

TOAST: TOPOLOGICAL ALGORITHM FOR SINGULARITY TRACKING

Anonymous authors

Paper under double-blind review

ABSTRACT

1 The manifold hypothesis, which assumes that data lie on or close to an unknown
 2 manifold of low intrinsic dimensionality, is a staple of modern machine learning
 3 research. However, recent work has shown that real-world data exhibit distinct
 4 non-manifold structures, which result in singularities that can lead to erroneous
 5 conclusions about the data. Detecting such singularities is therefore crucial as a
 6 precursor to interpolation and inference tasks. We address detecting singularities
 7 by developing (i) *persistent local homology*, a new topology-driven framework
 8 for quantifying the intrinsic dimension of a data set locally, and (ii) *Euclidicity*, a
 9 topology-based multi-scale measure for assessing the ‘manifoldness’ of individual
 10 points. We show that our approach can reliably identify singularities of complex
 11 spaces, while also capturing singular structures in real-world data sets.

12 1 INTRODUCTION

13 The ever-increasing amount and complexity of real-world data necessitate the development of new
 14 methods to extract less complex—but still *meaningful*—representations of the underlying data. One
 15 approach to this problem is via dimensionality reduction techniques, where the data is assumed to
 16 be of strictly lower dimension than its number of features. Traditional algorithms in this field such
 17 as PCA are restricted to linear descriptions of data, and are therefore of limited use for complex,
 18 non-linear data sets that often appear in practice. By contrast, non-linear dimensionality reduc-
 19 tion algorithms, such as UMAP (McInnes et al., 2018), *t*-SNE (van der Maaten & Hinton, 2008),
 20 or autoencoders (Kingma & Welling, 2019) share one common assumption: the underlying data is
 21 supposed to be close to a manifold with small intrinsic dimension, i.e. while the input data may have
 22 a large ambient dimension N , there is a n -dimensional manifold with $n \ll N$ that best describes the
 23 data. For some data sets, this *manifold hypothesis* is appropriate: certain natural images are known
 24 to be well-described by a manifold, for instance (Carlsson, 2009), enabling the use of specialised
 25 autoencoders for visualisation (Moor et al., 2020). However, recent research shows evidence that
 26 the manifold hypothesis does not necessarily hold for complex data sets (Brown et al., 2022), and
 27 that manifold learning techniques tend to fail for non-manifold data (Rieck & Leitte, 2015; Scoccola
 28 & Perea, 2022). These failures are typically the result of *singularities*, i.e. regions of a space that
 29 violate the properties of a manifold. For example, the ‘pinched torus,’ an object obtained by com-
 30 pressing a neighbourhood of a random point in a torus to a single point, fails to satisfy the manifold
 31 hypothesis at the ‘pinch point:’ this point, unlike all other points of the ‘pinched torus,’ does *not*
 32 have a neighbourhood homeomorphic to \mathbb{R}^2 (see Fig. 1 for an illustration).

33 Since singularities—unlike outliers that arise from incorrect labels, for example—may carry relevant
 34 information (Jakubowski et al., 2020), we address the shortcomings of existing dimensionality re-
 35 duction methods by assuming an agnostic view on any given data set. Instead of trying to prescribe
 36 the rigid requirements of a manifold, we consider intrinsic dimensionality to be a fundamentally
 37 *local phenomenon*: we permit dimensionality to vary across points in the data set, and, more im-
 38 portantly, across the *scale* of locality to be considered. The only assumption we make is that the
 39 data is of significantly lower dimension than the dimension of the ambient space. This perspective
 40 enables us to assess the deviation of individual points from idealised non-singular spaces, resulting
 41 in a measure of the *Euclidicity* of a point. Our method is based on a local version of topological data
 42 analysis (TDA), a method from computational topology that is capable of quantifying the shape of
 43 a data set on multiple scales (Edelsbrunner & Harer, 2010).

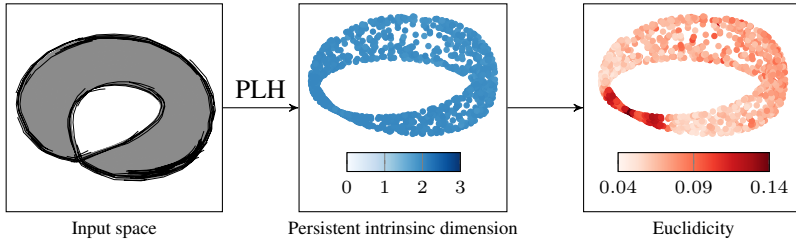


Figure 1: Overview of our method. Using *persistent local homology* (PLH), we derive a *persistent intrinsic dimension* and, subsequently, a *Euclidity* score that measures the deviation from a space to a Euclidean model space. Here, *Euclidity* highlights the singularity at the ‘pinch point.’ Please refer to Section 4 for more details.

