathbf{x}_k}(1)$, then the population covariance has its first d eigenvalues equal to $(d+2)^{-1}$ and the remaining $p-d$ eigenvalues equal to zero. Here the scaling factor R is estimated as the distance between \mathbf{x}_k and its K -th nearest neighbor \mathbf{x}_k^K .

3.1.2 The MADA estimator

The manifold-adaptive dimension estimator (**MADA**) proposed by [Farahmand et al. \(2007\)](#) is a refinement of an idea originally introduced by [Pettis et al. \(1979\)](#). It builds on the fact that the volume of a d -dimensional ball grows proportionally to r^d —a principle underlying the construction of many other estimators as well (e.g., [Fan et al. \(2009\)](#); [Kleindessner and Luxburg \(2015\)](#)).

Under the flatness assumption, the probability that a random data point \mathbf{X} falls within $B_{\mathbf{x}_k}(R)$ is

$$\mathbb{P}(\mathbf{X} \in B_{\mathbf{x}_k}(R)) = f_{\mathbf{X}}(\mathbf{x}_k) V_d R^d,$$

where $V_d = \pi^{d/2}/\Gamma(d/2+1)$ is the volume of the unit ball in \mathbb{R}^d .

Given \mathbf{x}_k and the two neighbors $\mathbf{x}_k^{\ell_1}$ and $\mathbf{x}_k^{\ell_2}$, with $\ell_1 = \lceil K/2 \rceil$ and $\ell_2 = K$, let

$$r_1 = \|\mathbf{x}_k - \mathbf{x}_k^{\ell_1}\|, \quad \text{and} \quad r_2 = \|\mathbf{x}_k - \mathbf{x}_k^{\ell_2}\|.$$

That is, r_1 and r_2 are the distances from \mathbf{x}_k to its ℓ_1 -th and ℓ_2 -th nearest neighbors, respectively.

Taking logarithms, we obtain

$$\log \mathbb{P}(\mathbf{X} \in B_{\mathbf{x}_k}(R)) = \log(f_{\mathbf{X}}(\mathbf{x}_k) V_d) + d \log R,$$

which holds for both r_1 and r_2 . Approximating the probabilities by empirical proportions gives

$$\log \frac{K}{2n} \doteq \log(f_{\mathbf{X}}(\mathbf{x}_k)V_d) + d \log r_1, \quad \log \frac{K}{n} \doteq \log(f_{\mathbf{X}}(\mathbf{x}_k)V_d) + d \log r_2.$$

Eliminating the nuisance term $\log(f_{\mathbf{X}}(\mathbf{x}_k)V_d)$ yields the local estimator

$$\hat{d}_k = \frac{\log 2}{\log r_2 - \log r_1}.$$

A global estimate \hat{d} is then obtained by averaging the \hat{d}_k 's across all neighborhoods or, alternatively, by a majority vote.

3.1.3 The MLE estimator

The maximum likelihood estimator (**MLE**) proposed by [Levina and Bickel \(2004\)](#) is one of the most influential methods for intrinsic dimension estimation. Although maximum likelihood arguments also appear in the derivation of many other estimators, the approach of [Levina and Bickel \(2004\)](#) is widely recognized under the name **MLE**.

Under the flatness assumption, let $N_{\mathbf{x}_k}(r)$ denote the number of data points falling within the d -dimensional ball $B_{\mathbf{x}_k}(r) \subseteq B_{\mathbf{x}_k}(R)$ for $0 < r \leq R$. As r increases from 0 to R , the distribution of $N_{\mathbf{x}_k}(r)$ can be approximated by an inhomogeneous Poisson process with rate function

$$\lambda_k(r) = n f(\mathbf{x}_k) d r^{d-1} V_d,$$

where $V_d = \pi^{d/2}/\Gamma(d/2 + 1)$ is the volume of the d -dimensional unit ball. This form arises because the expected number of data points in $B_{\mathbf{x}_k}(r)$ is $n f(\mathbf{x}_k)$ times the volume of the ball.

With this approximation, the log-likelihood for the observed neighborhood distances is

$$\ell(d, \theta) = \sum_{\ell=1}^K \log \lambda_k(\|\mathbf{x}_k^\ell - \mathbf{x}_k\|) - \int_0^R \lambda_k(t) dt,$$

where $\theta = \log(nf(\mathbf{x}_k))$ is treated as a nuisance parameter. Maximization with respect to d yields the local estimator

$$\hat{d}_k = \left[\frac{1}{K} \sum_{\ell=1}^K \log \frac{R}{\|\mathbf{x}_k^\ell - \mathbf{x}_k\|} \right]^{-1},$$

where R is typically taken as the distance from \mathbf{x}_k to its K -th nearest neighbor \mathbf{x}_k^K . A global estimate \hat{d} is then obtained by averaging over all local estimates.

Several refinements of **MLE** have since been proposed. For example, [Gupta and Huang \(2012\)](#) introduce a regularized likelihood formulation, while [Gomtsyan et al. \(2019\)](#) modify the rate function to reduce reliance on the flatness assumption.

3.1.4 The DanCo estimator

A major limitation of **MLE** and many other flatness-based estimators, as will be shown in our experiments, is that their performance deteriorates in the form of underestimation as the intrinsic dimension d increases. The reason is that as the intrinsic dimension d grows, the local geometry of the manifold becomes increasingly complex. As we discussed earlier, fully characterizing the local structure requires up to $(p - d)d$ principal curvatures. Since p has to grow with d as well, the number of data points needed within a neighborhood for reliable estimation increases rapidly. When the sample size n is fixed, we might not be able to choose a large enough neighborhood size K since a very big neighborhood will jeopardize the flatness assumption.

The dimensionality estimator based on angle and norm concentration, or **DanCo** ([Ceruti et al., 2014](#)), is primarily designed to improve estimation performance on high-dimensional manifolds from a small sample. It builds upon two of its predecessors **IDEA** and **MIND_{KL}** ([Rozza et al., 2011, 2012](#)). **DanCo** employs two key ideas. First, it leverages both the norm distribution and the angle distribution of data points within a neighborhood. Second, it estimates the Kullback–Leibler (KL) divergence between the estimated likelihood and a

simulated likelihood obtained from data under idealized conditions. The later approach mitigates the issue of data sparsity that affects maximum likelihood estimation in high dimensions.

Under the flatness assumption, for a given \mathbf{x}_k and its neighbors, it can be shown that the ratio of the minimum to maximum distances (norms), defined as

$$r_k = \frac{\|\mathbf{x}_k^1 - \mathbf{x}_k\|}{\|\mathbf{x}_k^K - \mathbf{x}_k\|},$$

is distributed according to the density function

$$f_{\text{norm}}(r; d) = K d r^{d-1} (1 - r^{d-1})^{K-1}.$$

Hence, given the values r_1, \dots, r_n computed from the sample, we can obtain the maximum likelihood estimate of d , denoted as \hat{d}_{norm} .

On the other hand, motivated by the observation that random vectors drawn uniformly from the unit sphere in high-dimensional spaces tend to be nearly orthogonal, the authors conjecture that for sufficiently large d , the mutual angles among the vectors $\mathbf{x}_k^\ell - \mathbf{x}_k$ will concentrate around some theoretical mean ν (not to be confused with the Hausdorff measure) and can be modeled using the von Mises distribution with density

$$f_{\text{angle}}(\theta; \tau, \nu, d) = \frac{e^{\tau \cos(\theta - \nu)}}{2\pi I_0(d)},$$

where I_0 is the modified Bessel function of the first kind (Mardia and Jupp, 2009). Although d appears as a parameter in the above density, obtaining an analytical form for its maximum likelihood estimate is intractable. Instead, **DanCo** only estimates the parameters τ and ν locally using the $K(K-1)/2$ pairwise angles within each neighborhood, then averages the local estimates to obtain global estimates $\hat{\tau}$ and $\hat{\nu}$.

Next, for $q = 1, \dots, p$, one simulates a dataset of size n consisting of points uniformly distributed in the q -dimensional unit ball, and compute estimates \hat{d}_q , $\hat{\tau}_q$, and $\hat{\nu}_q$ using the

same procedures for norm and angle distributions. Then, because of the independence between norms and angles, the Kullback–Leibler divergence between the empirical and simulated densities can be estimated as

$$\text{KL}(f_{\text{norm}}(r; \hat{d}_{\text{norm}}), f_{\text{norm}}(r; \hat{d}_q)) + \text{KL}(f_{\text{angle}}(\theta; \hat{\tau}, \hat{\nu}, q), f_{\text{angle}}(\theta; \hat{\tau}_q, \hat{\nu}_q, q)),$$

in which $\text{KL}(f_1, f_2) = \int_{\mathbb{R}} f_1(x) \log[f_1(x)/f_2(x)]dx$. Because of the flatness assumption, we expect the KL divergence to attain its minimum when $q = d$. This value of q is then taken as the **DanCo** estimate of the intrinsic dimension.

A major drawback of **DanCo** is that it is arguably the most computationally intensive manifold dimension estimator, especially when the sample size n and the ambient dimension p are large. This is because the method requires simulating data for all dimensions from 1 up to p , and comparing their corresponding estimates with those derived from the observed sample. The problem can become serious if the manifold is nonlinearly embedded into a space with much larger dimensions, as is the case for some real-world datasets (see Section 5). In this case, besides the computational cost, to obtain an accurate estimate, the method needs to produce a larger KL divergence estimate for every $q \neq d$, which is quite unlikely.

To mitigate this, one can limit the maximum intrinsic dimension considered or, as in the authors’ **Matlab** implementation (Lombardi, 2020), use spline interpolation on a coarser set of simulated datasets spanning different sample sizes and intrinsic dimensions (with some fixed neighborhood size). This avoids simulating every dimension by smoothly estimating the KL-divergence curve, stabilizing the estimate by reducing noise. Additionally, if the initial maximum likelihood estimate is very small, the method skips the KL-divergence comparison and relies on the initial estimate to prevent unstable corrections and improve efficiency.

3.1.5 The TLE estimator

Like **DanCo**, the intrinsic dimension estimator with tight localities (**TLE**) from [Amsaleg et al. \(2019\)](#) also aims to tackle the problem of insufficient data points when estimating in neighborhoods. It builds on a previously proposed manifold dimension estimator proposed by the same authors ([Amsaleg et al., 2018](#)).

TLE is constructed on the concept of the so-called local intrinsic dimension (LID), which is a population parameter defined for any given point on the manifold and is determined by the local geometric structure as well as the data distribution. With the help of statistical theory about extreme values, for a given \mathbf{x}_k and its neighbors, the LID at \mathbf{x}_k can be estimated as

$$\hat{d}_k = - \left(\frac{1}{K} \sum_{\ell=1}^K \ln \frac{\|\mathbf{x}_k^\ell - \mathbf{x}_k\|}{\|\mathbf{x}_k^K - \mathbf{x}_k\|} \right)^{-1}.$$

It is easy to show that under the flatness assumption, the LID at \mathbf{x}_k is the same as the intrinsic dimension d , hence, the above expression can be used as a manifold dimension estimator as well, and a global estimate can simply be taken as the average of the \hat{d}_k s.

The extreme value theory does not change the problem that \hat{d}_k will suffer if the neighborhood size is too small. To tackle the issue, **TLE** further makes use of the fact that different points in the neighborhood roughly have the same LID, and improves the estimation by obtaining more estimates within each neighborhood. More specifically, \hat{d}_k is modified as

$$\hat{d}_k = - \left(\frac{1}{2K(K-1)} \sum_{\mathbf{v} \neq \mathbf{w}} \left[\log \frac{d(\mathbf{v}, \mathbf{w})}{R} + \log \frac{d(2\mathbf{x}_k - \mathbf{v}, \mathbf{w})}{R} \right] \right)^{-1},$$

where \mathbf{v}, \mathbf{w} are taken from the K nearest neighbors of \mathbf{x}_k , R is estimated as $\|\mathbf{x}_k - \mathbf{x}_k^K\|$,

and $d(\mathbf{v}, \mathbf{w})$ is defined as⁸

$$d(\mathbf{v}, \mathbf{w}) = \sqrt{[\mathbf{u} \cdot (\mathbf{x}_k - \mathbf{v})]^2 + R\mathbf{u} \cdot (\mathbf{w} - \mathbf{v}) - \mathbf{u} \cdot (\mathbf{x}_k - \mathbf{v})}, \text{ where } \mathbf{u} = \frac{R(\mathbf{w} - \mathbf{v})}{R^2 - \|\mathbf{x}_k - \mathbf{v}\|^2},$$

which is the radius of the circle centred on the line segment between \mathbf{x}_k and \mathbf{v} and passing through \mathbf{v} and \mathbf{w} .

3.1.6 The TwoNN estimator

We conclude flatness-based estimators with the two nearest neighbor estimators (**TwoNN**) proposed by [Facco et al. \(2017\)](#), which pushes the flatness assumption to its extreme by only making use of the two nearest neighbors of every \mathbf{x}_k , hence avoiding curvatures as much as possible. Let $r_k^1 = \|\mathbf{x}_k - \mathbf{x}_k^1\|$ and $r_k^2 = \|\mathbf{x}_k - \mathbf{x}_k^2\|$ be the distances of \mathbf{x}_k to its two nearest neighbors. The key insight of **TwoNN** is that under the flatness assumption, the distribution function of the ratio of the distances, $\mu_k = r_k^2/r_k^1$, is

$$F(\mu) = (1 - \mu^{-d})1_{[0, \infty]}(\mu),$$

which not only depends on d , but also is independent of the data distribution. As a result, we do not need any manipulation like in [Pettis et al. \(1979\)](#) and **MADA**, and can obtain a global estimate of d directly by employing the following relation

$$-\frac{\log[1 - F(\mu)]}{\log \mu} = d.$$

All ratios μ_k s in the sample are obtained first to construct an empirical distribution function for F , and \hat{d} is estimated through the least squares method with $\log[1 - \hat{F}(\mu_k)]$ s as responses and $-\log \mu_k$ s as predictors. Our experiments show that such a simple and elegant design works very well in many cases, making it a top performer among all flatness-based estimators.

⁸The formula assumes $\|\mathbf{v} - \mathbf{x}_k\| < R$, when $\|\mathbf{x}_k - \mathbf{v}\| = R$, the formula should be $d(\mathbf{w}, \mathbf{v}) = R\|\mathbf{w} - \mathbf{v}\|^2 / [2(\mathbf{x}_k - \mathbf{v}) \cdot (\mathbf{w} - \mathbf{v})]$.

3.2 Other Estimators

Flatness-based estimators behave poorly when the flatness assumption is significantly violated, which typically occurs on manifolds with complex curvature, such as those that are nonlinearly embedded in Euclidean space.

3.2.1 The CA-PCA estimator

Our experimental results indicate that among the flatness-based estimators, **TwoNN** is the most robust to curvature—likely because it minimizes curvature effects by relying solely on the two nearest neighbors of each data point. At the other end of the design spectrum are methods that explicitly account for curvature. A recent example of this approach is the curvature-adjusted PCA estimator, **CA-PCA** (Gilbert and O’Neill, 2025), which incorporates curvature information directly into the estimation procedure.

Recall the threshold-free local PCA estimator,

$$\hat{d}_k = \operatorname{argmin}_{1 \leq q \leq p} \sum_{j=1}^p \left(\frac{1}{R^2} \hat{\lambda}_j - \frac{1}{j+2} 1_{j \leq q} \right)^2,$$

which is based on the fact that under the flatness assumption (within the unit ball), the true covariance matrix Σ_k has its first d eigenvalues equal to $1/(d+2)$ and the $p-d$ remaining ones equal to 0 (Lim et al., 2024).

If the flatness assumption holds reasonably well, the sample covariance matrix $\hat{\Sigma}_k$ computed from the neighborhood of \mathbf{x}_k should closely resemble this theoretical structure, which will not be the case if the assumption is seriously violated.

The main idea of **CA-PCA** is to adjust the eigenvalues to account for curvature. Recall the local graph representation for manifolds: in a neighborhood of \mathbf{x}_k , points can be expressed in new coordinates as

$$(\mathbf{x}'_{1:d}, g(\mathbf{x}'_{1:d})),$$

where $\mathbf{x}'_{1:d}$ denotes coordinates with respect to directions within the tangent space $T_{\mathbf{x}_k} M$, and g is a sufficiently differentiable function. The flatness assumption corresponds to taking $g \doteq 0$, which is equivalent to its first-order Taylor approximation and thus ignoring curvature.

To incorporate curvature, we can use a second-order Taylor polynomial to approximate each component of g , which leads to

$$g_\ell(\mathbf{x}'_{1:d}) \doteq q_\ell(\mathbf{x}'_{1:d}) = \mathbf{x}'_{1:d}^\top Q_\ell \mathbf{x}'_{1:d}, \quad \ell = 1, \dots, p-d,$$

where Q_ℓ is a symmetric $d \times d$ matrix capturing principal curvatures in the ℓ -th normal direction.

In other words, instead of assuming the \mathbf{x}_k^ℓ s are uniformly distributed in some d -dimensional ball within the tangent space, we assume they are uniformly distributed in the *tangent quadratic space* defined by the q_ℓ s, as illustrated in Figure 6, which should be a better approximation.

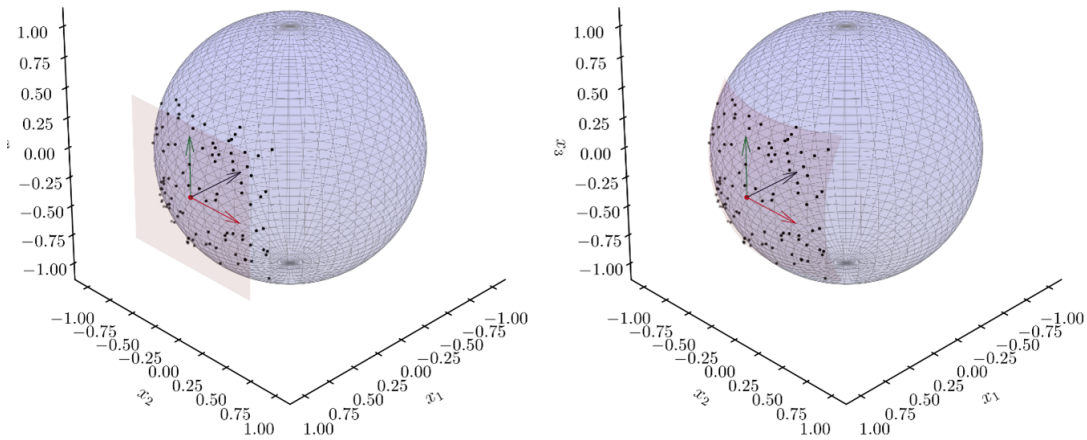


Figure 6: Approximate a neighborhood with the tangent space (left) and the tangent quadratic space (right).

Gilbert and O’Neill (2025) show that, with the quadratic approximation, the theoretical eigenvalues satisfy

$$\frac{1}{d+2} \doteq \lambda_j + \frac{3d+4}{d(d+4)} \sum_{i=d+1}^p \lambda_i.$$

Using this property, **CA-PCA** replaces the true eigenvalues with their sample version, and estimates the intrinsic dimension as

$$\hat{d}_k = \operatorname{argmin}_{1 \leq q \leq p} \left\{ \sqrt{\sum_{j=1}^q \left(\frac{1}{q+2} - \frac{1}{R^2} \left(\hat{\lambda}_j + \frac{3q+4}{q(q+4)} \sum_{i=q+1}^p \hat{\lambda}_{ki} \right) \right)^2} + \frac{2}{R^2} \sum_{j=q+1}^p \hat{\lambda}_j \right\},$$

where R can be estimated using $\|\mathbf{x}_k^K - \mathbf{x}_k\|^9$, and the last term serves as a penalty for neighborhoods where the quadratic approximation is of poor quality.