44 **Our contributions.** We present a *universal framework for detecting singular regions in data*. This
 45 framework is agnostic with respect to geometric or stochastic properties of the underlying data and
 46 only requires a notion of intrinsic dimension of neighbourhoods. Our approach is based on a novel
 47 formulation of persistent local homology (PLH), a multi-parameter tool that detects the shape of
 48 local neighbourhoods of a given point in the data set, making use of multiple scales of locality.
 49 We employ PLH in two different capacities: (i) We use PLH to estimate the intrinsic dimension
 50 of a point locally. This enables us to assess how complex a given data set is, both in terms of the
 51 magnitude of the intrinsic dimension and in terms of the variance of its intrinsic dimension across
 52 individual points. (ii) Given the intrinsic dimension of the neighbourhood of a point, we use PLH to
 53 measure *Euclidity*, a novel quantity that we define to measure the deviation of a point from being
 54 Euclidean. We also provide theoretical guarantees on the approximation quality for certain classes of
 55 spaces and show the utility of our proposed method experimentally on several data sets.

56 2 BACKGROUND: PERSISTENT HOMOLOGY AND STRATIFIED SPACES

57 We first provide an overview of persistent homology and stratified spaces, as well as their relation
 58 to *local homology*. The former concept constitutes a generic framework for assessing complex data
 59 at multiple scales by measuring its topological characteristics such as ‘holes’ and ‘voids,’ while the
 60 latter will subsequently serve as a general setting to describe singularities, in which our framework
 61 admits advantageous properties.

62 **Persistent homology.** Persistent homology is a method for computing topological features at dif-
 63 ferent scales, capturing an intrinsic notion of relevance in terms of spatial scale parameters. Given
 64 a finite metric space (\mathbb{X}, d) , the *Vietoris–Rips complex* at step t is defined as the abstract simplicial
 65 complex $\mathcal{V}(\mathbb{X}, t)$, in which an abstract k -simplex (x_0, \dots, x_k) of points in \mathbb{X} is spanned if and only
 66 if $d(x_i, x_j) \leq t$ for all $0 \leq i \leq j \leq k$.¹ For $t_1 \leq t_2$, the inclusions $\mathcal{V}(\mathbb{X}, t_1) \hookrightarrow \mathcal{V}(\mathbb{X}, t_2)$ yield
 67 a filtration, i.e. a sequence of nested simplicial complexes, which we denote by $\mathcal{V}(\mathbb{X}, \bullet)$. Applying
 68 the i th homology functor to this collection of spaces and inclusions between them induces maps
 69 on the homology level $f_i^{t_1, t_2} : H_i(\mathcal{V}(\mathbb{X}, t_1)) \rightarrow H_i(\mathcal{V}(\mathbb{X}, t_2))$ for any $t_1 \leq t_2$. The i th *persistent*
 70 *homology* (PH) of \mathbb{X} with respect to the Vietoris–Rips construction is defined to be the collection
 71 of all these i th homology groups, together with the respective induced maps between them, and
 72 denoted by $\text{PH}_i(\mathbb{X}; \mathcal{V})$. PH can therefore be viewed as a tool that keeps track of topological fea-
 73 tures such as holes and voids on multiple scales. For a more comprehensive introduction to PH in
 74 the context of machine learning, see Hensel et al. (2021). The so-called ‘creation’ and ‘destruc-
 75 tion’ times of these features are summarised in a *persistence diagram* $\mathcal{D} \subset \mathbb{R} \times \mathbb{R} \cup \{\infty\}$, where
 76 any point $(b, d) \in \mathcal{D}$ corresponds to a homology class that arises at filtration step b , and lasts un-
 77 til filtration step d . The difference $|d - b|$ is referred to as the lifetime or eponymous *persistence*
 78 of this homology class. There are several distance measures for comparing persistence diagrams,
 79 one of them being the bottleneck distance d_B . For two persistence diagrams $\mathcal{D}, \mathcal{D}'$, it is defined as
 80 $d_B(\mathcal{D}, \mathcal{D}') := \inf_{\gamma} \sup_{x \in \mathcal{D}} \|x - \gamma(x)\|_{\infty}$, where γ ranges over all bijections between \mathcal{D} and \mathcal{D}' .

¹For readers familiar with persistent homology, we depart from the usual convention of using ϵ as the threshold parameter since we will require it to denote the scale of our persistent local homology calculations.

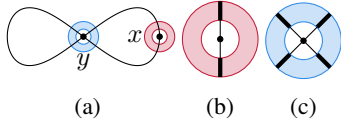


Figure 2: (a): Non-manifold space. (b): Annulus around a regular point x . (c): Annulus around a singular point. The neighbourhood around y is different from all others.

thus intrinsically capable of describing a wider class of spaces, we argue that stratified spaces are the right tool to analyse real-world data. Subsequently, we define stratified spaces in the setting of simplicial complexes. A stratified simplicial complex² of dimension 0 is a finite set of points with the discrete topology. A stratified simplicial complex of dimension n is an n -dimensional simplicial complex X , together with a filtration of closed subcomplexes $X = X_n \supset X_{n-1} \supset X_{n-2} \supset \dots \supset X_{-1} = \emptyset$ such that $X_i \setminus X_{i-1}$ is an i -dimensional manifold for all i , and such that every point $x \in X$ possesses a distinguished local neighbourhood $U \cong \mathbb{R}^k \times c^\circ L$ in X , where L is a compact stratified simplicial complex of dimension $n - k - 1$ and c° refers to the open cone construction (see Appendix A.1). If X is a manifold, then independently of the point under consideration, L is given by a sphere since for a manifold, *any* point admits a local neighbourhood that is homeomorphic to \mathbb{R}^n . This observation will serve as the primary motivation for our *Euclidicity* measure in Section 4.2.

Local homology. Local homology serves as a tool to quantify topological properties of infinitesimal small neighbourhoods of a fixed point. For a topological space X and $x \in X$, its i th local homology group is defined as $H_i(X, X \setminus x) := \varinjlim_U H_i(X, X \setminus U)$, where the direct system is given by the induced maps on homology that arise via the inclusion of (small) neighbourhoods of x .³ When X is a simplicial complex, we may view x as a vertex in X , using subdivision if necessary. Its *star* $\text{St}(x)$ is defined to be the union of simplices in X that have x as a face, whereas its *link* $\text{Lk}(x)$ consists of all simplices in $\text{St}(x)$ that do not have x as a face. Using excision and the long exact homology sequence (see Appendix A.3), we have

$$H_i(X, X \setminus x) \cong \tilde{H}_{i-1}(\text{Lk}(x)). \quad (1)$$

The **key takeaway** here is that the homology of $\text{Lk}(x)$ will usually differ from the homology of a sphere, once $\text{Lk}(x)$ is not homotopy-equivalent to a sphere. For example, when x is an isolated singularity in a stratified simplicial complex X of dimension n , then its distinguished neighbourhood is given by $U \cong c^\circ L$. Thus, $\text{Lk}(x) = L$ and $H_i(X, X \setminus x) = \tilde{H}_{i-1}(L)$ by Eq. (1), which is usually different from $\tilde{H}_{i-1}(S^{n-1})$, when x does not admit a Euclidean neighbourhood. This observation motivates and justifies using local homology for detecting non-Euclidean neighbourhoods.

3 RELATED WORK

Methods from topological data analysis have recently attracted much attention in machine learning, particularly due to persistent homology, which captures global topological properties of the underlying data set on different scales. We give a brief overview of existing methods in the emerging field of topology-driven singularity detection, outlining the differences to our approach below. Several works already assume a local perspective on homology to detect information about the intrinsic dimensionality of the data or the presence of certain singularities. Rieck et al. (2020) define persistent intersection homology via known stratifications, whereas Fasy & Wang (2016) and Bendich (2008), for instance, both present persistent versions of local homology. By contrast, Stolz et al. (2020) follow a different approach, where local homology is approximated as the absolute homology of

²Here, we actually mean the *geometric realisation* of the corresponding simplicial complex; by abuse of notation we may denote both objects by the term ‘simplicial complex.’

³Heuristically, a local homology class can be thought of as a homology class of an infinitesimal small punctured neighbourhood of a point.

128 a small annulus of a given neighbourhood, resulting in an algorithm for geometric anomaly detec-
 129 tion (which requires knowing the intrinsic dimension of the data set). Bendich et al. (2007) employ
 130 persistence vineyards, i.e. continuous families of persistence diagrams, to assess the local homology
 131 of a point in a stratified space, whereas Dey et al. (2014) use local homology to estimate the (global)
 132 intrinsic dimension of hidden, possibly noisy manifolds. While manifold learning is concerned with
 133 the development of algorithms that extract geometric information under the assumption that the
 134 given data lie on a manifold, Brown et al. (2022) recently introduced the idea to assume data spaces
 135 to consist of a *union of manifolds*. Intrinsic dimension is thus allowed to vary across connected
 136 components of the data space, but singularities are excluded under this assumption, whereas our
 137 framework detects the correct intrinsic dimension for large classes of singular spaces. Birdal et al.
 138 (2021) define a global persistent homology dimension for describing neural networks; our persistent
 139 intrinsic dimension, by contrast, is *local* and may thus change across different points in the data set.

140 **Key differences to existing approaches.** Our approach crucially differs from existing approaches
 141 in essential components. In comparison to all aforementioned contributions, we capture additional
 142 local geometric information: *we consider multiple scales of locality in a persistent framework for*
 143 *local homology*. Concerning the overall construction, Stolz et al. (2020) is the closest to our method.
 144 However, the authors assume that the intrinsic dimension is known and the proposed algorithm uses
 145 a fixed scale, whereas our approach (i) operates in a multi-scale setting, (ii) provides local estimates
 146 of intrinsic dimensionality of the data space, and (iii) incorporates model spaces that serve as a com-
 147 parison. We can thus measure the deviation from an idealised manifold, requiring fewer assumptions
 148 on the structure of the input data (Section 5.4 demonstrates the benefits of this perspective).

149 4 METHODS

150 Our framework TOAST (Topological Algorithm for Singularity Tracking) consists of two parts:
 151 (i) a method to calculate a local intrinsic dimension of the data, and (ii) *Euclidicity*, a measure for
 152 assessing the multi-scale deviation from a Euclidean space. TOAST is based on the assumption that
 153 the intrinsic dimension of some given data is *not* necessarily constant across the data set, and is
 154 best described by *local measurements*, i.e. measurements in a small neighbourhood of a given point.
 155 Since there is no canonical choice for the magnitude of such a neighbourhood, TOAST is built on a
 156 multi-scale analysis of data. **Our main idea involves constructing a collection of local (punctured)**
 157 **neighbourhoods for varying locality scales, and subsequently recording their topological features.**
 158 **This procedure allows us to approximate local topological features (specifically, local homology)**
 159 **of a given point, which we use to measure the intrinsic dimensionality of a space. Moreover, by**
 160 **calculating the distance to Euclidean model spaces, we are capable of detecting singularities in a**
 161 **large range of input data sets.** Subsequently, we will briefly describe the ‘moving parts’ of TOAST;
 162 please refer to Appendix A.1 for a terminology list.