A global **CA-PCA** estimate is then obtained by averaging all the local estimates \hat{d}_k across the dataset.

3.2.2 The Wasserstein estimator

All estimators we have introduced so far make explicit use of the underlying manifold’s geometry to some extent. In contrast, the manifold dimension estimator based on the Wasserstein distance, **Wasserstein** (Block et al., 2022), incorporates the geometry only implicitly in its design.

Given two probability distributions μ and ν on \mathbb{R}^p , their Wasserstein distance is defined as

$$W_1(\mu, \nu) = \inf_{(\mathbf{X}, \mathbf{Y}) \sim \Gamma(\mu, \nu)} \mathbb{E} \left[\sum_{j=1}^p |X_j - Y_j| \right],$$

where \mathbf{X} and \mathbf{Y} are p -dimensional random vectors with marginal distributions μ and ν , respectively, and $\Gamma(\mu, \nu)$ denotes the set of all joint distributions with these marginals.

Intuitively, the Wasserstein distance measures the minimum expected transportation cost required to morph one probability distribution into another. When the marginals are

⁹The authors of **CA-PCA** actually suggest to use $(\|\mathbf{x}_k^K - \mathbf{x}_k\| + \|\mathbf{x}_k^{K-1} - \mathbf{x}_k\|)/2$, which is what we also use in our own implementation.

discrete—such as empirical distributions constructed from samples—the Wasserstein distance can be computed exactly by solving a discrete optimization problem (Villani, 2021).

Wasserstein is designed based on the theoretical result (Dudley, 1969) that the rate at which the difference between a distribution and its empirical estimate (measured in Wasserstein distance) converges as sample size increases is governed by the intrinsic dimension d of the manifold the distribution is supported on, rather than the ambient dimension p . Specifically, let P and P_n be the true and empirical distributions, it holds that

$$W_1(P_n, P) \propto n^{-\frac{1}{d}},$$

where the distance used in computing the Wasserstein metric is ideally the geodesic distance along the manifold.

The above result leads to the following expression of d which can be useful for the design of an estimator: given two samples of sizes m and αm ($\alpha > 1$ such that $\alpha m \in \mathbb{N}$) distributed according to P , we have

$$d = \frac{\log \alpha}{\log W_1(P_m, P) - \log W_1(P_{\alpha m}, P)}.$$

In practice, P is unknown. To overcome this, the authors propose splitting the full dataset into four parts of sizes $m, m, \alpha m$ and αm , which gives 4 empirical distribution estimates $P_m, P'_m, P_{\alpha m}$ and $P'_{\alpha m}$. Using P'_m and $P'_{\alpha m}$ to approximate P , the dimension can then be estimated as

$$\hat{d} = \frac{\log \alpha}{\log W_1(P_m, P'_m) - \log W_1(P_{\alpha m}, P'_{\alpha m})}.$$

For simplicity, Euclidean distances can be used as a substitute for geodesic distances when computing Wasserstein distances. Since \hat{d} is a global estimate, we should also generate multiple instances of it through repeated data splits, and to average these estimates in order to reduce uncertainty. For $d \geq 3$, under certain regularity conditions, it can be

shown that \hat{d} is close to d with high probability. However, verifying or establishing the conditions is nontrivial.

4 Empirical Evaluation

In this section, we present a comprehensive set of numerical experiments designed to evaluate the performance of the eight estimators introduced in Section 3. Among them, **DanCo** was previously identified as the top performer by [Campadelli et al. \(2015\)](#), while **MLE**, **MADA**, and **TwoNN** are well-established and widely used. By contrast, **TLE**, **Wasserstein**, and **CA-PCA** are more recent proposals that, to our knowledge, have not yet been systematically compared with earlier methods. We also include **Local PCA**, despite its frequent omission from benchmarks, using the threshold rule by [Fukunaga and Olsen \(1971\)](#).

Although R implementations of manifold dimension estimators are available, they remain relatively underdeveloped compared to their Python counterparts. See [You and Shung \(2022\)](#) for recent advances in the `Rdimtools` package ([You et al., 2022](#)), which consolidates and extends earlier efforts and will be further introduced in the next section. Moreover, Python implementations are generally faster, which is why we chose to perform all experiments in this section using Python. Specifically, we use the estimator implementations provided by the `scikit-dimension` (`skdim`) package ([Bac et al., 2021](#)) for **Local PCA**, **MLE**, **DanCo**, **MADA**, **TLE**, and **TwoNN**, while the **Wasserstein** and **CA-PCA** methods are implemented by ourselves. We also implement our own sampling and evalua-

tion procedures¹⁰.

We set the number of replicates to 100, meaning that each experiment was repeated 100 times on independently generated random samples under identical conditions. While prior intrinsic dimension studies such as [Levina and Bickel \(2004\)](#), [Campadelli et al. \(2015\)](#), and [Facco et al. \(2017\)](#) do not disclose their exact replicate counts—*some even appear to conduct the tests only once*—simulation-based evaluations in both statistics and machine learning commonly use between 100 and 200 replicates. This practice is consistent with general recommendations in the Monte Carlo literature [Robert and Casella \(2004\)](#); [Gentle \(2004\)](#); [Morris et al. \(2019\)](#), where 100 is typically sufficient for reliable comparative evaluations.

4.1 Impact of Individual Factors on Performance

The performance of a manifold dimension estimator can be affected by numerous factors. Among the most critical are the choice of hyperparameters (e.g., neighborhood size K), the sample size n , the ambient dimension p , the curvature of the manifold, and the level of noise.

Before undertaking large-scale comparisons across datasets and estimators, it is useful to benchmark each method under controlled variations of these factors. This approach clarifies the strengths, limitations, and sensitivities of the estimators.

In this subsection, we report results from a series of controlled experiments that examine how the estimators respond to changes in these key factors. Most of our controlled experiments are conducted on a 5-dimensional sphere, given it is a widely studied manifold that

¹⁰This refers to the full pipeline: selecting hyperparameters, repeatedly sampling from the target manifolds, applying the estimators, and recording the results. In the GitHub repository, any file starting with "run" corresponds to this evaluation process. The implementations of the estimators and sampling routines are located in `mymodules.py`.

serves as a standard benchmark for dimension estimation. Its geometry is well understood, and its curvature can be easily adjusted by varying the radius. Choosing dimension $d = 5$ offers a suitable level of difficulty for the estimators—challenging enough to be informative, yet not so high as to obscure meaningful differences in performance.

4.1.1 Effect of hyperparameters

The choice of hyperparameter values is of crucial importance to the performance of manifold dimension estimators (Campadelli et al., 2015). Among the eight methods considered, **TwoNN** is parameter-free. For **Local PCA**, **MLE**, **DanCo**, **MADA**, **TLE** and **CAPCA**, the main hyperparameter is the neighborhood size K , i.e., a neighborhood consisting of the K nearest neighbors (out of the total n points) of the selected point. By contrast, **Wasserstein** relies on a different type of hyperparameter: the subsample splitting ratio $\alpha \in (1, \infty)$. This parameter determines how a sample of size n is divided into four parts of sizes m , m , αm , and αm where $m = n/(2 + 2\alpha)$.

When the sample size n is fixed, choosing K involves a trade-off. On the one hand, smaller neighborhoods are desirable because many estimators rely on local approximations of the manifold, which are more accurate when the neighborhood is nearly flat (or quadratic). On the other hand, larger neighborhoods provide more points per estimate, which can improve stability and accuracy—particularly for estimators based on maximum likelihood, whose performance typically benefits from larger sample sizes.

The experiment designed to benchmark the estimators’ sensitivity to hyperparameter choices is summarized in Table 1 below; see Tables 15–17 for full results.

Table 1: Experiment design for evaluating the effect of K or α .

Component	Description
Manifold M	A 5-dimensional sphere linearly embedded in \mathbb{R}^{10} .
Sample	$n = 1,000$ data points uniformly distributed over M .
Hyperparameter	Neighborhood size K or regularization parameter α .
Factor	<ul style="list-style-type: none"> • $K \in \{5, 10, 20, 30, 40, 50, 100, 200, 500, 750\}$ (general) • $K \in \{2, 4, 6, \dots, 20\}$ for DanCo • $\alpha \in \{1.01, 1.2, 1.4, 1.6, 1.8, 2, 4, 6, 8, 10\}$

In Figure 7, we show how the estimates of **Local PCA**, **MLE**, **MADA**, **TLE** and **CA-PCA** vary with K . We observe that the estimates of **Local PCA** increase with K , while those of **MLE**, **MADA**, and **TLE** decrease. In contrast, **CA-PCA**'s estimates remain stable across the tested range of K , suggesting that accounting for curvature may improve the estimator's robustness to hyperparameter selection. One possible explanation for this phenomenon is that the quadratic space approximation **CA-PCA** relies on works well for a much larger range of neighborhood sizes. All five estimators tend to stabilize when K is between 20 and 50.

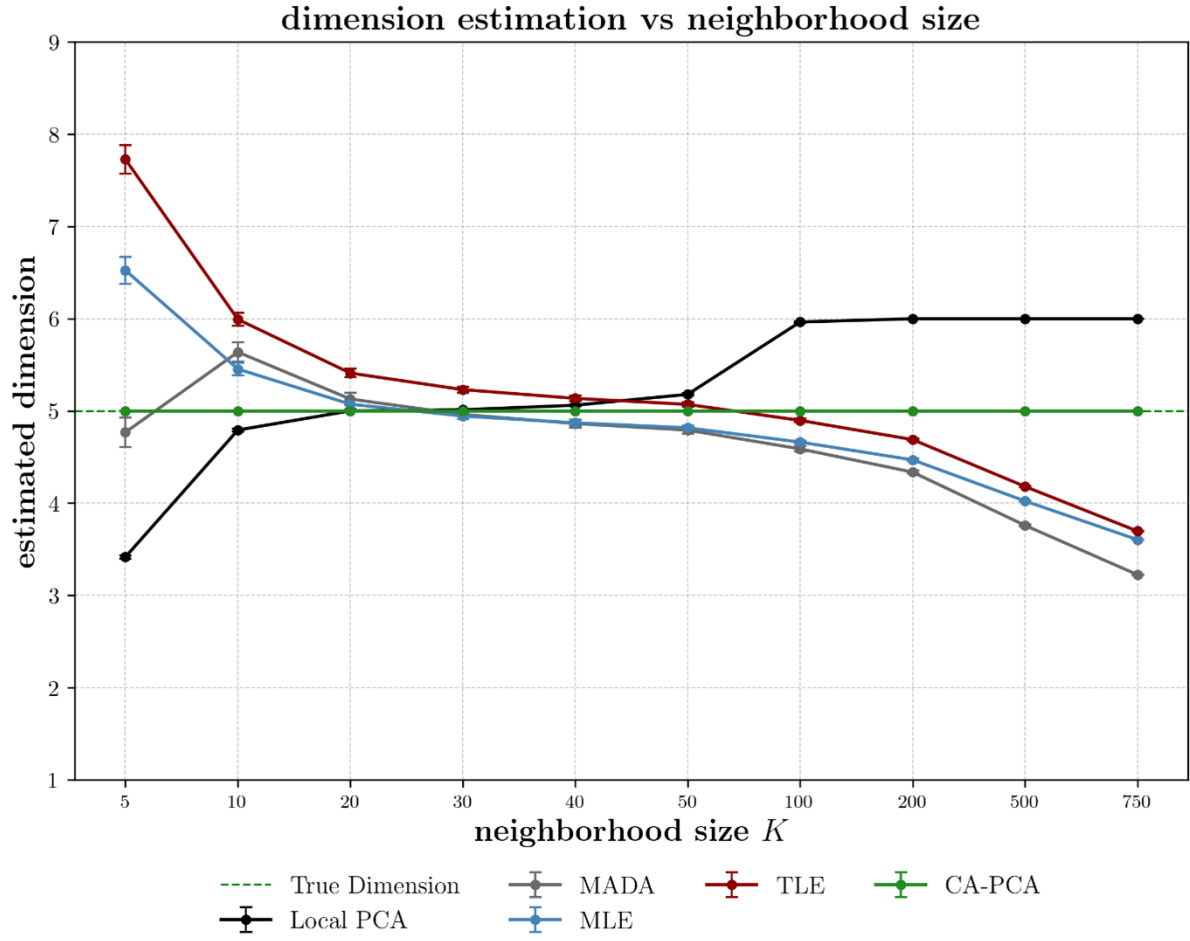
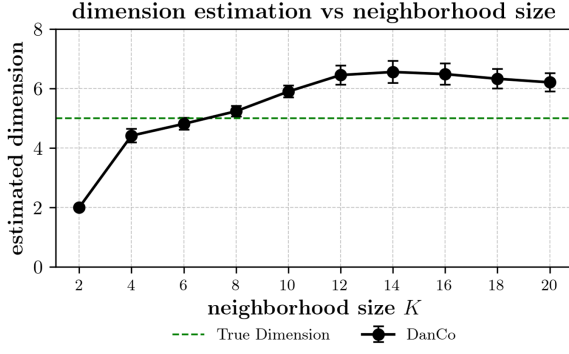
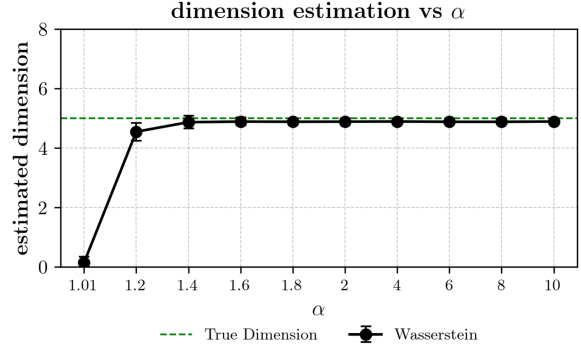


Figure 7: Dimension estimates versus K .

In comparison, **DanCo** requires a much smaller neighborhood size, and thus its results are plotted separately in Figure 8a. The estimates of **Wasserstein** against α are plotted in Figure 8b.



(a) **DanCo**



(b) **Wasserstein**

Figure 8: Dimension estimates versus K or α .

In this particular experiments, **DanCo** gives the most accurate estimate when $K \doteq 7$. Our other experiments show that in general, **DanCo** performs the best with neighborhood sizes around 10, which aligns with the setting used in [Ceruti et al. \(2014\)](#), although the authors did not provide an explanation for this choice.

The **Wasserstein** estimates stabilize for $\alpha > 1.4$.

4.1.2 Effect of sample size

For flatness-based estimators, when all other conditions are fixed, the performance of an estimator is expected to improve as the sample size n increases. The improvement can be understood from two complementary perspectives. First, a fixed neighborhood size K combined with a larger n leads to a higher concentration of data points within each neighborhood on the manifold. This local densification improves the quality of linear approximations by the tangent space, and thus enhances estimation accuracy, as illustrated in Figure 9.

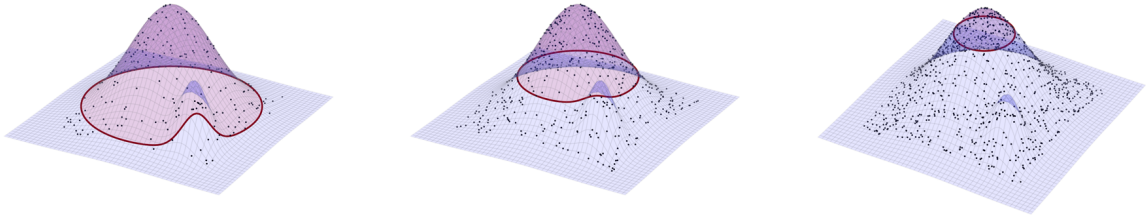


Figure 9: As the sample size n increases, a fixed-size neighborhood K more accurately captures the local structure of the underlying manifold.

At the same time, a larger sample size also allows us to choose larger neighborhood sizes, which can further improve the quality of local estimates by reducing variance without sacrificing locality.

Our experiment design to benchmark the estimators' performance against sample sizes is shown in Table 2, which focuses on the asymptotic mode with K fixed and $n \rightarrow \infty$, as it is more relevant for our choice of estimators; see Table 18 for full results. Another asymptotic mode which has been considered in the literature is $K \rightarrow \infty, n \rightarrow \infty$ with their ratio fixed (Costa and Hero, 2004b).

Table 2: Experiment design for evaluating the effect of n .

Component	Description
Manifold M	a 5-dimensional sphere linearly embedded in \mathbb{R}^{10}
Sample	$n = 1,000$ data points uniformly distributed on M
Hyperparameter	$K = 100$ (10 for DanCo), $\alpha = 5$
Factor	sample size $n \in \{200, 400, \dots, 1000, 2000, \dots, 5000, 10000\}$

We can see from Figure 10 that for all estimators, both the bias and the variance of the estimates converge towards 0 as n grows.

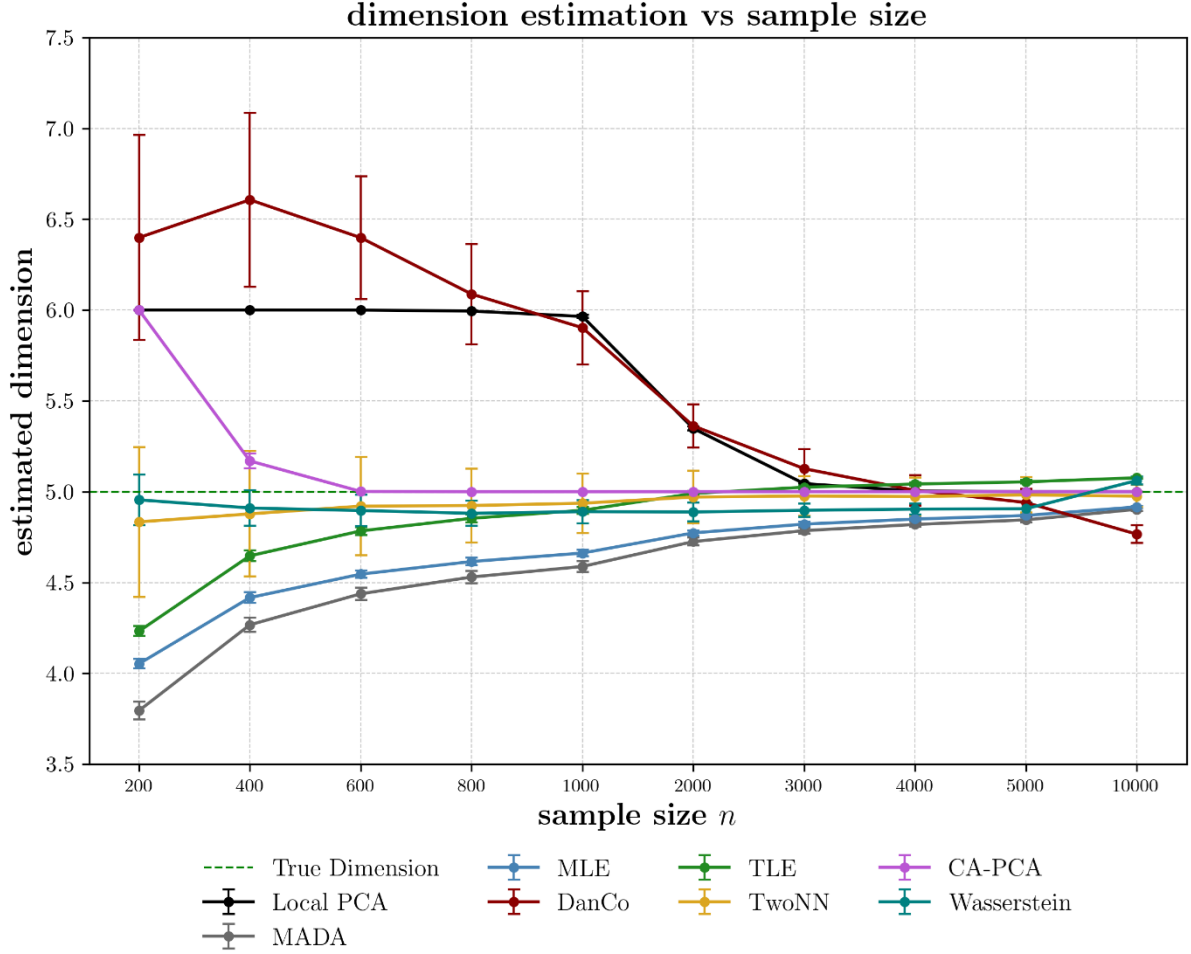


Figure 10: Dimension estimate versus sample size.