163 4.1 PERSISTENT INTRINSIC DIMENSION

164 For a finite metric space (\mathbb{X}, d) and $x \in \mathbb{X}$, let $B_r^s(x) := \{y \in \mathbb{X} \mid r \leq d(x, y) \leq s\}$ denote
 165 the intrinsic annulus of x in \mathbb{X} with respect to the parameters r and s . Moreover, let \mathcal{F} denote a
 166 procedure that takes as input a finite metric space and outputs an ascending filtration of topological
 167 spaces—such as a Vietoris–Rips filtration. By applying \mathcal{F} to the intrinsic annulus of x , we obtain
 168 a tri-filtration $(\mathcal{F}(B_r^s(x), t))_{r,s,t}$, where t corresponds to the respective filtration step that is deter-
 169 mined by \mathcal{F} . Note that this tri-filtration is covariant in s and t , but contravariant in r ; we denote it by
 170 $\mathcal{F}(B_\bullet^s(x), \bullet)$. Applying i th homology to this filtration yields a tri-parameter persistent module that
 171 we call i th **persistent local homology (PLH)** of x , denoted by $\text{PLH}_i(x; \mathcal{F}) := \text{PH}_i(\mathcal{F}(B_\bullet^s(x), \bullet))$.
 172 To the best of our knowledge, this is the first time that PLH is considered as a multi-parameter per-
 173 sistence module. Since the Vietoris–Rips filtration is the pre-eminent filtration in TDA, we will also
 174 use the abbreviated notation $\text{PLH}_i(x) := \text{PLH}_i(x; \mathcal{V})$.

175 Our PLH formulation enjoys stability properties similar to the seminal stability theorem in persistent
 176 homology (Cohen-Steiner et al., 2007), making it robust to small parameter changes (we assess
 177 empirical stability in Section 5.1).

178 **Theorem 1.** *Given a finite metric space \mathbb{X} and $x \in \mathbb{X}$, let $B_r^s(x)$ and $B_{r'}^{s'}(x)$ denote two*
 179 *intrinsic annuli with $|r - r'| \leq \epsilon_1$ and $|s - s'| \leq \epsilon_2$. Furthermore, let $\mathcal{D}, \mathcal{D}'$ denote the*

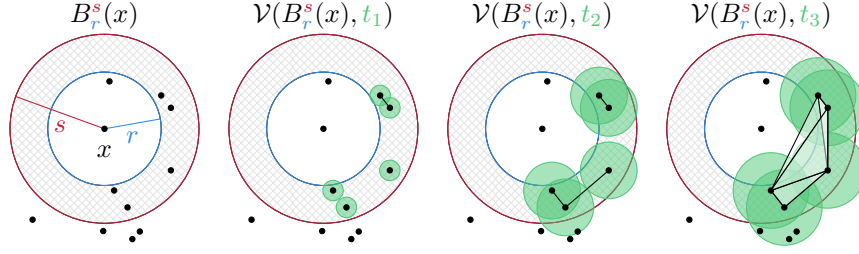


Figure 3: The intrinsic annulus $B_r^s(x)$ around a point x in a metric space (\mathbb{X}, d) , as well as three filtration steps with varying t parameters. By adjusting r and s , we obtain a tri-filtration.

180 persistence diagrams corresponding to $\text{PH}_i(B_r^s(x); \mathcal{V})$ and $\text{PH}_i(B_{r'}^s(x); \mathcal{V})$, respectively. Then
 181 $\frac{1}{2} d_B(\mathcal{D}, \mathcal{D}') \leq \max\{\epsilon_1, \epsilon_2\}$.

182 For a finite set of points $\mathbb{X} \subset \mathbb{R}^N$ and $x \in \mathbb{X}$, we define the **persistent intrinsic dimension (PID)**
 183 of x at scale ϵ as $i_x(\epsilon) := \max\{i \in \mathbb{N} \mid \text{PH}_{i-1}(B_r^s(x)) \neq \emptyset \text{ for some } r \text{ and } s \text{ with } s < \epsilon\}$. This
 184 measure serves as a multi-scale characterisation of the intrinsic dimension of a data set. In case our
 185 data set constitutes a manifold sample, it turns out that we can recover the correct dimension.

186 **Theorem 2.** Let $M \subset \mathbb{R}^N$ be an n -dimensional compact smooth manifold and let $\mathbb{X} :=$
 187 $\{x_1, \dots, x_S\}$ be a collection of uniform samples from M . For a sufficiently large S , there exist
 188 constants $\epsilon_1, \epsilon_2 > 0$ such that $i_x(\epsilon) = n$ for all $\epsilon_1 < \epsilon < \epsilon_2$ and any point $x \in \mathbb{X}$. Moreover, ϵ_1 can
 189 be chosen arbitrarily small by increasing S .

190 The implication of Theorem 2 is that $i_x(\epsilon)$ computes the correct intrinsic dimension of M in a certain
 191 range of values $\epsilon > 0$, provided the sample is sufficiently large. In particular, $i_x(\epsilon)$ persists in this
 192 range, which suggests to consider a collection of $i_x(\epsilon)$ for varying ϵ to analyse the intrinsic dimension
 193 of x . We also have the following corollary, which specifically addresses stratified spaces (such
 194 as the ‘pinched torus’), implying that our method can correctly detect the intrinsic dimension of
 195 individual strata. PID is thus capable of handling large classes of ‘non-manifold’ data sets.