In particular, **CA-PCA** converges much faster than **Local PCA**, consistent with our previous result that it is more robust with respect to K given fixed n . Besides, **TwoNN** and **Wasserstein** have very low bias even with small sample sizes, making them strong candidates when we do not have sufficient data.

We also notice that **DanCo** seems to be asymptotically biased downward.

4.1.3 Effect of ambient dimension

A high ambient dimension p can affect the performance of manifold dimension estimators in two distinct scenarios. As we discussed before, a manifold can be embedded *linearly*

into Euclidean spaces with much higher dimensions. In this case, the intrinsic structure of the manifold remains unchanged and the manifold still lies entirely within some lower-dimensional subspace; on the other hand, a nonlinear embedding into a higher-dimensional space results in a different manifold, typically with increased curvature and geometric complexity. In view of the dimension estimators, they will recover the intrinsic dimension of the image manifold, regardless of how it was embedded. Therefore, nonlinear embeddings—which alter the geometry of the manifold—are better studied when we consider curvature as the factor, and we focus on the effect of linear embeddings here.

Table 3: Experiment design for testing the effect of p .

Component	Description
Manifold M	a 5-dimensional sphere linearly embedded in \mathbb{R}^p
Sample	$n = 1000$ data points uniformly distributed on M
Hyperparameter	$K = 100$ ($K = 10$ for DanCo), $\alpha = 5$
Factor	ambient dimension $p \in \{6, 10, 15, \dots, 50\}$

Intuitively, a linear embedding into higher ambient dimensions should not significantly impact the performance of dimension estimators. This can be confirmed by their mathematical formulations. For example, the **MLE** estimate

$$\hat{d}_k = \left[\frac{1}{K} \sum_{\ell=1}^K \log \frac{R}{\|\mathbf{x}_k^\ell - \mathbf{x}_k\|} \right]^{-1}$$

depends only on norms of pairwise differences between data points in the neighborhood, which are invariant under linear embeddings. Our experiment confirms this analysis: as shown in Figure 11, the estimates from all estimators remain constant under linear embeddings into increasingly higher dimensional spaces¹¹. The experiment design is summarised in Table 3; see Table 19 for full results.

¹¹We observe fluctuations in the variance of **DanCo**, which we attribute to the randomness introduced

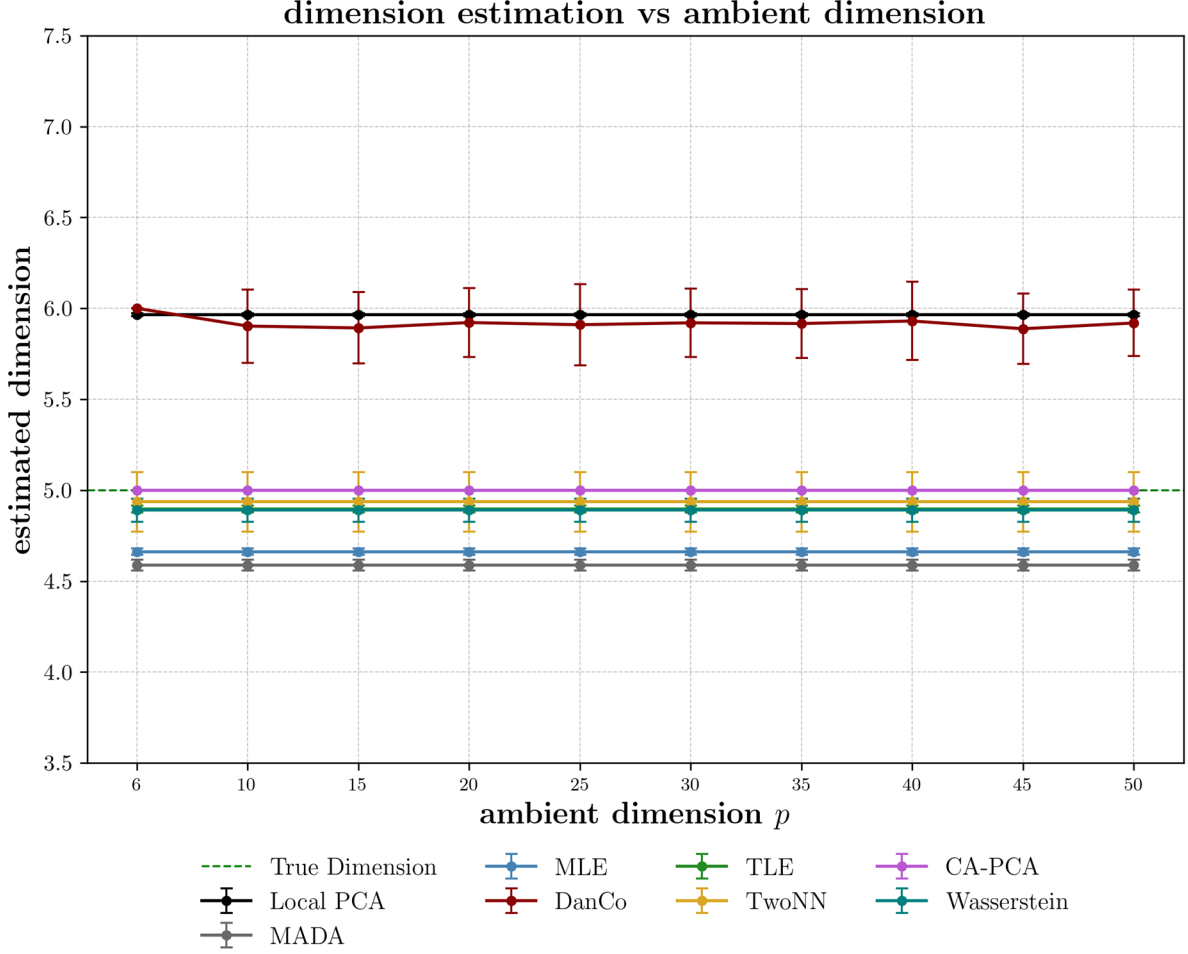


Figure 11: Dimension estimate versus ambient dimension.

4.1.4 Effect of curvature

Three aspects of curvature can affect the performance of the estimators: the magnitude of curvature, the number of principal curvatures, and the variation of curvature along the manifold. In this section, we evaluate their effects through three separate experiments.

Intuitively, as suggested in much of the prior literature, it may seem that larger curvature makes the manifold dimension estimation problem more difficult, since the manifold becomes less flat. However, our experiment shows that this relationship is not so straight-
by sampling from the ideal dataset; see the previous section for further details.

forward. The magnitude of the principal curvatures only significantly affects estimator performance when combined with other factors such as noise; by itself, it does not appear to increase the difficulty of the problem.

Table 4: Experiment design for evaluating the effect of curvature size.

Component	Description
Manifold M	a $d = 5$ -dimensional sphere linearly embedded in \mathbb{R}^{10}
Sample	$n = 1,000$ data points uniformly distributed on M
Hyperparameter	$K = 100$ ($K = 10$ for DanCo), $\alpha = 5$
Factor	radius R of the sphere, ranging from 0.01 to 100

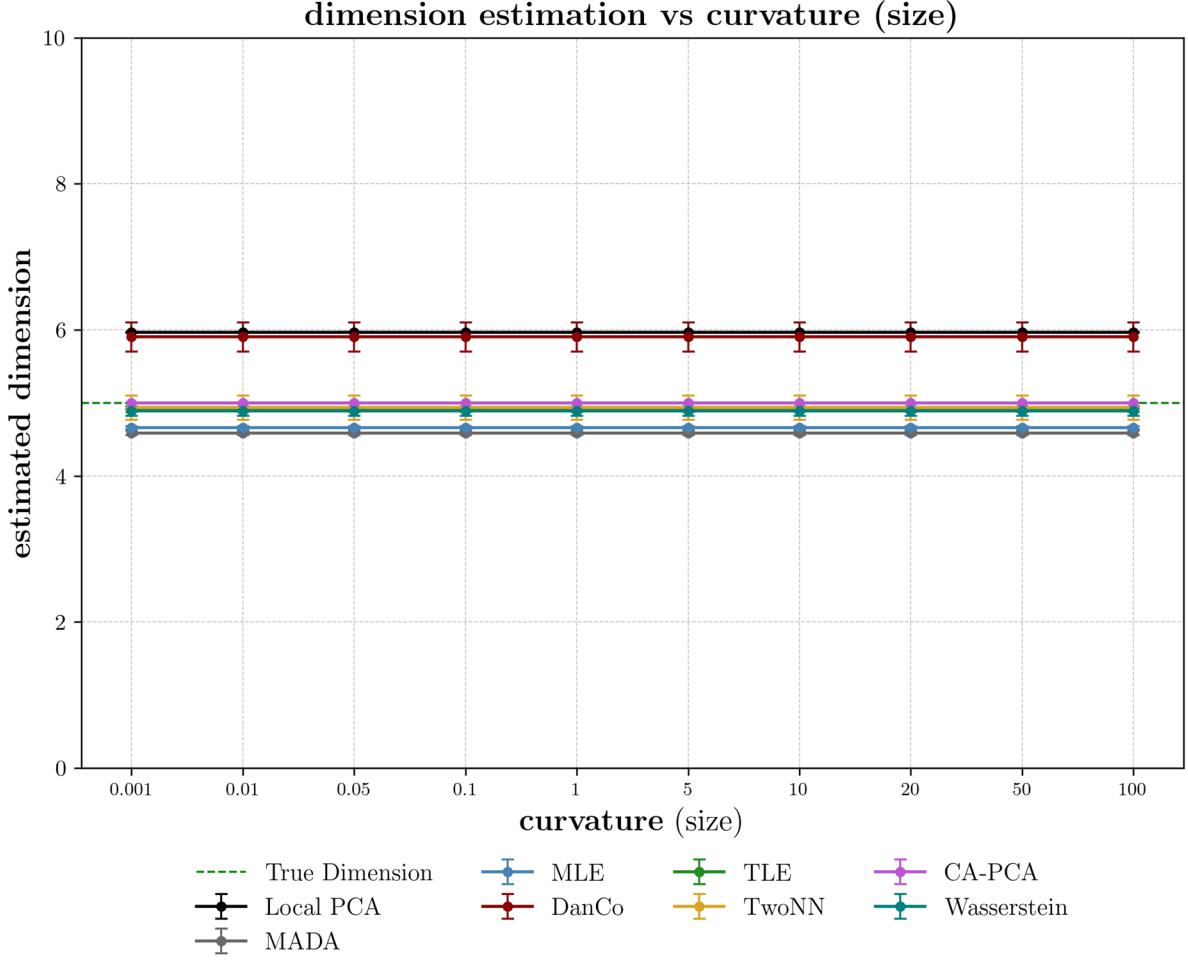


Figure 12: Dimension Estimate vs Curvature (size)

We start by varying the curvature size through the radius of the sphere, as indicated in Table 4, see Table 20 for full results. In Figure 12, we show the estimates of various estimators on a 5-dimensional sphere embedded in \mathbb{R}^{10} , with varying curvature values. The principal curvature of a sphere is everywhere $1/R$, so a smaller radius corresponds to a higher curvature. As shown in the figure, the performance of all eight estimators remains largely unaffected by changes in the curvature magnitude alone. The most likely explanation is the perfect symmetry of the sphere: regardless of its radius, for a fixed sample size n and neighborhood size K , each neighborhood covers the same proportion of the sphere, leading to identical estimates.

Our next experiment evaluates the effect of the number of principal curvatures. In this setting, we consider spheres of increasing intrinsic dimension; see Table 5 for details, and see Table 21 for full results..

Table 5: Experiment design for evaluating the effect of the number of principal curvatures.

Component	Description
Manifold M	a d -dimensional sphere linearly embedded in \mathbb{R}^{2d}
Sample	$n = 1,000$ data points uniformly distributed on M
Hyperparameter	$K = 100$ ($K = 10$ for DanCo), $\alpha = 5$
Factor	$d \in \{2, 4, \dots, 20\}$

Since a d -dimensional sphere resides in a $(d + 1)$ -dimensional subspace, the number of principal curvatures required to fully capture the local geometry at any point is $d \times 1 = d$.

As the curvature structure becomes more complex, more data points are generally required to accurately identify the underlying geometry. Therefore, we expect that for a fixed sample size n and neighborhood size K , the performance of manifold dimension estimators will deteriorate as d increases.

The results of this experiment are presented in Figure 13.

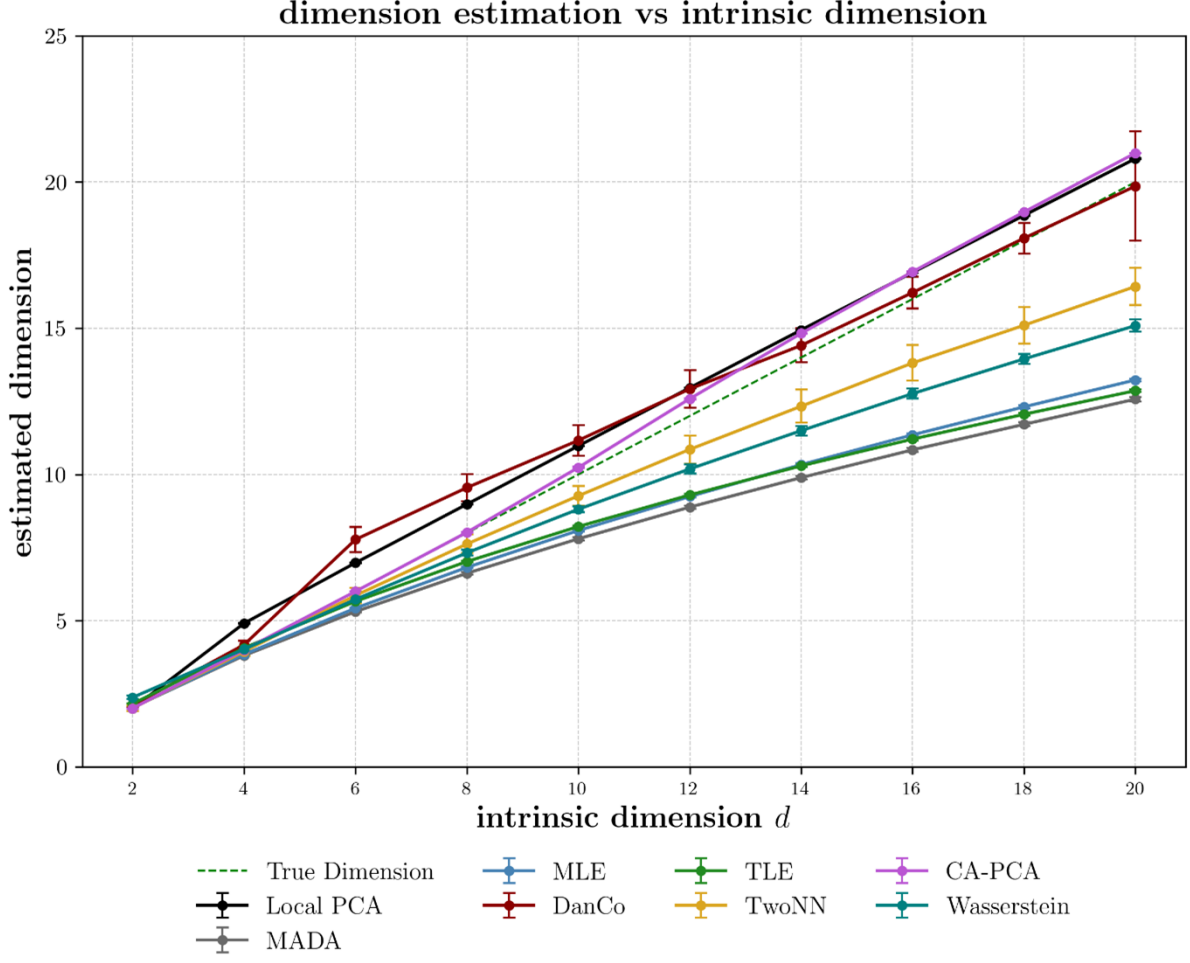


Figure 13: Dimension Estimate vs Curvature (d)

We observe that three estimators—**Local PCA**, **CA-PCA**, and **DanCo**—are quite robust to increasing intrinsic dimension. In contrast, the remaining five estimators exhibit increasing bias with respect to d , consistently underestimating the true dimension.

As previously discussed, **DanCo** is specifically designed to address the underestimation problem in high-dimensional settings, leveraging an estimated KL divergence. The two PCA-based methods likely benefit from their simplicity and minimal reliance on strong modeling assumptions: since PCA aims to estimate the tangent space directly under the flatness assumption, it can achieve reliable results with as few as d data points.

Our final experiment on curvature, described in Table 6, examines the effect of curvature variation along the manifold; see Table 22 for full results..

Table 6: Experiment design to test the effect of curvature variation.

Component	Description
Manifold M	a 3-dimensional deformed sphere in \mathbb{R}^6
Sample	$n = 1,000$ data points uniformly distributed on M
Hyperparameter	$K = 100$ ($K = 10$ for DanCo), $\alpha = 5$
Factor	$c \in \{0.0001, 0.001, 0.01, 0.05, 0.1, 0.2, 0.5, 1, 2, 4\}$

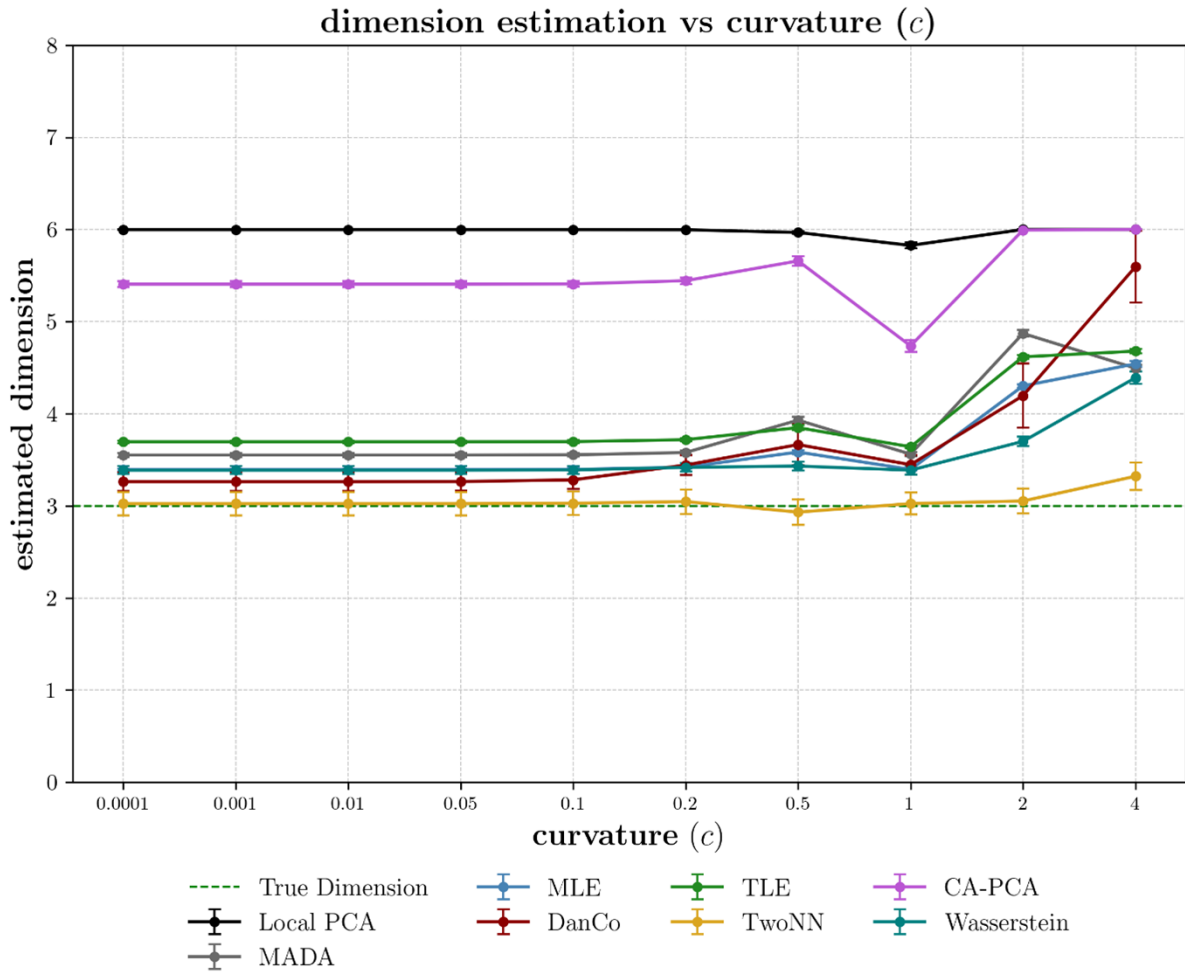


Figure 14: Dimension Estimate vs Curvature (c)

For this purpose, we use a deformed sphere as the underlying manifold, whose curvature complexity increases globally with the parameter c (see Figure 5). We expect the performance of the estimators to deteriorate as c increases, which is confirmed by the results shown in Figure 14.

Notably, the estimates from **Local PCA** and **CA-PCA** degrade more significantly than those of other estimators for the given neighborhood size, suggesting a potential weakness when dealing with nonlinear embeddings. However, **CA-PCA** still outperforms **Local PCA** in this setting.

We can conjecture that when $c = 1$, the geometry of the resulting manifold happens to align well with the assumptions or mechanisms built into the **CA-PCA** estimator, as well as the parameter setting, leading to the improved observed performance.

4.1.5 Effect of noise

It is quite unlikely that all data points in a sample of data lie exactly on some underlying manifold. A more plausible scenario is that they lie close to the manifold instead. A natural model for this setting assumes that data points originate from the manifold and are subsequently perturbed by some noise, which we assume for simplicity to be p -dimensional Gaussian $N(0, \sigma^2 I_p)$ in our experiment.

Almost all state-of-the-art manifold dimension estimators are designed under the assumption of a noiseless sample. Consequently, they are often not well-suited for application to noisy datasets, despite demonstrating strong performance on noise-free examples. The presence of noise significantly complicates the manifold dimension estimation problem, as even small perturbations can easily overwhelm the local structure upon which many estimators rely, and increasing sample size n alone will not necessarily be helpful.

The experiment benchmarking the estimators' performance against varying Gaussian noise

levels is described in Table 7, see Table 23 for full results.

Table 7: Experiment design to test the effect of noise.

Component	Description
Manifold M	a 5-dimensional unit sphere linearly embedded in \mathbb{R}^{10}
Sample	$n = 1,000$ data points uniformly distributed on M , with additive Gaussian noise from $N(0, \sigma^2 I_{10})$
Hyperparameter	$K = 100$ ($K = 10$ for DanCo), $\alpha = 5$
Factor	$\sigma \in \{0.00, 0.01, \dots, 0.09\}$

The results are shown in Figure 15.

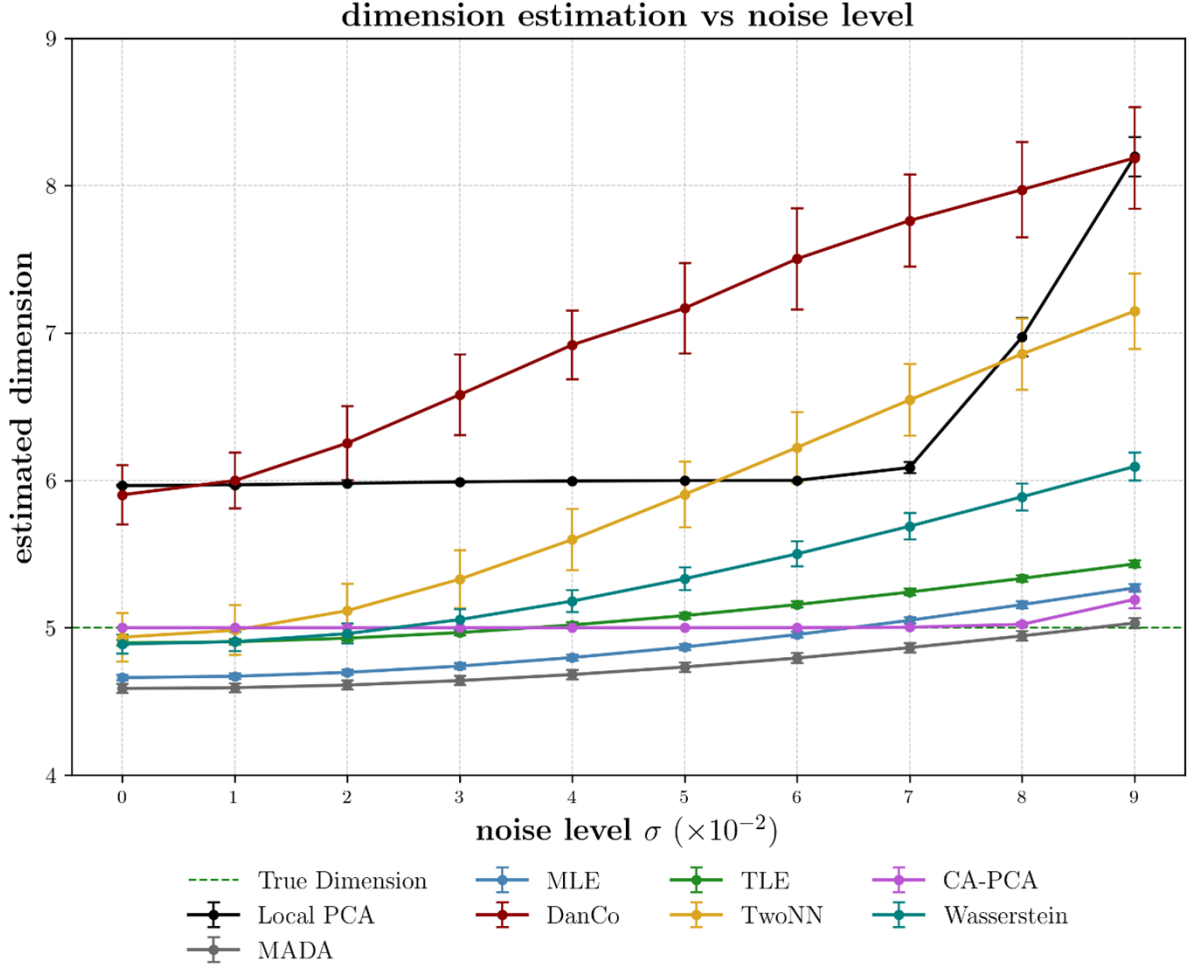


Figure 15: Dimension Estimate vs Noise Level (σ^2)

We observe that as the noise level increases, the estimates from all estimators tend to increase, and their performance deteriorates in terms of both bias and variance. This indicates that the estimators become less effective at capturing the true underlying structure under higher noise conditions.

In particular, **TwoNN** and **DanCo**, which have been top performers in noiseless scenarios, exhibit greater sensitivity to noise compared to the other methods. This highlights the challenging nature of manifold dimension estimation in noisy environments.

We have illustrated in Figure 12 that the size of curvature alone is not a driver of the

performance of manifold dimension estimators. It usually manifests itself when combined with other factors, and noise is one such example. The combination of curvature and noise leads to the so-called *noise curvature dilemma*. We have seen the choice of neighborhood size is important for most flatness-based estimators. For the flatness assumption to hold, the neighborhood K size is preferable to be small compared to the sample size n , especially if the manifold has large curvatures. However, when Gaussian noise is present, using a very small neighborhood may result in the noise to dominate any meaningful local structure, and we have to look at a large enough neighborhood instead to not be overwhelmed, hence the dilemma, which is illustrated in Figure 16a–16b.

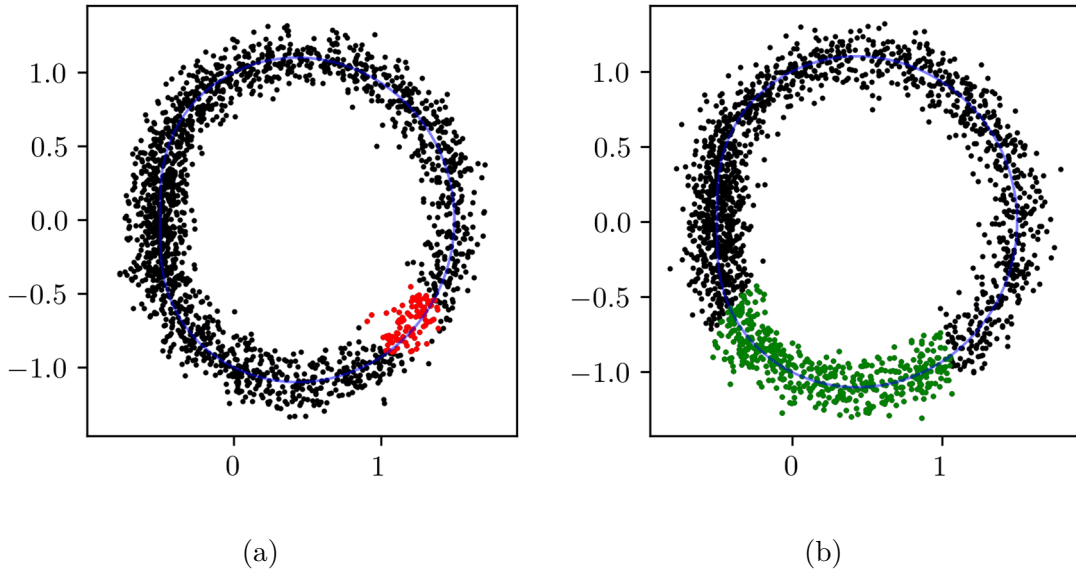


Figure 16: Demonstration of Noise-Curvature Dilemma: A small neighborhood can be overwhelmed by noise and fail to capture the true underlying geometry, while a large neighborhood is more robust to noise but will jeopardize the flatness assumption more.

In Figure 17, we superimpose on Figure 15 the estimates curved obtained using a sphere with a smaller radius—and thus larger curvature—while keeping all other conditions identical. The new estimates are shown as dashed curves.

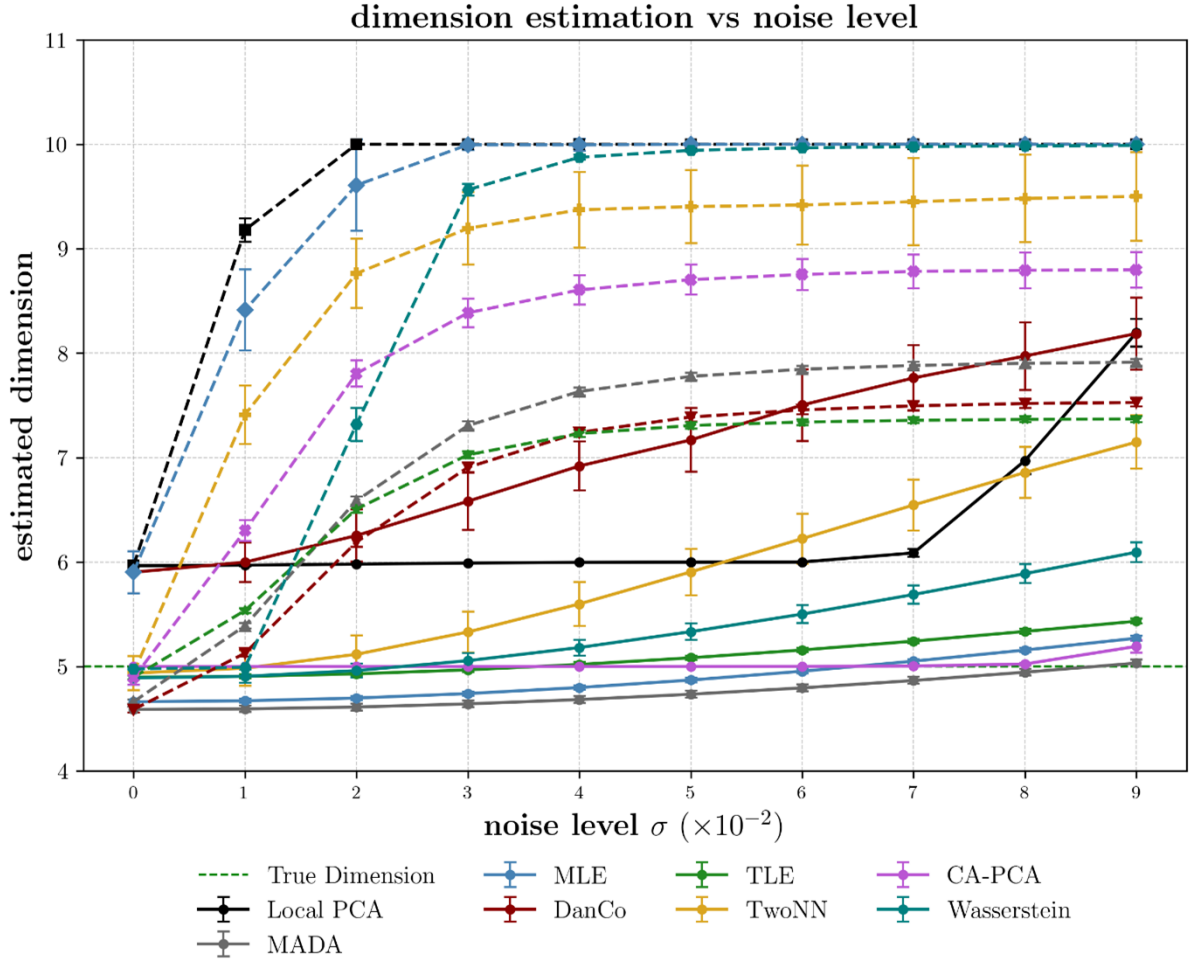


Figure 17: Noise Effect under High and Low Curvature.

We observe that when the data is noiseless, curvature size does not affect the estimates. However, as the noise level increases, estimates on manifolds with higher curvature deteriorate more rapidly in terms of both bias and variance compared to those on manifolds with lower curvature.

To mitigate the noise-curvature dilemma, some researchers have proposed shifting the focus away from recovering the true intrinsic dimension and instead estimating a dimension that is sufficient to describe the data at the chosen neighborhood scale. This notion, often referred to as the *effective dimension*, lacks a standardized definition but has gained attention as a practical alternative in noisy or high-curvature scenarios (Wang and Marron, 2008; Little

et al., 2009).

Another factor worth consideration is the data distribution on the manifold. Since most estimators assume at least local uniformity of the data, a highly irregular distribution is expected to also have significant impacts on their performance. For example, as illustrated in Figure 18a–18b, the flatness assumption is clearly violated in the neighborhood where the data density is low.

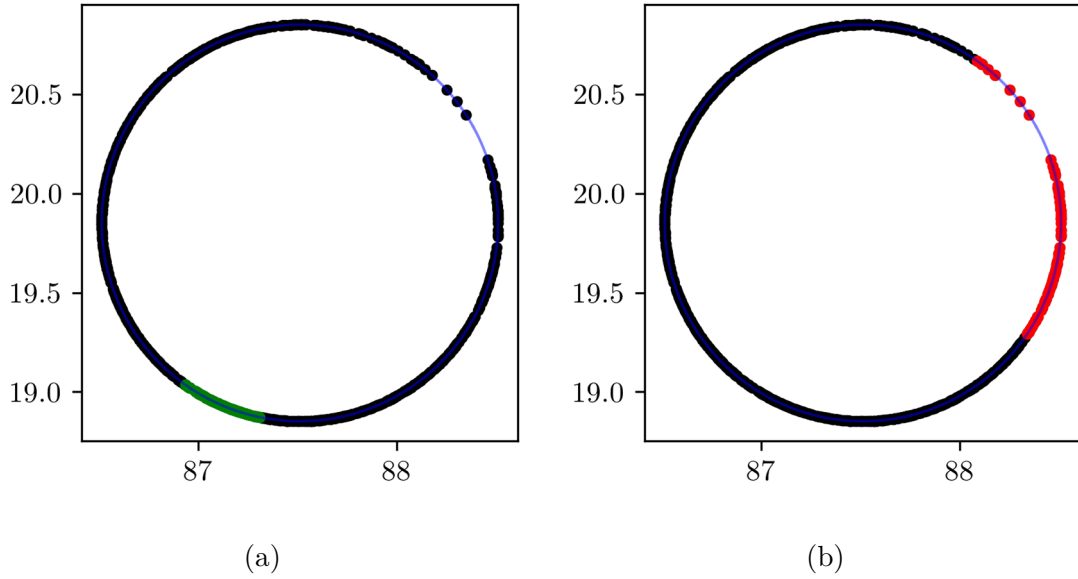


Figure 18: Effect of Data Distribution on Effective Neighborhood Sizes: A neighborhood taken from a region with high data density (green) is more local and thus better approximated by a linear space. In contrast, a neighborhood from a low-density region (red) spans a larger area, violating the local flatness assumption to a greater extent.

To better understand estimator robustness, it is more informative to compare their performance under both uniform and nonuniform distributions across a diverse set of manifolds. We defer this experiment to the next section, where we conduct a more systematic comparison of manifold dimension estimators across varying data distributions and manifold geometries.

4.2 Comparative Evaluation of Estimators

4.2.1 Choosing hyperparameter values

As we have shown, the performance of a manifold dimension estimator can vary significantly depending on its hyperparameter values, and the optimal setting may differ across experimental setups. Therefore, when comparing estimators across various manifolds, it is crucial to tune each estimator’s hyperparameters appropriately for each sample. However, in practice, the true intrinsic dimension is unknown, necessitating an algorithmic approach for per-sample hyperparameter tuning.