196 **Corollary 1.** Let $X = X_n \supset X_{n-1} \supset X_{n-2} \supset \dots \supset X_1 = \emptyset$ be an n -dimensional
 197 compact stratified simplicial complex, s.t. $X_i \setminus X_{i-1}$ is smooth for every i . For a fixed i , let
 198 $\mathbb{X}_i := \{x_1, \dots, x_S\}$ be a collection of uniform samples from $X_i \setminus X_{i-1}$. For a sufficiently large
 199 S , there are constants $\epsilon_1, \epsilon_2 > 0$ such that $i_x(\epsilon) = i$ for all $\epsilon_1 < \epsilon < \epsilon_2$ and any point $x \in \mathbb{X}_i$.
 200 Moreover, ϵ_1 can be chosen arbitrarily small by increasing S .

201 4.2 EUCLIDICITY

202 Knowledge about the intrinsic dimension of a neighbourhood is crucial for measuring to what extent
 203 such a neighbourhood deviates from being Euclidean. We refer to this deviation as *Euclidicity*, with
 204 the understanding that low values indicate Euclidean neighbourhoods while high values indicate
 205 singular regions of a data set. *Euclidicity* can be calculated without stringent assumptions on mani-
 206 foldness: let $\mathbb{X} \subset \mathbb{R}^N$ be a finite data set, $x \in \mathbb{X}$ a point, and assume that we are given an estimate n
 207 of the intrinsic dimension of x . In particular, the previously-described PID estimation procedure is
 208 applicable in this setting and may be used to obtain n , for example by calculating statistics on the
 209 set of $i_x(\epsilon)$ for varying locality parameters ϵ . Euclidicity, however, can also make use of other di-
 210 mensionality estimation procedures (see Camastra & Staiano (2016) for a survey). To assess how
 211 far a given neighbourhood of x is from being Euclidean, we compare it to a Euclidean model space
 212 by measuring the deviation of their corresponding persistent local homology features. We start by
 213 defining the Euclidean annulus $\text{EB}_r^s(x)$ of x for parameters r and s to be a set of random uniform
 214 samples of $\{y \in \mathbb{R}^n \mid r \leq d(x, y) \leq s\}$ such that $|\text{EB}_r^s(x)| = |B_r^s(x)|$. Here, r and s correspond
 215 to the inner and outer radius of the Euclidean annulus, respectively. For $r' \leq r$ and $s \leq s'$ we extend
 216 $\text{EB}_r^s(x)$ by sampling additional points to obtain $\text{EB}_{r'}^{s'}(x)$ with $|\text{EB}_{r'}^{s'}(x)| = |B_{r'}^{s'}(x)|$. Iterating this
 217 procedure leads to a tri-filtration $(\mathcal{F}(\text{EB}_r^s(x), t))_{r,s,t}$ for any filtration \mathcal{F} , following our description
 218 in Section 4.1. We now define the persistent local homology of a Euclidean model space as

$$\text{PLH}_i^{\text{E}}(x; \mathcal{F}) := \text{PH}_i(\mathcal{F}(\text{EB}_r^s(x), \bullet)). \quad (2)$$

219 Again, for a Vietoris–Rips filtration \mathcal{V} , we use a short-form notation, i.e. $\text{PLH}_i^{\mathbb{E}}(x) := \text{PLH}_i^{\mathbb{E}}(x; \mathcal{V})$.
 220 Notice that $\text{PLH}_i^{\mathbb{E}}(x)$ implicitly depends on the choice of intrinsic dimension n , and on the samples
 221 that are generated randomly. To remove the dependency on the samples, we consider $\text{PLH}_i^{\mathbb{E}}(x)$
 222 to be a sample of a random variable $\mathbf{PLH}_i^{\mathbb{E}}(x)$. Let $D(\cdot, \cdot)$ be a distance measure for 3-parameter
 223 persistence modules, such as the *interleaving distance*.⁴ We then define the **Euclidicity** of x , denoted
 224 by $\mathfrak{E}(x)$, as the expected value of these distances, i.e.

$$\mathfrak{E}(x) := \mathbb{E} \left[D \left(\text{PLH}_{n-1}(x), \mathbf{PLH}_{n-1}^{\mathbb{E}}(x) \right) \right]. \quad (3)$$

225 This quantity essentially assesses how far x is from admitting a regular Euclidean neighbourhood.

226 **Implementation.** Calculating $\mathfrak{E}(x)$ requires different choices, namely (i) a range of locality scales,
 227 (ii) a filtration, and (iii) a distance metric between filtrations D . Using a grid Γ of possible radii (r, s)
 228 with $r < s$, we approximate Eq. (3) using the *mean of the bottleneck distances of fibred Vietoris–Rips*
 229 *barcodes*, i.e.

$$\mathfrak{E}(x) \approx D(\text{PLH}_i(x), \text{PLH}_i^{\mathbb{E}}(x)) := \frac{1}{C} \sum_{(r,s) \in \Gamma} d_{\text{B}}(\text{PH}_i(\mathcal{V}(B_r^s(x), \bullet)), \text{PH}_i(\mathcal{V}(\text{EB}_r^s(x), \bullet))), \quad (4)$$