Motivated by the patterns observed in Figure 7, we propose an algorithm for selecting hyperparameter values that yield stable estimator performance. The procedure identifies regions of the hyperparameter space where the standard deviation of the estimates is minimized or when the mean estimate remains approximately constant. If no such stable region is detected, the algorithm defaults to reporting the average across all candidate hyperparameter values. The full algorithm is given below.

Algorithm 2 Choosing a Value of the Hyperparameter K (or α)

```
1: Input: Dimension estimates  $\hat{d}_\ell$  for  $\ell = 1, \dots, k_{\max}$ , corresponding to hyperparameters  
    $K_\ell$   
2: Output: Indices  $k_1$  and  $k_2$ , defining the selected stable range for  $K$   
3: Initialize  $k_1 \leftarrow 1$ ,  $k_2 \leftarrow 3$ ,  $k^* \leftarrow 0$   
4: Initialize  $s_{\min} \leftarrow \infty$ ,  $s_{\max} \leftarrow 0$   
5: for  $k = 1$  to  $k_{\max} - 2$  do  
6:    $s \leftarrow$  standard deviation of  $\hat{d}_k, \hat{d}_{k+1}, \hat{d}_{k+2}$   
7:   if  $s < s_{\min}$  then  
8:      $s_{\min} \leftarrow s$   
9:      $k^* \leftarrow k$   
10:  end if  
11:  if  $s > s_{\max}$  then  
12:     $s_{\max} \leftarrow s$   
13:  end if  
14: end for  
15: if  $s_{\max} > 1.25s_{\min}$  then  
16:    $k_1 \leftarrow k^*$ ,  $k_2 \leftarrow k^* + 2$   
17: else  
18:    $k_1 \leftarrow 1$ ,  $k_2 \leftarrow k_{\max}$   
19: end if  
20: return  $k_1, k_2$ 
```

4.2.2 Experimental setup

We generate 100 random samples of size n (both $n = 500$ and $n = 2000$ are considered) on 18 manifolds. Each manifold has a specific intrinsic dimension d (ranging from

$\{1, 2, 3, 5, 10, 20\}$), and a distinct geometry: sphere, ball, Gaussian density surface, deformed sphere, cylinder, helix, Swiss roll, Mobius strip, torus, or hyperbolic surface. The 18 manifolds are detailed in Table 8. For each manifold, points are sampled uniformly and also according to a skewed Beta distribution in the coordinate space, following [Campadelli et al. \(2015\)](#) as a natural way to model realistic nonuniformity. Both noiseless and noisy settings are considered, the latter obtained by adding Gaussian noise $N(0, I_p)$.

Table 8: The 18 manifolds considered.

Manifold	Description
M_{11}	A 5-dimensional sphere embedded in \mathbb{R}^{10} with $R = 1$.
M_{12}	A 10-dimensional sphere embedded in \mathbb{R}^{20} with $R = 1$.
M_{13}	A 20-dimensional sphere embedded in \mathbb{R}^{40} with $R = 1$.
M_{21}	A 5-dimensional ball embedded in \mathbb{R}^{10} with $R = 1$.
M_{22}	A 10-dimensional ball embedded in \mathbb{R}^{20} with $R = 1$.
M_{23}	A 20-dimensional ball embedded in \mathbb{R}^{40} with $R = 1$.
M_{31}	A 5-dimensional Gaussian density surface embedded in \mathbb{R}^{10} corresponding to $N(0, 0.25I_5)$.
M_{32}	A 10-dimensional Gaussian density surface embedded in \mathbb{R}^{20} corresponding to $N(0, 0.25I_{10})$.
M_{33}	A 20-dimensional Gaussian density surface embedded in \mathbb{R}^{40} corresponding to $N(0, 0.25I_{20})$.
M_{41}	A 3-dimensional deformed sphere with $c = 0.01$.
M_{42}	A 3-dimensional deformed sphere with $c = 0.1$.
M_{43}	A 3-dimensional deformed sphere with $c = 1$.
M_5	A 2-dimensional cylinder embedded in \mathbb{R}^4 .
M_6	A 1-dimensional helix embedded in \mathbb{R}^3 .
M_7	A 2-dimensional Swiss roll embedded in \mathbb{R}^4 .
M_8	A 2-dimensional Mobius strip embedded in \mathbb{R}^4 .
M_9	A 2-dimensional torus embedded in \mathbb{R}^4 .
M_{10}	A 2-dimensional hyperbolic surface embedded in \mathbb{R}^4 .

These 18 manifolds cover a wide variety of cases, including not only standard toy man-

ifolds commonly found in textbooks (M_5 to M_{10} ; see [Lee \(2003\)](#) for example), but also more challenging scenarios such as high-dimensional settings (M_{11} to M_{33}) and nonlinear embeddings (M_{41} to M_{43}).

Although many of these manifolds can be embedded in \mathbb{R}^{d+1} , we further apply a linear embedding into \mathbb{R}^{2d} to increase the potential for estimation error and to evaluate the robustness of the dimension estimators to this factor.

All computations were carried out on the Katana compute node at UNSW, configured via a batch script requesting a single chunk of 16 `x86_64` CPU cores and 16GB of RAM. Experiments were split into six tasks based on sample size, data distribution, and noise configuration. The total runtime required to generate the numerical results in the following section is approximately 120 hours, using Python version 3.6.8.

4.2.3 Results

The results are presented in six numerical tables, each comprising 18 rows (corresponding to the manifolds under study) and 8 columns (corresponding to the dimension estimators). Each cell reports the mean and standard deviation computed from 100 replications. For ease of presentation, the complete set of tables is relegated to the Appendix. In what follows, we restrict attention to the principal findings, which are examined through a detailed analysis of these tables.

Tables [25](#) and [26](#) present results on uniformly distributed, noiseless samples with respective sample sizes $n = 500$ and $n = 2000$. Beyond what is already known from individual factor analyses, we uncover the following insights.

Despite being excluded from many prior comparative studies, **Local PCA** consistently delivers strong performance when $n = 2000$, achieving the smallest bias and variance across almost all manifolds. However, it struggles when $n = 500$, particularly on manifolds

exhibiting complex curvature such as M_{41} to M_{43} .

CA-PCA further enhances this performance, not only when $n = 500$ but also on nonlinearly embedded manifolds like M_{41} to M_{43} . When $n = 2000$, it achieves near-unbiased estimates with zero variance, highlighting the strength of the PCA procedure.

Manifolds characterized by high curvature or complex geometry, such as M_{41} to M_{43} , present challenges for nearly all estimators when $n = 500$. While estimation improves when $n = 2000$, substantial bias remains for most methods.

High-dimensional manifolds tend to be underestimated. This trend, already observed with spheres previously, is further confirmed by results on balls and Gaussian density surfaces. Notable exceptions include **Local PCA**, **CA-PCA**, and **DanCo**, which handle high-dimensional cases better.

DanCo and **TwoNN** are also competitive. When $n = 500$, **TwoNN** achieves the best overall performance on low-dimensional manifolds but tends to underestimate dimensions in high-dimensional settings. **DanCo** performs slightly worse than **TwoNN** in low dimensions, in particular the deformed spheres, but gets better with high-dimensional manifolds. As the sample size increases from $n = 500$ to $n = 2000$, **TwoNN** shows greater improvement, becoming nearly unbiased on the manifolds where it performs well. However, it continues to face challenges with high-dimensional cases.

The remaining estimators, **MLE**, **MADA**, **TLE**, and **Wasserstein**, can be classified as bottom-tier. They all tend to underestimate the dimension of high-dimensional manifolds and struggle with nonlinearly embedded manifolds. Even for simple manifolds with $n = 2000$, a relatively large bias still persists.

Tables 27 and 28 present the experimental results under nonuniform data distributions, where nonuniformity is introduced by sampling from skewed Beta distributions in the pa-

parameter space. Our findings are summarized below.

The bias of all estimators increases under nonuniform sampling, especially for high-dimensional manifolds. However, an interesting exception occurs with M_{41} to M_{43} , where reduced bias is observed. We hypothesize that this is due to the complex curvature of these manifolds—that is, curvature varies significantly across different regions. Under nonuniform sampling, data tends to concentrate in one or a few high-density regions, potentially skipping areas with high curvature. This effectively simplifies the estimation task. Since we do not fully control the nonuniformity, such outcomes can occur.

With nonuniform data distributions, the underestimation problem for high-dimensional spheres and balls becomes significantly more severe. This time, even **Local PCA**, **DanCo**, and **CA-PCA** exhibit underestimation, although **CA-PCA** underestimates significantly less than **Local PCA**. For the PCA-based methods, this underestimation is intuitive: with nonuniform sampling, a neighborhood may include data points that are far from the center (see Figure 18b for example). The principal components corresponding to these distant points can dominate the local structure and pull down the estimated dimension.

Despite the increased difficulty introduced by nonuniformity, the best-performing estimators remain **Local PCA**, **CA-PCA**, **DanCo**, and **TwoNN**. However, in this setting, they perform well primarily on low-dimensional manifolds and also fail noticeably on high-dimensional ones.

Lastly, when additive Gaussian noise is present, as shown in Tables 29 and 30, the performance of all manifold dimension estimators deteriorates, primarily in the form of increased bias. In particular, our previous top performers—**Local PCA**, **DanCo**, **TwoNN**, and **CA-PCA**—exhibit the most significant degradation. Conversely, the previous underdogs—**MLE**, **MADA**, **TLE**, and **Wasserstein**—do not degrade as severely. A possible explanation is that these estimators do not rely on geometric information as delicate as the others

and therefore are less overwhelmed by noise.

Notably, increasing the sample size does not lead to performance improvement in the presence of noise. We also observe an apparent decrease in bias for high-dimensional manifolds compared to the noiseless case for most estimators; however, this reduction appears to be an artifact of the noise rather than a genuine improvement. This is evidenced by **TwoNN**, whose bias shifts from negative to positive as the sample size increases. For the same reason, although **MADA** appears to perform best on most manifolds with low to moderate dimensions, we believe this may not accurately reflect its true effectiveness.

5 Real-World Data Examples

In this section, we further assess the performance of the eight estimators on three real-world datasets, which are publicly available on the GitHub repository accompanying this article. As statisticians often prefer working in R, we provide results obtained from both the `Python` and R implementations (noting that the R implementation is considerably slower). For all three datasets, hyperparameters are selected according to Algorithm 2. We used the `Rdimtools` R package for **MLE**, **MADA**, and **TwoNN**. Although a **DanCo** implementation is available in this package, it appears to contain bugs and consistently returns an estimate of 2. Thus, we coded in R the five other estimators.

Since all three of our real-world datasets have ambient dimension much larger than the intrinsic dimension, as discussed in Section 3, **DanCo** is expected to struggle. Indeed, our experiments with **DanCo**, implemented in both `Python` and R, produced dimensionality estimates significantly higher than expected—often close to the ambient dimension. Moreover, the runtime was prohibitively long, particularly with the R implementation. These discrepancies led to results that differ substantially from those reported in prior studies such as [Ceruti et al. \(2014\)](#) and [Campadelli et al. \(2015\)](#), which in fact used the mechanisms

we introduced before, and thus no longer reflect pure **DanCo** outputs. Consequently, we cite their reported estimates (in brackets) for comparative purposes.

5.1 The ISOMAP Dataset

The ISOMAP dataset ([Balasubramanian and Schwartz, 2002](#)) consists of 698 grayscale images of a face sculpture, each of size 64×64 , as shown in Figure 19.

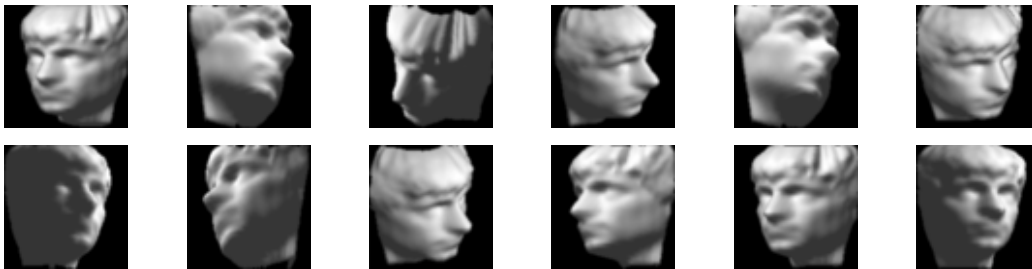


Figure 19: Sample images from ISOMAP dataset.

It is widely believed that the vectorized dataset has an intrinsic dimension of $d = 3$, corresponding to the two degrees of freedom for camera position and one for lighting condition.

Table 9 reports the estimation results, while Table 10 presents the computation times of each estimator for both the Python and R implementations.

Table 9: Dimension Estimate on ISOMAP

	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wass.
Python	10.05	5.86	4.34	(4.00)	4.70	3.49	25.71	4.13
R	10.10	6.04	4.36	(4.00)	4.71	3.61	25.85	4.46

Table 10: Time Used for Estimation on ISOMAP (in seconds)

	Local PCA	MADA	MLE	TLE	TwoNN	CA-PCA	Wass.
Python	35.26	7.83	1.55	11.72	0.05	42.56	9.26
R	129.10	120.02	120.51	1789.34	11.49	131.43	106.73

5.2 The MNIST Dataset

The now well known MNIST dataset ([LeCun et al., 2002](#)) contains images of handwritten digits from 0 to 9. From this, we extract 6,742 grayscale images of digit 1, each of size 28×28 , as shown in Figure 20.

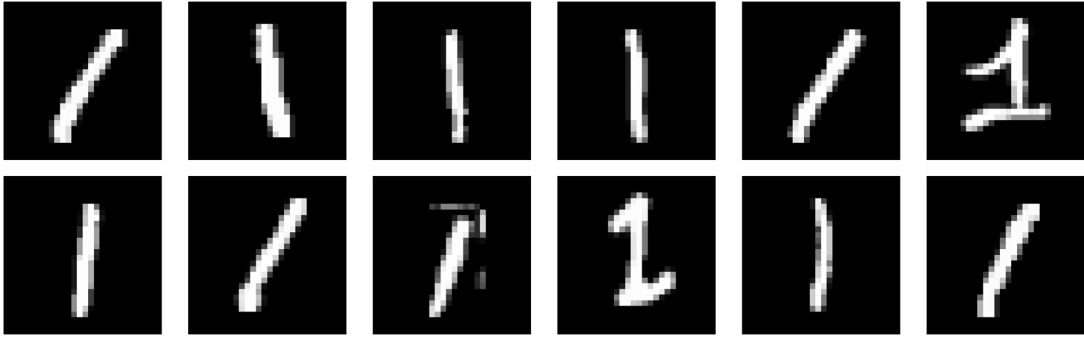


Figure 20: Sample images from MNIST dataset.

The manifold hypothesis is believed to hold for this subset, with an intrinsic dimension estimated to lie between 8 and 11.

Table 11 reports the estimation results, while Table 12 presents the computation times of each estimator for both the Python and R implementations.

Table 11: Dimension Estimate on MNIST

	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wass.
Python	15.41	8.67	8.97	(9.98)	7.92	12.98	84.66	9.49
R	15.65	8.67	8.97	(9.98)	5.42	13.20	84.15	9.60

Table 12: Time Used for Estimation on MNIST (in seconds)

	Local PCA	MADA	MLE	TLE	TwoNN	CA-PCA	Wass.
Python	251.94	146.68	14.10	99.37	3.19	429.07	163.94
R	1282.67	1763.15	1675.46	9536.70	162.71	16641.17	95155.58

5.3 The ISOLET Dataset

The ISOLET dataset ([Asuncion et al., 2007](#)) includes 6,240 vectors of length 616. Each vector represents frequency characteristics extracted from audio recordings of individuals pronouncing letters from the English alphabet, as illustrated in Figure 21.

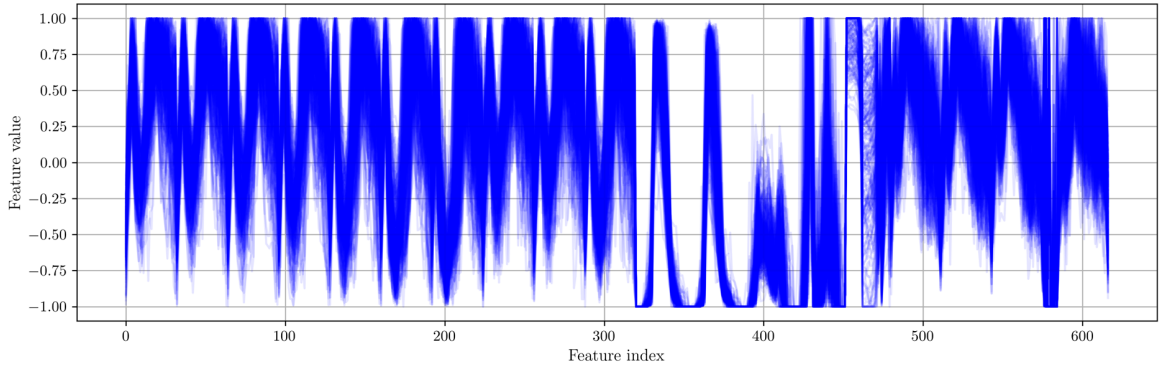


Figure 21: Frequency characteristics for letter “A” from ISOLET dataset.

The manifold hypothesis is again believed to apply, with an estimated intrinsic dimension between 16 and 22.

Table 13 reports the estimation results, while Table 14 presents the computation times of each estimator for both the Python and R implementations.

Table 13: Dimension Estimate on ISOLET

	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wass.
Python	34.44	13.63	14.29	(19.00)	12.58	9.11	54.22	13.38
R	34.70	13.64	14.28	(19.00)	12.57	9.82	54.00	13.36

Table 14: Time Used for Estimation on ISOLET (in seconds)

	Local PCA	MADA	MLE	TLE	TwoNN	CA-PCA	Wass.
Python	233.68	105.04	11.76	78.37	2.68	385.66	136.54
R	985.18	1062.51	958.78	4608.27	95.29	11976.68	53613.68

5.4 Analysis of the results for the three datasets

Among all estimators, (the adapted) **DanCo** produces estimates within the expected ranges for all three datasets. While **MLE**, **MADA**, **TLE**, **TwoNN**, and **Wasserstein** perform well on the **ISOMAP** and **MNIST** datasets, they underestimate the dimension of **ISOLET**, which is known to have higher intrinsic dimensionality. This behavior is consistent with our findings on the synthetic datasets.

In contrast, **Local PCA** and **CA-PCA**, despite being top performers on simulated data, perform considerably worse on these real-world datasets. This is likely because the underlying manifolds are embedded in spaces with much higher ambient dimensionality. Since both methods rely—at least in part—on identifying the number of significant eigenvalues, the accumulation of many small (but nonzero) eigenvalues can mislead the estimators and degrade their performance. This highlights the potential limitations of applying the two

methods in practice.

6 Conclusion

In this article, we reviewed the fundamental constructions of manifolds and their sampling necessary for the manifold dimension estimation problem, introduced eight representative estimators, and conducted a comprehensive empirical study of their performance. This is the most up-to-date survey on this important topic, featuring extensive controlled experiments and comparisons.

Perhaps the most valuable takeaway is for practitioners. Our experiments identify **CA-PCA** and **TwoNN** as the overall top performers, making them strong default choices for real-world datasets. However, caution is warranted: **TwoNN** tends to underestimate the dimension when the true intrinsic dimensionality is high, while **CA-PCA** may overestimate when the manifold is nonlinearly embedded in a high-dimensional ambient space. When these two estimators produce substantially different values, and computational cost is not a major constraint, we recommend considering **DanCo** for a third opinion. We also suggest exploring a wide range of neighborhood sizes for the estimators to help identify stable estimates.

Our experimental results lead to a perhaps unsurprising conclusion: for a problem of this complexity, simpler methods often perform better. Given the flexibility of the manifold model, any strong assumption made by an estimator about the underlying structure is more likely to be violated than satisfied. For researchers developing new manifold dimension estimators, we recommend moving beyond the extensively exploited flatness assumption. Instead, focusing on methods that explicitly incorporate curvature and are robust to practical challenges such as noise and non-uniform sampling may yield more broadly applicable and reliable results.

7 Acknowledgements

This research includes computations using the computational cluster Katana supported by Research Technology Services at UNSW Sydney.

This paper is part of Z. Bi's Ph.D. thesis.

8 Disclosure statement

The authors do not have any conflicts of interest to declare.

9 Data Availability Statement

All code and data used in this research is available through <https://github.com/loong-bi/manifold-dimension-estimation/tree/main>.

Appendix

Claim. A d -dimensional deformed sphere M fills the entire \mathbb{R}^{2d} ; in other words, it is not contained in any proper linear subspace.

Proof. It suffices to show the collection of tangent vectors of M spans \mathbb{R}^{2d} . From the definition of M , it has tangent vectors in the form of

$$\mathbf{v}_j(\mathbf{u}) = d\varphi^{-1}|_{\mathbf{u}}(\mathbf{e}_j) \in \mathbb{R}^{2d} \text{ with } \mathbf{u} \in (-\pi, \pi)^d,$$

where $\{\mathbf{e}_j\}_{j=1}^d$ are standard basis of \mathbb{R}^d ,

$$\begin{cases} v_{jj}(\mathbf{u}) = -2\pi cr \sin(2c\pi u_k) \cos(2\pi u_j) - [R + r \cos(2c\pi u_j)] \sin(2\pi u_j), \\ v_{j,j+d}(\mathbf{u}) = -2\pi cr \sin(2c\pi u_j) \sin(2\pi u_j) + [R + r \cos(2c\pi u_j)] \cos(2\pi u_j) \end{cases}$$

for $j = 1, \dots, d$ and $v_{ji}(\mathbf{u}) = v_{j,i+d}(\mathbf{u}) = 0$ for $i \neq j$. When $2c$ is an integer, we take $2d$ tangent vectors to be $\mathbf{v}_j(\pm \mathbf{e}_j)$ for $j = 1, \dots, d$. Otherwise, we set $\mathbf{v}_j(\pm \mathbf{e}_j/2)$. In both cases, one can verify that the resulting $2d$ tangent vectors are linearly independent. This establishes that the tangent vectors span \mathbb{R}^{2d} , and the proof is complete. \square

Tables of results

Table 15: Dimension estimation vs neighborhood size K (general).

K	5	10	20	30	40	50	100
Local PCA	3.42 ± 0.02	4.79 ± 0.02	5.00 ± 0.00	5.01 ± 0.00	5.06 ± 0.01	5.18 ± 0.01	5.96 ± 0.01
MADA	4.77 ± 0.16	5.64 ± 0.11	5.13 ± 0.07	4.96 ± 0.05	4.86 ± 0.04	4.79 ± 0.04	4.59 ± 0.03
MLE	6.53 ± 0.15	5.45 ± 0.07	5.07 ± 0.05	4.94 ± 0.03	4.87 ± 0.03	4.82 ± 0.03	4.66 ± 0.02
TLE	7.73 ± 0.15	5.99 ± 0.07	5.41 ± 0.05	5.23 ± 0.03	5.14 ± 0.03	5.07 ± 0.03	4.90 ± 0.02
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00

K	200	500	1000
Local PCA	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00
MADA	4.34 ± 0.02	3.76 ± 0.01	3.22 ± 0.01
MLE	4.47 ± 0.01	4.02 ± 0.01	3.60 ± 0.00
TLE	4.69 ± 0.01	4.18 ± 0.01	3.70 ± 0.00
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00

Table 16: Dimension estimation vs neighborhood size K (**DanCo**).

K	2	4	6	8	10	12	14
DanCo	2.00 ± 0.00	4.42 ± 0.23	4.81 ± 0.20	5.24 ± 0.18	5.90 ± 0.20	6.46 ± 0.32	6.56 ± 0.38

K	16	18	20
DanCo	6.49 ± 0.36	6.33 ± 0.32	6.21 ± 0.31

Table 17: Dimension estimation vs α (**Wasserstein**).

K	1.01	1.2	1.4	1.6	1.8	2	4
Wasserstein	0.14 ± 0.20	4.55 ± 0.30	4.87 ± 0.22	4.89 ± 0.14	4.88 ± 0.12	4.89 ± 0.11	4.89 ± 0.07

K	6	8	10
Wasserstein	4.88 ± 0.07	4.88 ± 0.06	4.89 ± 0.05

Table 18: Dimension estimation vs sample size n .

n	200	400	600	800	1000	2000	3000
Local PCA	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	5.99 ± 0.00	5.96 ± 0.01	5.35 ± 0.01	5.04 ± 0.00
MADA	3.80 ± 0.05	4.27 ± 0.04	4.44 ± 0.03	4.53 ± 0.03	4.59 ± 0.03	4.73 ± 0.02	4.79 ± 0.02
MLE	4.05 ± 0.03	4.42 ± 0.03	4.55 ± 0.02	4.62 ± 0.02	4.66 ± 0.02	4.77 ± 0.01	4.82 ± 0.01
DanCo	6.40 ± 0.56	6.61 ± 0.48	6.40 ± 0.34	6.09 ± 0.28	5.90 ± 0.20	5.36 ± 0.12	5.13 ± 0.11
TLE	4.23 ± 0.03	4.65 ± 0.03	4.78 ± 0.02	4.85 ± 0.02	4.90 ± 0.02	4.99 ± 0.01	5.02 ± 0.01
TwoNN	4.83 ± 0.41	4.88 ± 0.34	4.92 ± 0.27	4.92 ± 0.20	4.94 ± 0.16	4.97 ± 0.14	4.98 ± 0.11
CA-PCA	6.00 ± 0.00	5.17 ± 0.04	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.95 ± 0.14	4.91 ± 0.10	4.90 ± 0.09	4.88 ± 0.07	4.89 ± 0.06	4.89 ± 0.05	4.90 ± 0.04

n	4000	5000	10000
Local PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
MADA	4.82 ± 0.01	4.84 ± 0.01	4.90 ± 0.01
MLE	4.85 ± 0.01	4.87 ± 0.01	4.92 ± 0.01
DanCo	5.01 ± 0.08	4.94 ± 0.08	4.77 ± 0.05
TLE	5.04 ± 0.01	5.05 ± 0.01	5.08 ± 0.01
TwoNN	4.97 ± 0.10	4.98 ± 0.10	4.98 ± 0.06
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.90 ± 0.03	4.91 ± 0.03	5.06 ± 0.02

Table 19: Dimension estimation vs ambient dimension p .

p	6	10	15	20	25	30	35
Local PCA	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01
MADA	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03
MLE	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02
DanCo	6.00 ± 0.00	5.90 ± 0.20	5.89 ± 0.20	5.92 ± 0.19	5.91 ± 0.22	5.92 ± 0.19	5.92 ± 0.19
TLE	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02
TwoNN	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06

p	40	45	50
Local PCA	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01
MADA	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03
MLE	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02
DanCo	5.89 ± 0.22	5.89 ± 0.19	5.92 ± 0.18
TLE	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02
TwoNN	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06

Table 20: Dimension estimation vs curvature (size).

R	0.001	0.01	0.05	0.1	1	5	10
Local PCA	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01
MADA	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03
MLE	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02
DanCo	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20
TLE	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02
TwoNN	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06

R	20	50	100
Local PCA	5.96 ± 0.01	5.96 ± 0.01	5.96 ± 0.01
MADA	4.59 ± 0.03	4.59 ± 0.03	4.59 ± 0.03
MLE	4.66 ± 0.02	4.66 ± 0.02	4.66 ± 0.02
DanCo	5.90 ± 0.20	5.90 ± 0.20	5.90 ± 0.20
TLE	4.90 ± 0.02	4.90 ± 0.02	4.90 ± 0.02
TwoNN	4.94 ± 0.16	4.94 ± 0.16	4.94 ± 0.16
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.89 ± 0.06	4.89 ± 0.06	4.89 ± 0.06

Table 21: Dimension estimation vs intrinsic dimension d .

d	2	4	6	8	10	12	14
Local PCA	2.00 ± 0.00	4.91 ± 0.01	6.98 ± 0.00	8.98 ± 0.00	10.98 ± 0.01	12.96 ± 0.01	14.94 ± 0.01
MADA	2.03 ± 0.01	3.81 ± 0.02	5.31 ± 0.03	6.62 ± 0.04	7.80 ± 0.05	8.88 ± 0.05	9.89 ± 0.06
MLE	2.02 ± 0.01	3.85 ± 0.01	5.43 ± 0.02	6.82 ± 0.03	8.08 ± 0.03	9.24 ± 0.03	10.34 ± 0.04
DanCo	2.11 ± 0.06	4.18 ± 0.13	7.77 ± 0.44	9.55 ± 0.47	11.16 ± 0.53	12.92 ± 0.64	14.41 ± 0.58
TLE	2.16 ± 0.01	4.07 ± 0.02	5.66 ± 0.02	7.02 ± 0.03	8.22 ± 0.03	9.30 ± 0.03	10.29 ± 0.04
TwoNN	2.01 ± 0.09	3.95 ± 0.16	5.86 ± 0.26	7.62 ± 0.32	9.27 ± 0.33	10.86 ± 0.48	12.33 ± 0.57
CA-PCA	2.00 ± 0.00	4.00 ± 0.00	6.00 ± 0.00	8.02 ± 0.01	10.24 ± 0.02	12.59 ± 0.02	14.82 ± 0.01
Wasserstein	2.37 ± 0.06	4.02 ± 0.05	5.73 ± 0.07	7.32 ± 0.10	8.81 ± 0.10	10.19 ± 0.15	11.50 ± 0.16

d	16	18	20
Local PCA	16.91 ± 0.01	18.86 ± 0.01	20.80 ± 0.01
MADA	10.84 ± 0.06	11.71 ± 0.07	12.58 ± 0.07
MLE	11.35 ± 0.04	12.31 ± 0.05	13.23 ± 0.05
DanCo	16.22 ± 0.54	18.08 ± 0.52	19.86 ± 1.87
TLE	11.21 ± 0.04	12.06 ± 0.04	12.87 ± 0.04
TwoNN	13.81 ± 0.61	15.10 ± 0.62	16.43 ± 0.63
CA-PCA	16.93 ± 0.01	18.98 ± 0.00	20.99 ± 0.00
Wasserstein	12.76 ± 0.17	13.95 ± 0.17	15.09 ± 0.21

Table 22: Dimension estimation vs curvature (variation).

c	0.0001	0.001	0.01	0.05	0.1	0.2	0.5
Local PCA	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	5.97 ± 0.01
MADA	3.55 ± 0.02	3.55 ± 0.02	3.55 ± 0.02	3.55 ± 0.02	3.55 ± 0.02	3.58 ± 0.02	3.93 ± 0.04
MLE	3.39 ± 0.01	3.39 ± 0.01	3.39 ± 0.01	3.39 ± 0.01	3.39 ± 0.01	3.42 ± 0.01	3.58 ± 0.02
DanCo	3.26 ± 0.10	3.26 ± 0.10	3.26 ± 0.10	3.26 ± 0.10	3.28 ± 0.10	3.44 ± 0.11	3.66 ± 0.22
TLE	3.70 ± 0.01	3.70 ± 0.01	3.70 ± 0.01	3.70 ± 0.01	3.70 ± 0.01	3.72 ± 0.01	3.85 ± 0.03
TwoNN	3.02 ± 0.13	3.02 ± 0.13	3.02 ± 0.13	3.02 ± 0.13	3.03 ± 0.13	3.05 ± 0.13	2.93 ± 0.14
CA-PCA	5.41 ± 0.03	5.41 ± 0.03	5.41 ± 0.03	5.41 ± 0.03	5.41 ± 0.03	5.44 ± 0.04	5.66 ± 0.05
Wasserstein	3.39 ± 0.04	3.39 ± 0.04	3.39 ± 0.04	3.39 ± 0.04	3.39 ± 0.04	3.42 ± 0.04	3.43 ± 0.05

c	1	2	4
Local PCA	5.83 ± 0.03	6.00 ± 0.00	6.00 ± 0.00
MADA	3.56 ± 0.02	4.87 ± 0.03	4.49 ± 0.04
MLE	3.39 ± 0.01	4.30 ± 0.02	4.54 ± 0.03
DanCo	3.45 ± 0.11	4.20 ± 0.35	5.59 ± 0.39
TLE	3.64 ± 0.02	4.62 ± 0.02	4.68 ± 0.03
TwoNN	3.02 ± 0.12	3.05 ± 0.14	3.32 ± 0.15
CA-PCA	4.74 ± 0.06	5.99 ± 0.00	6.00 ± 0.00
Wasserstein	3.39 ± 0.05	3.70 ± 0.05	4.39 ± 0.07

Table 23: Dimension estimation vs Gaussian noise level σ ($R = 1$).

σ	0.00	0.01	0.02	0.03	0.04	0.05	0.06
Local PCA	5.96 ± 0.01	5.97 ± 0.01	5.98 ± 0.01	5.99 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00
MADA	4.59 ± 0.03	4.59 ± 0.03	4.61 ± 0.03	4.64 ± 0.03	4.68 ± 0.03	4.73 ± 0.03	4.79 ± 0.03
MLE	4.66 ± 0.02	4.67 ± 0.02	4.70 ± 0.02	4.74 ± 0.02	4.80 ± 0.02	4.87 ± 0.02	4.95 ± 0.02
DanCo	5.90 ± 0.20	6.00 ± 0.19	6.25 ± 0.25	6.58 ± 0.27	6.92 ± 0.23	7.17 ± 0.31	7.50 ± 0.34
TLE	4.90 ± 0.02	4.91 ± 0.02	4.93 ± 0.02	4.97 ± 0.02	5.02 ± 0.02	5.08 ± 0.02	5.16 ± 0.02
TwoNN	4.94 ± 0.16	4.98 ± 0.17	5.12 ± 0.18	5.33 ± 0.20	5.60 ± 0.21	5.90 ± 0.22	6.22 ± 0.24
CA-PCA	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00
Wasserstein	4.89 ± 0.06	4.91 ± 0.06	4.96 ± 0.07	5.06 ± 0.07	5.18 ± 0.08	5.33 ± 0.08	5.50 ± 0.08

σ	0.07	0.08	0.09
Local PCA	6.09 ± 0.04	6.97 ± 0.13	8.20 ± 0.13
MADA	4.87 ± 0.03	4.94 ± 0.03	5.03 ± 0.03
MLE	5.05 ± 0.02	5.16 ± 0.02	5.27 ± 0.02
DanCo	7.76 ± 0.31	7.97 ± 0.32	8.19 ± 0.35
TLE	5.24 ± 0.02	5.34 ± 0.02	5.43 ± 0.02
TwoNN	6.55 ± 0.24	6.86 ± 0.24	7.15 ± 0.26
CA-PCA	5.00 ± 0.00	5.02 ± 0.01	5.19 ± 0.06
Wasserstein	5.69 ± 0.09	5.89 ± 0.09	6.09 ± 0.09

Table 24: Dimension estimation vs Gaussian noise level σ^2 ($R = 0.1$).

σ	0.00	0.01	0.02	0.03	0.04	0.05	0.06
Local PCA	5.96 ± 0.01	9.18 ± 0.11	10.00 ± 0.00	10.00 ± 0.00	10.00 ± 0.00	10.00 ± 0.00	10.00 ± 0.00
MADA	4.59 ± 0.03	5.13 ± 0.03	6.19 ± 0.04	6.91 ± 0.05	7.24 ± 0.04	7.39 ± 0.04	7.46 ± 0.04
MLE	4.66 ± 0.02	5.39 ± 0.02	6.59 ± 0.04	7.31 ± 0.04	7.63 ± 0.04	7.78 ± 0.04	7.85 ± 0.03
DanCo	5.90 ± 0.20	8.41 ± 0.39	9.61 ± 0.44	9.99 ± 0.05	9.99 ± 0.06	10.00 ± 0.00	10.00 ± 0.00
TLE	4.90 ± 0.02	5.54 ± 0.02	6.51 ± 0.04	7.03 ± 0.03	7.23 ± 0.03	7.31 ± 0.03	7.34 ± 0.03
TwoNN	4.94 ± 0.16	7.41 ± 0.28	8.76 ± 0.33	9.19 ± 0.35	9.37 ± 0.36	9.40 ± 0.35	9.42 ± 0.38
CA-PCA	4.98 ± 0.01	4.99 ± 0.00	7.32 ± 0.16	9.56 ± 0.05	9.87 ± 0.02	9.94 ± 0.01	9.97 ± 0.01
Wasserstein	4.89 ± 0.06	6.30 ± 0.10	7.80 ± 0.13	8.39 ± 0.14	8.61 ± 0.14	8.70 ± 0.14	8.75 ± 0.15

σ	0.07	0.08	0.09
Local PCA	10.00 ± 0.00	10.00 ± 0.00	10.00 ± 0.00
MADA	7.50 ± 0.04	7.52 ± 0.04	7.53 ± 0.04
MLE	7.88 ± 0.03	7.90 ± 0.03	7.91 ± 0.03
DanCo	10.00 ± 0.00	10.00 ± 0.00	10.00 ± 0.00
TLE	7.36 ± 0.03	7.37 ± 0.03	7.37 ± 0.03
TwoNN	9.45 ± 0.42	9.48 ± 0.42	9.50 ± 0.42
CA-PCA	9.98 ± 0.01	9.98 ± 0.01	9.99 ± 0.00
Wasserstein	8.78 ± 0.16	8.79 ± 0.17	8.80 ± 0.17

Table 25: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 500$ uniform samples. Bold indicates the estimate with minimal MSE.