230 where C is equal to the number of summands and $\text{PLH}_i^{\mathbb{E}}(x)$ refers to a sample from a Euclidean
 231 annulus of the same size as the intrinsic annulus around x . Eq. (4) can be implemented using
 232 effective persistent homology calculation methods (Bauer, 2021), thus permitting an integration into
 233 existing TDA and machine learning frameworks (The GUDHI Project, 2015; Tauzin et al., 2020).
 234 Appendix A.4 provides pseudocode implementations, while Section 5 discusses how to pick these
 235 parameters in practice. We make one specific instantiation of our framework publicly available.⁵

236 **Properties.** The main appeal of our formulation is that calculating both PID and Euclidicity does
 237 not require strong assumptions about the input data. Treating dimension as a local quantity that is
 238 allowed to vary across multiple scales leads to beneficial expressivity properties. As we showed
 239 in Section 4.1, our method is *guaranteed* to yield the right values for manifolds and stratified sim-
 240 plicial complexes. This property substantially increases the practical applicability and expressivity,
 241 enabling our framework to handle unions of manifolds of varying dimensions, for instance. We
 242 require only a basic assumption, namely that the intrinsic dimension n of the given data space is
 243 significantly lower than the ambient dimension N , making Euclidicity broadly applicable. Similar
 244 to curvature, Euclidicity makes use of the fact that one can compare data to ‘model spaces,’ allowing
 245 for different future adjustments.

246 **Limitations.** Our implementation of Euclidicity makes use of the Vietoris–Rips complex, which is
 247 known to grow exponentially with increasing dimensionality. While all calculations of Eq. (3) can be
 248 performed *in parallel*—thus substantially improving scalability vis-à-vis persistent homology on the
 249 complete input data set, both in terms of dimensions and in terms of samples—the memory require-
 250 ments for a full Vietoris–Rips complex construction may still prevent our method to be applicable
 251 for certain high-dimensional data sets. This can be mitigated by selecting a different filtration (Anai
 252 et al., 2020; Sheehy, 2013); our proofs do not assume a specific filtration, and we leave the treatment
 253 of filtration-specific theoretical properties for future work. Finally, we remark that the reliability of
 254 the Euclidicity score depends on the validity of the intrinsic dimension; otherwise, the comparison
 255 does not take place with respect to the appropriate model space.

256 5 EXPERIMENTS

257 We demonstrate the expressivity of our proposed TOAST procedure in different settings, empiri-
 258 cally showing that it (i) calculates the correct intrinsic dimension, and (ii) detects singularities when
 259 analysing data sets with known singular points. We also conduct a comparison with one-parameter
 260 approaches, showcasing how our multi-scale approach results in more stable outcomes. Finally, we
 261 analyse Euclidicity scores of benchmark datasets, giving evidence that our technique can be used as
 262 a measure for the geometric complexity of data.

⁴In our implementation, we will approximate this distance via the bottleneck distance.

⁵See the supplementary materials for the code and experiments.

263 5.1 PARAMETER SELECTION

264 Since Eq. (3) intrinsically incorporates multiple scales of locality, we need to specify an upper bound
 265 for the radii $(r_{\min}, r_{\max}, s_{\min}, s_{\max})$ that define the respective annuli in practice. Given a point x ,
 266 we found the following procedure to be useful in practice: we set s_{\max} , i.e. the maximum of the
 267 outer radius, to the distance to the k th nearest neighbour of a point, and r_{\min} , i.e. the minimum inner
 268 radius, to the smallest non-zero distance to a neighbour of x . Finally, we set the minimum outer
 269 radius s_{\min} and the maximum inner radius r_{\max} to the distance to the $\lfloor \frac{k}{3} \rfloor$ th nearest neighbour.
 270 While we find $k = 50$ to yield sufficient results, spaces with a high intrinsic dimension may require
 271 larger values. The advantage of using such a parameter selection procedure is that it works in a
 272 data-driven manner, accounting for differences in density. Since our approach is inherently *local*,
 273 we need to find a balance between sample sizes that are sufficiently large to contain topological
 274 information, while at the same time being sufficiently small to retain a local perspective. We found
 275 the given range to be an appropriate choice in practice. As for the number of steps, we discretise
 276 the parameter range using 20 steps by default. Higher numbers are advisable when there are large
 277 discrepancies between the radii, for instance when $s_{\max} \gg r_{\max}$.

278 5.2 PERSISTENT INTRINSIC DIMENSION IS EXPRESSIVE

	METHOD	MIN	$\mu \pm \sigma$	MAX
1D	lpca	1.00	1.42±0.78	3.00
	twoNN	0.83	1.00±0.07	1.20
	DANCo	1.00	1.00±0.01	1.16
	PID	1.00	1.12±0.24	1.97
2D	lpca	2.00	2.88±0.32	3.00
	twoNN	1.01	1.90±0.36	2.53
	DANCo	1.00	2.10±0.32	3.00
	PID	1.52	1.95±0.06	2.08

291 Table 1: Dimensionality estimates for
 292 the concatenation of S^1 and S^2 .

293
 294 Euclidicity calculations support *any* dimensionality estimator; since such estimators do not come with
 295 strong guarantees such as Theorem 2, their choice must be ultimately driven by the data set at hand.
 296 See Appendix A.6 for a more detailed analysis of these estimates.

297 **Stability.** In practice, the sample density may not be sufficiently high for Theorem 2 to apply. This
 298 means that there may appear artefact homological features in dimensions *higher* than the intrinsic
 299 dimension of a given space. We thus only consider features that exceed a certain persistence thresh-
 300 old in comparison to the persistence of features of lower dimension: for any data point x and the
 301 respective intrinsic annulus $B_r^s(x)$, we eliminate all topological features whose lifetimes are smaller
 302 than the maximum lifetime of features in one dimension below. This results in markedly stable
 303 estimates of intrinsic dimension, which are less prone to overestimations.

304 5.3 EUCLIDICITY CAPTURES SINGULARITIES

305 Fig. 1 shows that Euclidicity is capable of detecting the singularity of the ‘pinched torus.’ Of partic-
 306 ular relevance is the fact that Euclidicity also highlights that points in the vicinity of the singular
 307 point are *not* fully regular. This is an important property for practical applications since it implies
 308 that Euclidicity can detect such *isolated singularities* even in the presence of sampling errors.

309 Besides the pinched torus, another prototypical example of singular spaces is given by $S^n \vee S^n$, the
 310 wedge of two n -dimensional spheres. Intuitively, $S^n \vee S^n$ is obtained by two n -dimensional spheres
 311 that are glued together at a certain point. Denoting the gluing point by x_0 , for a suitable triangulation
 312 of $X = S^n \vee S^n$, this space is naturally stratified by $X \supset \{x_0\}$. Next, we apply TOAST to samples

⁶Method names are taken from the `scikit-dimension` toolkit. See Appendix A.6 for more details.

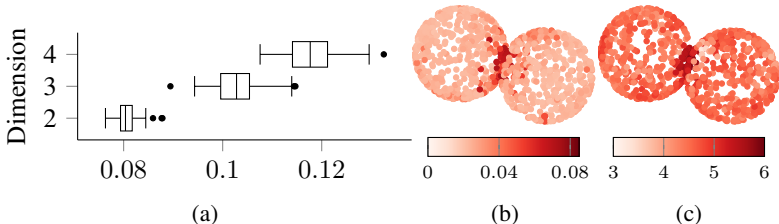


Figure 4: (a): Euclidicity scores of wedged spheres for different dimensions. High values indicate singular points/neighbourhoods. The Euclidicity of the singular point always constitutes a clear positive outlier. In 2D, *Euclidicity* (b) results in a clearly-delineated singular region when compared to a single-parameter score (c).

313 of such wedged spheres of dimensions 2, 3 and 4, calculating their respective Euclidicity scores.
 314 Since larger intrinsic dimensions require higher sample sizes to maintain the same density, we start
 315 with a sample size of 20000 in dimension 2 and increase it consecutively by a factor of 10. We
 316 then calculate Euclidicity of 50 random samples in the respective data set, and additionally for the
 317 singular point x_0 . Fig. 4a shows the results of our experiments. We observe that the singular point
 318 possesses a *significantly higher Euclidicity* score than the random samples. Moreover, we find that
 319 Euclidicity scores of non-singular points exhibit a high degree of variance across the data, which
 320 is caused by the fact that the sampled data does not perfectly fit the underlying space the points
 321 are being sampled from. This strengthens our main argument: assessing whether a specific point is
 322 Euclidean does not require a binary decision but a continuous measure such as Euclidicity.

323 **Stability.** As predicted by Theorem 1, Euclidicity estimates are stable in practice. We first note that
 324 Euclidicity is *robust towards sampling*: repeating the calculations for the ‘pinched torus’ over differ-
 325 ent batches results in highly similar distributions that are not distinguishable according to Tukey’s
 326 range test (Tukey, 1949) at the $\alpha = 0.05$ confidence level. Moreover, choosing larger locality scales
 327 still enables us to detect singularities at higher computational costs and incorporating larger parts of
 328 the point cloud. Please refer to Appendix A.5 for a more detailed discussion of this aspect.

329 5.4 EUCLIDICITY IS MORE EXPRESSIVE THAN SINGLE-PARAMETER APPROACHES

330 Our Euclidicity measure leads to significantly more stable results than a comparable one-parameter
 331 approach for geometry-based anomaly detection (Stolz et al., 2020): Fig. 4b and Fig. 4c compare
 332 multi-parameter Euclidicity with one-parameter Euclidicity for 20000 samples of $S^2 \vee S^2$. The
 333 constant-scale approach results in many points with high anomaly scores that in fact *do* admit a Eu-
 334 clidean neighbourhood. We quantify this by analysing the empirical distributions of anomaly scores
 335 of the two data spaces (see Appendix A.8 for more details), with the one-parameter method ex-
 336 hibiting a much larger variance than our multi-parameter Euclidicity measure. The multi-parameter
 337 distribution shows that the mass is concentrated around the mean, but also contains outliers with
 338 high Euclidicity scores. These outliers correspond to points in the data space whose distance to the
 339 singular point is small. We thus conclude that Euclidicity scores increase once one approaches the
 340 singularity—which is *not* the case for single-parameter methods with a fixed locality scale. In fact,
 341 the main advantage of Euclidicity is that it implicitly incorporates information about the scale on
 342 which a given data point admits a Euclidean neighbourhood.