96

Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
$M_{11}(5)$	5.03 ± 0.01	4.65 ± 0.05	4.78 ± 0.04	6.18 ± 0.26	5.06 ± 0.04	4.90 ± 0.28	5.00 ± 0.00	4.88 ± 0.08
$M_{12}(10)$	10.10 ± 0.01	7.93 ± 0.10	8.25 ± 0.07	11.16 ± 0.37	8.48 ± 0.07	9.17 ± 0.53	10.21 ± 0.02	8.67 ± 0.16
$M_{13}(20)$	19.76 ± 0.02	12.89 ± 0.15	13.54 ± 0.10	20.13 ± 0.37	13.33 ± 0.10	15.90 ± 0.90	20.68 ± 0.02	14.73 ± 0.35
$M_{21}(5)$	5.00 ± 0.00	4.24 ± 0.05	4.41 ± 0.04	4.97 ± 0.13	4.58 ± 0.04	4.64 ± 0.30	4.91 ± 0.03	4.52 ± 0.07
$M_{22}(10)$	9.99 ± 0.01	7.45 ± 0.08	7.75 ± 0.06	10.03 ± 0.17	7.88 ± 0.05	8.79 ± 0.51	9.93 ± 0.03	8.28 ± 0.26
$M_{23}(20)$	19.10 ± 0.04	12.42 ± 0.15	13.05 ± 0.12	20.01 ± 0.28	12.79 ± 0.11	15.53 ± 0.92	19.93 ± 0.02	14.44 ± 0.50
$M_{31}(5)$	5.00 ± 0.00	4.21 ± 0.06	4.38 ± 0.06	4.97 ± 0.13	4.54 ± 0.05	4.55 ± 0.27	5.00 ± 0.01	4.40 ± 0.07
$M_{32}(10)$	10.00 ± 0.00	7.29 ± 0.08	7.55 ± 0.07	9.55 ± 0.39	7.56 ± 0.06	8.28 ± 0.51	10.01 ± 0.02	7.82 ± 0.15
$M_{33}(20)$	19.43 ± 0.04	12.04 ± 0.13	12.69 ± 0.10	19.59 ± 0.54	12.14 ± 0.09	14.75 ± 0.83	20.07 ± 0.05	13.66 ± 0.28
$M_{41}(3)$	5.83 ± 0.02	3.53 ± 0.04	3.36 ± 0.03	3.42 ± 0.12	3.67 ± 0.03	3.05 ± 0.17	3.24 ± 0.02	3.54 ± 0.05
$M_{42}(3)$	5.83 ± 0.02	3.53 ± 0.04	3.36 ± 0.03	3.44 ± 0.12	3.67 ± 0.03	3.05 ± 0.17	3.25 ± 0.02	3.55 ± 0.05
$M_{43}(3)$	5.52 ± 0.05	3.45 ± 0.04	3.36 ± 0.03	3.58 ± 0.14	3.63 ± 0.03	3.06 ± 0.16	3.44 ± 0.03	3.58 ± 0.04
$M_5(2)$	2.00 ± 0.00	1.86 ± 0.03	1.91 ± 0.02	2.16 ± 0.03	2.01 ± 0.02	1.98 ± 0.12	1.95 ± 0.01	2.35 ± 0.06
$M_6(1)$	2.75 ± 0.04	1.11 ± 0.02	1.09 ± 0.01	1.01 ± 0.02	1.21 ± 0.01	1.00 ± 0.07	1.00 ± 0.00	2.21 ± 0.05
$M_7(2)$	2.99 ± 0.01	2.96 ± 0.07	2.14 ± 0.03	2.18 ± 0.03	2.30 ± 0.02	1.97 ± 0.12	2.78 ± 0.04	2.52 ± 0.14
$M_8(2)$	2.00 ± 0.00	1.83 ± 0.02	1.88 ± 0.01	2.18 ± 0.03	1.97 ± 0.01	1.96 ± 0.12	1.93 ± 0.01	2.29 ± 0.07
$M_9(2)$	3.00 ± 0.00	2.34 ± 0.03	2.19 ± 0.02	2.16 ± 0.03	2.35 ± 0.02	2.00 ± 0.13	2.00 ± 0.01	2.15 ± 0.11
$M_{10}(2)$	2.07 ± 0.01	2.11 ± 0.03	2.11 ± 0.02	2.18 ± 0.03	2.14 ± 0.01	1.99 ± 0.12	1.95 ± 0.01	2.44 ± 0.07

Table 26: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 2,000$ uniform samples. Bold indicates the estimate with minimal MSE.

Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
$M_{11}(5)$	5.00 ± 0.00	4.86 ± 0.03	4.92 ± 0.02	6.54 ± 0.17	5.15 ± 0.02	4.97 ± 0.14	5.00 ± 0.00	4.89 ± 0.04
$M_{12}(10)$	10.01 ± 0.00	8.68 ± 0.05	8.89 ± 0.04	10.92 ± 0.14	9.20 ± 0.04	9.40 ± 0.24	10.00 ± 0.00	9.06 ± 0.15
$M_{13}(20)$	19.70 ± 0.01	14.61 ± 0.08	15.12 ± 0.06	19.97 ± 0.24	15.13 ± 0.06	16.74 ± 0.52	20.17 ± 0.01	15.52 ± 0.14
$M_{21}(5)$	5.00 ± 0.00	4.54 ± 0.03	4.62 ± 0.02	5.00 ± 0.06	4.76 ± 0.02	4.76 ± 0.14	4.93 ± 0.01	4.56 ± 0.03
$M_{22}(10)$	9.99 ± 0.00	8.20 ± 0.05	8.41 ± 0.04	10.00 ± 0.11	8.61 ± 0.04	8.96 ± 0.22	9.91 ± 0.02	8.45 ± 0.08
$M_{23}(20)$	19.00 ± 0.03	14.09 ± 0.08	14.59 ± 0.06	19.98 ± 0.15	14.51 ± 0.05	16.17 ± 0.43	19.79 ± 0.03	14.94 ± 0.13
$M_{31}(5)$	5.00 ± 0.00	4.41 ± 0.02	4.50 ± 0.02	5.11 ± 0.08	4.62 ± 0.02	4.68 ± 0.13	5.00 ± 0.00	4.42 ± 0.03
$M_{32}(10)$	10.00 ± 0.00	7.86 ± 0.05	8.09 ± 0.04	9.71 ± 0.31	8.17 ± 0.03	8.61 ± 0.25	10.00 ± 0.00	7.99 ± 0.07
$M_{33}(20)$	19.38 ± 0.02	13.65 ± 0.07	14.09 ± 0.05	20.87 ± 0.37	13.71 ± 0.05	15.41 ± 0.38	19.99 ± 0.01	14.37 ± 0.14
$M_{41}(3)$	3.11 ± 0.01	3.22 ± 0.02	3.19 ± 0.01	2.99 ± 0.05	3.40 ± 0.01	3.03 ± 0.09	3.00 ± 0.00	3.30 ± 0.02
$M_{42}(3)$	3.12 ± 0.01	3.22 ± 0.02	3.19 ± 0.01	3.01 ± 0.05	3.41 ± 0.01	3.03 ± 0.09	3.00 ± 0.00	3.30 ± 0.02
$M_{43}(3)$	3.17 ± 0.01	3.23 ± 0.02	3.20 ± 0.01	3.76 ± 0.09	3.39 ± 0.01	3.03 ± 0.09	3.01 ± 0.00	3.28 ± 0.02
$M_5(2)$	2.00 ± 0.00	1.95 ± 0.01	1.97 ± 0.01	2.14 ± 0.02	2.03 ± 0.01	1.97 ± 0.05	1.97 ± 0.00	2.26 ± 0.06
$M_6(1)$	1.00 ± 0.00	1.01 ± 0.01	1.02 ± 0.00	1.00 ± 0.00	1.10 ± 0.00	1.00 ± 0.03	1.00 ± 0.00	2.09 ± 0.04
$M_7(2)$	2.03 ± 0.00	1.99 ± 0.01	1.99 ± 0.01	2.16 ± 0.01	2.09 ± 0.01	1.97 ± 0.06	1.98 ± 0.01	2.44 ± 0.03
$M_8(2)$	2.00 ± 0.00	1.94 ± 0.01	1.97 ± 0.01	2.15 ± 0.01	2.02 ± 0.01	1.97 ± 0.05	1.96 ± 0.00	1.70 ± 0.10
$M_9(2)$	2.00 ± 0.00	2.09 ± 0.01	2.08 ± 0.01	2.15 ± 0.02	2.19 ± 0.01	2.00 ± 0.06	2.00 ± 0.00	2.29 ± 0.10
$M_{10}(2)$	2.00 ± 0.00	2.06 ± 0.01	2.05 ± 0.01	2.16 ± 0.01	2.10 ± 0.01	1.98 ± 0.06	1.98 ± 0.00	2.31 ± 0.05

Table 27: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 500$ nonuniform samples. Bold indicates the estimate with minimal MSE.

	Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
\mathcal{M}	$M_{11}(5)$	4.94 ± 0.02	4.52 ± 0.07	4.71 ± 0.06	5.69 ± 0.15	4.21 ± 0.06	4.91 ± 0.25	5.23 ± 0.04	4.59 ± 0.14
	$M_{12}(10)$	5.79 ± 0.11	4.60 ± 0.09	4.85 ± 0.08	6.60 ± 0.19	4.23 ± 0.06	5.87 ± 0.31	9.10 ± 0.32	4.94 ± 0.17
	$M_{13}(20)$	5.89 ± 0.12	4.59 ± 0.10	4.83 ± 0.10	6.58 ± 0.19	4.22 ± 0.07	5.73 ± 0.35	15.75 ± 0.76	4.84 ± 0.21
	$M_{21}(5)$	4.99 ± 0.01	4.81 ± 0.08	4.95 ± 0.06	6.04 ± 0.14	4.17 ± 0.05	5.01 ± 0.24	5.48 ± 0.06	4.77 ± 0.14
	$M_{22}(10)$	5.98 ± 0.10	4.96 ± 0.12	5.20 ± 0.11	6.46 ± 0.16	4.23 ± 0.06	6.00 ± 0.31	10.34 ± 0.34	4.79 ± 0.33
	$M_{23}(20)$	5.97 ± 0.12	4.94 ± 0.11	5.17 ± 0.09	6.64 ± 0.18	4.23 ± 0.06	5.99 ± 0.33	18.32 ± 0.80	4.75 ± 0.34
	$M_{31}(5)$	5.00 ± 0.00	4.57 ± 0.06	4.76 ± 0.05	5.50 ± 0.17	4.75 ± 0.04	4.90 ± 0.30	5.05 ± 0.02	4.72 ± 0.10
	$M_{32}(10)$	10.00 ± 0.00	7.44 ± 0.07	7.90 ± 0.07	10.04 ± 0.27	7.60 ± 0.05	8.91 ± 0.42	10.38 ± 0.07	8.16 ± 0.19
	$M_{33}(20)$	18.81 ± 0.06	11.84 ± 0.13	12.72 ± 0.11	21.74 ± 1.29	11.76 ± 0.12	15.31 ± 0.89	21.74 ± 0.30	13.77 ± 0.38
	$M_{41}(3)$	4.03 ± 0.07	3.17 ± 0.04	3.20 ± 0.03	3.40 ± 0.12	3.31 ± 0.03	3.01 ± 0.19	2.95 ± 0.01	3.33 ± 0.08
	$M_{42}(3)$	4.05 ± 0.07	3.17 ± 0.04	3.20 ± 0.03	3.40 ± 0.12	3.31 ± 0.03	3.01 ± 0.19	2.95 ± 0.01	3.33 ± 0.08
	$M_{43}(3)$	3.88 ± 0.06	3.04 ± 0.04	3.15 ± 0.03	3.38 ± 0.12	3.22 ± 0.03	3.02 ± 0.20	3.00 ± 0.01	3.34 ± 0.07
	$M_5(2)$	2.00 ± 0.00	1.83 ± 0.02	1.89 ± 0.01	2.23 ± 0.05	1.93 ± 0.01	2.00 ± 0.13	1.88 ± 0.01	2.33 ± 0.08
	$M_6(1)$	1.61 ± 0.08	1.43 ± 0.15	1.15 ± 0.02	1.00 ± 0.00	1.22 ± 0.01	0.99 ± 0.06	1.03 ± 0.01	2.23 ± 0.05
	$M_7(2)$	2.59 ± 0.04	2.67 ± 0.11	2.29 ± 0.05	2.19 ± 0.04	2.28 ± 0.02	1.99 ± 0.11	2.13 ± 0.02	2.40 ± 0.15
	$M_8(2)$	2.00 ± 0.00	1.70 ± 0.03	1.78 ± 0.02	2.26 ± 0.05	1.83 ± 0.02	2.00 ± 0.13	1.86 ± 0.02	2.23 ± 0.08
	$M_9(2)$	2.69 ± 0.05	2.20 ± 0.04	2.17 ± 0.02	2.20 ± 0.03	2.22 ± 0.02	1.99 ± 0.12	1.98 ± 0.01	2.07 ± 0.15
	$M_{10}(2)$	2.02 ± 0.01	1.94 ± 0.03	2.02 ± 0.02	2.22 ± 0.04	2.00 ± 0.01	2.01 ± 0.12	1.93 ± 0.01	2.38 ± 0.08

Table 28: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 2,000$ nonuniform samples. Bold indicates the estimate with minimal MSE.

Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
$M_{11}(5)$	4.99 \pm 0.00	4.98 \pm 0.03	5.03 \pm 0.03	6.63 \pm 0.14	4.68 \pm 0.02	5.01 \pm 0.14	5.06 \pm 0.01	4.71 \pm 0.06
$M_{12}(10)$	6.42 \pm 0.06	5.46 \pm 0.06	5.71 \pm 0.05	6.97 \pm 0.06	5.05 \pm 0.03	6.42 \pm 0.19	8.29 \pm 0.16	5.54 \pm 0.09
$M_{13}(20)$	6.41 \pm 0.06	5.42 \pm 0.05	5.66 \pm 0.04	6.94 \pm 0.06	5.03 \pm 0.03	6.36 \pm 0.19	11.44 \pm 0.26	5.47 \pm 0.10
$M_{21}(5)$	5.00 \pm 0.00	5.17 \pm 0.04	5.18 \pm 0.02	6.38 \pm 0.13	4.67 \pm 0.02	5.05 \pm 0.14	5.14 \pm 0.02	4.80 \pm 0.07
$M_{22}(10)$	6.60 \pm 0.07	5.80 \pm 0.07	6.03 \pm 0.06	7.01 \pm 0.04	5.10 \pm 0.04	6.58 \pm 0.17	9.19 \pm 0.14	5.90 \pm 0.20
$M_{23}(20)$	6.51 \pm 0.06	5.74 \pm 0.06	5.97 \pm 0.05	6.99 \pm 0.04	5.08 \pm 0.03	6.48 \pm 0.18	13.84 \pm 0.31	5.71 \pm 0.08
$M_{31}(5)$	5.00 \pm 0.00	4.75 \pm 0.03	4.89 \pm 0.02	5.59 \pm 0.09	4.89 \pm 0.02	4.94 \pm 0.15	5.01 \pm 0.00	4.73 \pm 0.04
$M_{32}(10)$	10.00 \pm 0.00	8.32 \pm 0.05	8.66 \pm 0.04	10.02 \pm 0.06	8.47 \pm 0.03	9.20 \pm 0.27	10.09 \pm 0.02	8.49 \pm 0.09
$M_{33}(20)$	18.79 \pm 0.02	13.70 \pm 0.07	14.43 \pm 0.06	21.77 \pm 0.86	13.55 \pm 0.06	16.21 \pm 0.42	20.47 \pm 0.07	14.71 \pm 0.18
$M_{41}(3)$	3.02 \pm 0.00	3.13 \pm 0.02	3.12 \pm 0.01	3.18 \pm 0.06	3.21 \pm 0.01	3.01 \pm 0.08	3.00 \pm 0.00	3.18 \pm 0.04
$M_{42}(3)$	3.02 \pm 0.00	3.13 \pm 0.02	3.12 \pm 0.01	3.18 \pm 0.06	3.21 \pm 0.01	3.01 \pm 0.08	3.00 \pm 0.00	3.18 \pm 0.04
$M_{43}(3)$	3.01 \pm 0.00	3.13 \pm 0.02	3.13 \pm 0.01	3.41 \pm 0.07	3.17 \pm 0.01	3.02 \pm 0.08	3.00 \pm 0.00	3.17 \pm 0.04
$M_5(2)$	2.00 \pm 0.00	2.02 \pm 0.01	2.03 \pm 0.01	2.18 \pm 0.01	2.05 \pm 0.01	2.01 \pm 0.06	1.97 \pm 0.00	2.24 \pm 0.06
$M_6(1)$	1.01 \pm 0.00	1.19 \pm 0.07	1.06 \pm 0.00	1.00 \pm 0.00	1.11 \pm 0.00	1.00 \pm 0.03	1.00 \pm 0.00	2.09 \pm 0.04
$M_7(2)$	2.02 \pm 0.00	2.42 \pm 0.06	2.14 \pm 0.02	2.18 \pm 0.02	2.14 \pm 0.01	2.01 \pm 0.06	2.00 \pm 0.00	2.46 \pm 0.04
$M_8(2)$	2.00 \pm 0.00	1.99 \pm 0.01	2.02 \pm 0.01	2.18 \pm 0.03	2.04 \pm 0.01	2.00 \pm 0.06	1.95 \pm 0.01	2.20 \pm 0.06
$M_9(2)$	2.01 \pm 0.00	2.11 \pm 0.01	2.10 \pm 0.01	2.18 \pm 0.02	2.13 \pm 0.01	2.01 \pm 0.06	2.00 \pm 0.00	2.35 \pm 0.07
$M_{10}(2)$	2.00 \pm 0.00	2.10 \pm 0.01	2.11 \pm 0.01	2.21 \pm 0.03	2.10 \pm 0.01	2.01 \pm 0.06	1.98 \pm 0.00	2.29 \pm 0.05

Table 29: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 500$ uniform samples perturbed by Gaussian noise. Bold indicates the estimate with minimal MSE.

8	Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
	$M_{11}(5)$	7.19 ± 0.07	5.22 ± 0.07	5.59 ± 0.05	7.49 ± 0.24	5.80 ± 0.05	6.99 ± 0.35	7.35 ± 0.32	6.06 ± 0.11
	$M_{12}(10)$	13.27 ± 0.08	9.17 ± 0.10	9.67 ± 0.08	14.69 ± 0.41	9.61 ± 0.07	11.98 ± 0.61	19.06 ± 0.06	10.70 ± 0.23
	$M_{13}(20)$	25.73 ± 0.11	15.41 ± 0.15	16.37 ± 0.12	34.36 ± 1.01	15.32 ± 0.11	20.84 ± 1.13	35.93 ± 0.07	18.90 ± 0.56
	$M_{21}(5)$	8.77 ± 0.08	5.19 ± 0.06	5.65 ± 0.05	8.22 ± 0.22	5.74 ± 0.05	7.52 ± 0.43	9.89 ± 0.03	6.07 ± 0.09
	$M_{22}(10)$	15.20 ± 0.12	8.99 ± 0.10	9.53 ± 0.08	14.65 ± 0.57	9.35 ± 0.07	12.25 ± 0.63	19.50 ± 0.04	10.88 ± 0.39
	$M_{23}(20)$	25.89 ± 0.12	15.25 ± 0.16	16.24 ± 0.13	35.01 ± 1.25	15.11 ± 0.12	21.03 ± 0.97	36.23 ± 0.08	18.97 ± 0.76
	$M_{31}(5)$	5.04 ± 0.02	4.45 ± 0.06	4.74 ± 0.06	6.02 ± 0.15	4.91 ± 0.06	5.66 ± 0.33	5.06 ± 0.03	4.80 ± 0.08
	$M_{32}(10)$	10.00 ± 0.00	7.51 ± 0.09	7.82 ± 0.07	10.16 ± 0.28	7.79 ± 0.06	8.93 ± 0.50	10.22 ± 0.08	8.18 ± 0.15
	$M_{33}(20)$	19.55 ± 0.03	12.30 ± 0.14	12.99 ± 0.10	20.42 ± 0.57	12.39 ± 0.09	15.28 ± 0.76	21.15 ± 0.24	14.08 ± 0.28
	$M_{41}(3)$	5.90 ± 0.01	3.59 ± 0.04	3.47 ± 0.03	4.04 ± 0.14	3.78 ± 0.03	3.47 ± 0.18	3.44 ± 0.03	3.68 ± 0.06
	$M_{42}(3)$	5.90 ± 0.01	3.59 ± 0.04	3.48 ± 0.03	4.07 ± 0.14	3.78 ± 0.03	3.48 ± 0.18	3.46 ± 0.03	3.69 ± 0.06
	$M_{43}(3)$	5.79 ± 0.03	3.57 ± 0.05	3.58 ± 0.03	4.73 ± 0.19	3.83 ± 0.03	3.88 ± 0.19	3.89 ± 0.05	3.85 ± 0.05
	$M_5(2)$	3.93 ± 0.01	2.44 ± 0.05	2.68 ± 0.04	4.00 ± 0.00	2.77 ± 0.03	3.85 ± 0.21	3.52 ± 0.10	3.02 ± 0.07
	$M_6(1)$	2.97 ± 0.01	1.58 ± 0.04	1.67 ± 0.03	2.94 ± 0.09	2.01 ± 0.03	2.77 ± 0.15	1.73 ± 0.04	2.51 ± 0.06
	$M_7(2)$	2.99 ± 0.01	2.96 ± 0.07	2.17 ± 0.03	2.27 ± 0.04	2.33 ± 0.03	2.18 ± 0.11	2.78 ± 0.04	2.55 ± 0.14
	$M_8(2)$	3.94 ± 0.02	2.22 ± 0.03	2.45 ± 0.03	4.00 ± 0.02	2.55 ± 0.02	3.75 ± 0.20	3.20 ± 0.10	2.78 ± 0.09
	$M_9(2)$	3.00 ± 0.00	2.40 ± 0.03	2.35 ± 0.02	2.89 ± 0.13	2.55 ± 0.02	2.98 ± 0.15	2.10 ± 0.01	2.31 ± 0.12
	$M_{10}(2)$	3.19 ± 0.06	2.23 ± 0.04	2.35 ± 0.03	3.86 ± 0.16	2.37 ± 0.02	3.29 ± 0.18	2.15 ± 0.04	2.67 ± 0.07

Table 30: Mean dimension estimates (\pm standard deviation) for 18 manifolds (Table 8) with true dimension d , based on 100 replicates of $n = 2,000$ uniform samples perturbed by Gaussian noise. Bold indicates the estimate with minimal MSE.

Manifold (d)	Local PCA	MADA	MLE	DanCo	TLE	TwoNN	CA-PCA	Wasserstein
$M_{11}(5)$	8.40 ± 0.03	5.90 ± 0.04	6.22 ± 0.03	10.00 ± 0.02	6.38 ± 0.03	7.77 ± 0.20	9.78 ± 0.04	6.53 ± 0.05
$M_{12}(10)$	14.56 ± 0.05	10.50 ± 0.05	10.98 ± 0.05	14.89 ± 0.21	10.96 ± 0.04	13.03 ± 0.32	19.60 ± 0.02	11.83 ± 0.20
$M_{13}(20)$	26.17 ± 0.05	18.08 ± 0.09	19.00 ± 0.08	36.49 ± 0.90	17.99 ± 0.07	22.68 ± 0.64	35.38 ± 0.04	20.44 ± 0.21
$M_{21}(5)$	9.61 ± 0.02	6.14 ± 0.04	6.53 ± 0.03	10.00 ± 0.00	6.58 ± 0.03	8.25 ± 0.24	9.99 ± 0.00	7.06 ± 0.11
$M_{22}(10)$	16.21 ± 0.06	10.46 ± 0.06	11.01 ± 0.05	15.02 ± 0.13	10.84 ± 0.04	13.33 ± 0.30	19.83 ± 0.01	11.77 ± 0.13
$M_{23}(20)$	26.38 ± 0.06	18.00 ± 0.10	18.94 ± 0.08	37.51 ± 0.96	17.81 ± 0.08	22.81 ± 0.58	35.68 ± 0.04	20.46 ± 0.20
$M_{31}(5)$	5.87 ± 0.05	4.81 ± 0.03	5.05 ± 0.02	6.97 ± 0.09	5.18 ± 0.02	6.34 ± 0.17	5.08 ± 0.02	5.27 ± 0.11
$M_{32}(10)$	10.00 ± 0.00	8.22 ± 0.05	8.50 ± 0.04	10.49 ± 0.33	8.55 ± 0.03	9.47 ± 0.27	10.09 ± 0.02	8.48 ± 0.07
$M_{33}(20)$	19.52 ± 0.02	14.03 ± 0.06	14.52 ± 0.05	22.82 ± 0.62	14.07 ± 0.05	16.13 ± 0.37	20.27 ± 0.04	14.90 ± 0.14
$M_{41}(3)$	3.85 ± 0.02	3.32 ± 0.02	3.36 ± 0.01	4.46 ± 0.12	3.58 ± 0.01	3.93 ± 0.10	3.03 ± 0.00	3.55 ± 0.02
$M_{42}(3)$	3.89 ± 0.02	3.33 ± 0.02	3.37 ± 0.01	4.52 ± 0.13	3.59 ± 0.01	3.96 ± 0.10	3.04 ± 0.01	3.56 ± 0.02
$M_{43}(3)$	4.55 ± 0.03	3.45 ± 0.02	3.55 ± 0.01	5.04 ± 0.19	3.73 ± 0.01	4.47 ± 0.14	3.53 ± 0.04	3.78 ± 0.03
$M_5(2)$	4.00 ± 0.00	3.26 ± 0.02	3.46 ± 0.02	4.00 ± 0.00	3.49 ± 0.02	4.00 ± 0.12	3.98 ± 0.00	3.10 ± 0.16
$M_6(1)$	2.85 ± 0.01	1.66 ± 0.01	1.88 ± 0.01	3.00 ± 0.00	2.18 ± 0.01	3.01 ± 0.08	2.54 ± 0.01	2.47 ± 0.06
$M_7(2)$	2.13 ± 0.01	2.01 ± 0.01	2.04 ± 0.01	2.39 ± 0.02	2.14 ± 0.01	2.50 ± 0.07	1.99 ± 0.01	2.49 ± 0.03
$M_8(2)$	4.00 ± 0.00	3.01 ± 0.03	3.33 ± 0.02	4.00 ± 0.00	3.37 ± 0.02	3.97 ± 0.11	3.96 ± 0.01	2.63 ± 0.16
$M_9(2)$	3.52 ± 0.02	2.25 ± 0.01	2.41 ± 0.01	3.98 ± 0.07	2.56 ± 0.01	3.54 ± 0.10	2.17 ± 0.02	2.69 ± 0.08
$M_{10}(2)$	3.87 ± 0.01	2.45 ± 0.02	2.72 ± 0.02	4.00 ± 0.00	2.80 ± 0.02	3.76 ± 0.10	3.07 ± 0.04	2.76 ± 0.04

References

- Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. *arXiv preprint arXiv:2012.13255*, 2020.
- Giovanni S Alberti, Johannes Hertrich, Matteo Santacesaria, and Silvia Sciutto. Manifold learning by mixture models of vaes for inverse problems. *Journal of Machine Learning Research*, 25(202):1–35, 2024.
- Laurent Amsaleg, Oussama Chelly, Teddy Furon, Stéphane Girard, Michael E Houle, Ken-ichi Kawarabayashi, and Michael Nett. Extreme-value-theoretic estimation of local intrinsic dimensionality. *Data Mining and Knowledge Discovery*, 32(6):1768–1805, 2018.
- Laurent Amsaleg, Oussama Chelly, Michael E Houle, Ken-Ichi Kawarabayashi, Miloš Radovanović, and Weeris Treeratanajaru. Intrinsic dimensionality estimation within tight localities. In *Proceedings of the 2019 SIAM international conference on data mining*, pages 181–189. SIAM, 2019.
- Arthur Asuncion, David Newman, et al. Uci machine learning repository, 2007.
- Anil Aswani, Peter Bickel, and Claire Tomlin. Regression on manifolds: Estimation of the exterior derivative. *The Annals of Statistics*, 39(1):48–81, 2011.
- Jonathan Bac and Andrei Zinovyev. Local intrinsic dimensionality estimators based on concentration of measure. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.
- Jonathan Bac, Evgeny M Mirkes, Alexander N Gorban, Ivan Tyukin, and Andrei Zinovyev. Scikit-dimension: a python package for intrinsic dimension estimation. *Entropy*, 23(10):1368, 2021.

- Mukund Balasubramanian and Eric L Schwartz. The isomap algorithm and topological stability. *Science*, 295(5552):7–7, 2002.
- Robert Bennett. The intrinsic dimensionality of signal collections. *IEEE Transactions on Information Theory*, 15(5):517–525, 1969.
- James AD Binnie, John Harvey, Jakub Malinowski, Ka Man Yim, et al. A survey of dimension estimation methods. *arXiv preprint arXiv:2507.13887*, 2025.
- Adam Block, Zeyu Jia, Yury Polyanskiy, and Alexander Rakhlin. Intrinsic dimension estimation using wasserstein distance. *Journal of Machine Learning Research*, 23(313):1–37, 2022.
- Charles Bouveyron, Gilles Celeux, and Stéphane Girard. Intrinsic dimension estimation by maximum likelihood in isotropic probabilistic pca. *Pattern Recognition Letters*, 32(14):1706–1713, 2011.
- Jörg Bruske and Gerald Sommer. Intrinsic dimensionality estimation with optimally topology preserving maps. *IEEE Transactions on pattern analysis and machine intelligence*, 20(5):572–575, 2002.
- Francesco Camastra and Antonino Staiano. Intrinsic dimension estimation: Advances and open problems. *Information Sciences*, 328:26–41, 2016.
- Francesco Camastra and Alessandro Vinciarelli. Estimating the intrinsic dimension of data with a fractal-based method. *IEEE Transactions on pattern analysis and machine intelligence*, 24(10):1404–1407, 2002.
- Paola Campadelli, Elena Casiraghi, Claudio Ceruti, and Alessandro Rozza. Intrinsic dimension estimation: Relevant techniques and a benchmark framework. *Mathematical Problems in Engineering*, 2015(1):759567, 2015.

- Kevin M Carter, Raviv Raich, and Alfred O Hero III. On local intrinsic dimension estimation and its applications. *IEEE Transactions on Signal Processing*, 58(2):650–663, 2009.
- Claudio Ceruti, Simone Bassis, Alessandro Rozza, Gabriele Lombardi, Elena Casiraghi, and Paola Campadelli. Danco: An intrinsic dimensionality estimator exploiting angle and norm concentration. *Pattern Recognition*, 47(8):2569–2581, 2014. ISSN 0031-3203.
- Ngai Hang Chan, Ye Lu, and Chun Yip Yau. Factor modelling for high-dimensional time series: inference and model selection. *Journal of Time Series Analysis*, 38(2):285–307, 2017.
- Ming-Yen Cheng and Hau-tieng Wu. Local linear regression on manifolds and its geometric interpretation. *Journal of the American Statistical Association*, 108(504):1421–1434, 2013.
- Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
- Jose A Costa and Alfred O Hero. Geodesic entropic graphs for dimension and entropy estimation in manifold learning. *IEEE Transactions on Signal Processing*, 52(8):2210–2221, 2004a.
- Jose A Costa and Alfred O Hero. Learning intrinsic dimension and intrinsic entropy of high-dimensional datasets. In *2004 12th European Signal Processing Conference*, pages 369–372. IEEE, 2004b.
- Persi Diaconis, Susan Holmes, and Mehrdad Shahshahani. Sampling from a manifold. In *Advances in modern statistical theory and applications: a Festschrift in honor of Morris L. Eaton*, volume 10, pages 102–126. Institute of Mathematical Statistics, 2013.

- Richard Mansfield Dudley. The speed of mean glivenko-cantelli convergence. *The Annals of Mathematical Statistics*, 40(1):40–50, 1969.
- Elena Facco, Maria d’Errico, Alex Rodriguez, and Alessandro Laio. Estimating the intrinsic dimension of datasets by a minimal neighborhood information. *Scientific reports*, 7(1):12140, 2017.
- Kenneth Falconer. *Fractal geometry: mathematical foundations and applications*. John Wiley & Sons, 2013.
- Mingyu Fan, Hong Qiao, and Bo Zhang. Intrinsic dimension estimation of manifolds by incising balls. *Pattern Recognition*, 42(5):780–787, 2009.
- Mingyu Fan, Nannan Gu, Hong Qiao, and Bo Zhang. Intrinsic dimension estimation of data by principal component analysis. *arXiv preprint arXiv:1002.2050*, 2010.
- Amir Massoud Farahmand, Csaba Szepesvári, and Jean-Yves Audibert. Manifold-adaptive dimension estimation. In *Proceedings of the 24th international conference on Machine learning*, pages 265–272, 2007.
- Vitaly V Fedorchuk. The fundamentals of dimension theory. In *General Topology I: Basic Concepts and Constructions Dimension Theory*, pages 91–192. Springer, 1990.
- Daniel Floryan and Michael D Graham. Data-driven discovery of intrinsic dynamics. *Nature Machine Intelligence*, 4(12):1113–1120, 2022.
- Keinosuke Fukunaga and David R Olsen. An algorithm for finding intrinsic dimensionality of data. *IEEE Transactions on computers*, 100(2):176–183, 1971.
- James E. Gentle. *Random Number Generation and Monte Carlo Methods*. Springer, New York, 2nd edition, 2004. doi: 10.1007/978-0-387-27605-5.

- Anna C Gilbert and Kevin O’Neill. Ca-pca: Manifold dimension estimation, adapted for curvature. *SIAM Journal on Mathematics of Data Science*, 7(1):355–383, 2025.
- Marina Gomtsyan, Nikita Mokrov, Maxim Panov, and Yury Yanovich. Geometry-aware maximum likelihood estimation of intrinsic dimension. In *Asian Conference on Machine Learning*, pages 1126–1141. PMLR, 2019.
- Peter Grassberger and Itamar Procaccia. Measuring the strangeness of strange attractors. *Physica D: nonlinear phenomena*, 9(1-2):189–208, 1983.
- Mithun Das Gupta and Thomas S Huang. Regularized maximum likelihood for intrinsic dimension estimation. *arXiv preprint arXiv:1203.3483*, 2012.
- Matthias Hein and Jean-Yves Audibert. Intrinsic dimensionality estimation of submanifolds in rd. In *Proceedings of the 22nd international conference on Machine learning*, pages 289–296, 2005.
- Kerstin Johnsson, Charlotte Soneson, and Magnus Fontes. Low bias local intrinsic dimension estimation from expected simplex skewness. *IEEE transactions on pattern analysis and machine intelligence*, 37(1):196–202, 2014.
- Jochen Einbecka Zakiah Kalantana and Z Einbecka. Intrinsic dimensionality estimation for high-dimensional data sets: New approaches for the computation of correlation dimension. *J. Emerg. Technol. Web Intell*, 5(2):91–97, 2013.
- Balázs Kégl. Intrinsic dimension estimation using packing numbers. *Advances in neural information processing systems*, 15, 2002.
- Eamonn Keogh and Abdullah Mueen. Curse of dimensionality. In *Encyclopedia of machine learning*, pages 257–258. Springer, 2011.

- Jisu Kim, Alessandro Rinaldo, and Larry Wasserman. Minimax rates for estimating the dimension of a manifold. *Journal of Computational Geometry*, 10(1):42–95, 2019.
- Matthäus Kleindessner and Ulrike Luxburg. Dimensionality estimation without distances. In *Artificial Intelligence and Statistics*, pages 471–479. PMLR, 2015.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 2002.
- John M Lee. Smooth manifolds. In *Introduction to smooth manifolds*, pages 1–29. Springer, 2003.
- John M Lee. *Introduction to Riemannian manifolds*, volume 2. Springer, 2018.
- Elizaveta Levina and Peter Bickel. Maximum likelihood estimation of intrinsic dimension. *Advances in neural information processing systems*, 17, 2004.
- Liang Liao, Yanning Zhang, Stephen John Maybank, and Zhoufeng Liu. Intrinsic dimension estimation via nearest constrained subspace classifier. *Pattern recognition*, 47(3):1485–1493, 2014.
- Uzu Lim, Harald Oberhauser, and Vidit Nanda. Tangent space and dimension estimation with the wasserstein distance. *SIAM Journal on Applied Algebra and Geometry*, 8(3):650–685, 2024.
- Anna V Little, Yoon-Mo Jung, and Mauro Maggioni. Multiscale estimation of intrinsic dimensionality of data sets. In *AAAI fall symposium: manifold learning and its applications*, volume 9, page 04, 2009.
- Gabriele Lombardi. Intrinsic dimensionality estimation techniques. *MATLAB Central File Exchange*, 2020.
- Kanti V Mardia and Peter E Jupp. *Directional statistics*. John Wiley & Sons, 2009.

- George Marsaglia. Choosing a point from the surface of a sphere. *The Annals of Mathematical Statistics*, 43(2):645–646, 1972.
- Marina Meilă and Hanyu Zhang. Manifold learning: What, how, and why. *Annual Review of Statistics and Its Application*, 11(1):393–417, 2024.
- Tim P. Morris, Ian R. White, and Michael J. Crowther. Using simulation studies to evaluate statistical methods. *Statistical Methods in Medical Research*, 28(9):2747–2764, 2019. doi: 10.1177/0962280218817802.
- John Nash. The imbedding problem for riemannian manifolds. *Annals of mathematics*, 63(1):20–63, 1956.
- Karl W Pettis, Thomas A Bailey, Anil K Jain, and Richard C Dubes. An intrinsic dimensionality estimator from near-neighbor information. *IEEE Transactions on pattern analysis and machine intelligence*, pages 25–37, 1979.
- Christian P. Robert and George Casella. *Monte Carlo Statistical Methods*. Springer, New York, 2nd edition, 2004. doi: 10.1007/978-1-4757-4145-2.
- Alessandro Rozza, Gabriele Lombardi, Marco Rosa, Elena Casiraghi, and Paola Campadelli. Idea: Intrinsic dimension estimation algorithm. In *International Conference on Image Analysis and Processing*, pages 433–442. Springer, 2011.
- Alessandro Rozza, Gabriele Lombardi, Claudio Ceruti, Elena Casiraghi, and Paola Campadelli. Novel high intrinsic dimensionality estimators. *Machine learning*, 89(1):37–65, 2012.
- Kumar Sricharan, Raviv Raich, and Alfred O Hero. Optimized intrinsic dimension estimator using nearest neighbor graphs. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 5418–5421. IEEE, 2010.

- Erik Thordsen and Erich Schubert. Abid: angle based intrinsic dimensionality—theory and analysis. *Information Systems*, 108:101989, 2022.
- Gerard V Trunk. Statistical estimation of the intrinsic dimensionality of data collections. *Information and Control*, 12(5):508–525, 1968.
- Peter J. Verwee and Robert P. W. Duin. An evaluation of intrinsic dimensionality estimators. *IEEE Transactions on pattern analysis and machine intelligence*, 17(1):81–86, 1995.
- Cédric Villani. *Topics in optimal transportation*, volume 58. American Mathematical Soc., 2021.
- Xiaohui Wang and JS Marron. A scale-based approach to finding effective dimensionality in manifold learning. *Electronic Journal of Statistics*, 2, 2008.
- Xin Yang, Sebastien Michea, and Hongyuan Zha. Conical dimension as an intrinsic dimension estimator and its applications. In *Proceedings of the 2007 SIAM International Conference on Data Mining*, pages 169–179. SIAM, 2007.
- Kisung You and Dennis Shung. Rdimtools: An r package for dimension reduction and intrinsic dimension estimation. *Software Impacts*, 14:100414, 2022.
- Kisung You, Changhee Suh, and Dennis Shung. *Rdimtools: Dimension Reduction and Estimation Methods*, 2022. URL <https://CRAN.R-project.org/package=Rdimtools>. R package version 1.1.2.
- Zhi-Hua Zhou. *Machine learning*. Springer nature, 2021.