343 5.5 EUCLIDICITY CAPTURES GEOMETRIC COMPLEXITY OF HIGH-DIMENSIONAL SPACES

344 To test TOAST in an unsupervised setting, we calculate Euclidicity scores for the MNIST
 345 and FASHIONMNIST data sets, selecting mini-batches of 1000 samples from a subsample
 346 of 10000 random images of these data sets. Following Pope et al. (2021), we assume an
 347 intrinsic dimension of 10; moreover, we use $k = 50$ neighbours for local scale estima-
 348 tion. To ensure that our results are representative, we repeat all calculations for five dif-
 349 ferent subsamples. Euclidicity scores range from $[1.1, 5.3]$ for MNIST, and $[1.3, 5.6]$ for
 350 FASHIONMNIST. The scores of the two datasets appear to be following different distribu-
 351 tions (see Appendix A.7 for a visualisation and a more detailed depiction of the distributions).

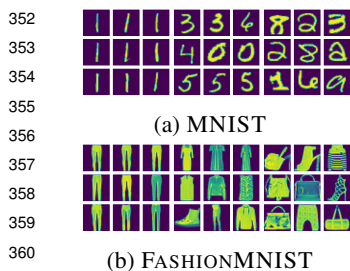


Figure 5: Left to right: low, median, high Euclidity.

Fig. 5 shows a selection of 9 images, corresponding to the lowest, median, and highest Euclidity scores, respectively. We observe that high Euclidity scores correspond to images with a high degree of non-linearity, whereas low Euclidity scores correspond to images that exhibit less complex structures: for MNIST, these are digits of ‘1.’ Interestingly, we observe the same phenomenon for FASHIONMNIST, where images with low Euclidity (‘pants’) possess less geometric complexity in contrast to images with high Euclidity. Since low Euclidity can also be seen as an indicator of how close a neighbourhood is to being *locally linear*, this finding hints at the existence of simple substructures in such data sets. Euclidity could thus be used as an unsupervised measure of geometric complexity.

5.6 EUCLIDITY CAPTURES LOWER-DIMENSIONAL STRUCTURES IN CYTOMETRY DATA

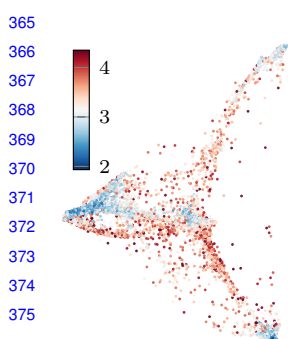


Figure 6: An embedding of the iPSC data with colours based on Euclidity highlights dense non-singular regions.

To highlight the utility of Euclidity in unsupervised representation learning, we calculate it on an induced pluripotent stem cell (iPSC) reprogramming data set (Zunder et al., 2015). The data set depicts a progression of so-called fibroblasts diverging, and splitting into two different lineages. Fig. 6 shows an embedding obtained via PHATE (Moon et al., 2019) and the Euclidity scores of the original data. We find that high Euclidity scores correspond to points that exhibit a *lower density* in the embedding, being in fact situated in lower-dimensional subspaces. Since lower-dimensional points in a space can be considered *singular* in the sense of stratified spaces, this is further evidence for Euclidity to be a useful tool for detecting non-manifold regions in data. Please refer to Appendix A.9 for more details.

6 DISCUSSION

We presented TOAST, a novel framework for locally estimating the intrinsic dimension (via PID, the persistent intrinsic dimension) and the ‘manifoldness’ (via Euclidity, a multi-scale measure of the deviation from Euclidean space) of point clouds. Our method is based on a novel formulation of persistent local homology as a multi-parameter approach, and we provide theoretical guarantees for it in a dense sample setting. Our experiments showed significant improvements of stability compared to geometry-based anomaly detection methods with fixed locality scales, and we found that Euclidity can detect singular regions in data sets with known singularities. Using high-dimensional benchmark data sets, we also observed that Euclidity can serve as an *unsupervised measure of geometric complexity*.

For future work, we envision two relevant research directions. First and foremost will be the inclusion of Euclidity into machine learning models to make them ‘singularity-aware.’ In light of our experiments in Section 5.5, we believe that Euclidity could be particularly useful in unsupervised scenarios, or provide an additional weight in classification settings (to ensure that singular examples are being given lower confidence scores). Moreover, Euclidity could be used in the detection of adversarial samples—a task for which knowledge about the underlying topology of a space is known to be crucial (Jang et al., 2020). As a second direction, we want to further improve the properties of Euclidity itself. To this end, we plan to investigate if incorporating custom distance measures for three-parameter persistence modules, i.e. different metrics for Eq. (4), lead to improved results in terms of stability, expressivity, or computational efficiency. Moreover, we hypothesise that replacing the Vietoris–Rips filtration by other constructions (de Silva & Carlsson, 2004) could prove beneficial in reducing the number of samples for calculating Euclidity. Along these lines, we also plan to derive theoretical results that relate specific filtrations and the expressivity of the corresponding Euclidity measure. Another direction for future research concerns the approximation of a manifold from inherently singular data, i.e. finding the *best* manifold approximation to a given data set with singularities. This way, singularities could be resolved during *the training phase of models*, provided an appropriate loss function exists. Euclidity may thus serve as a metric for assessing data sets, paving the way towards more trustworthy and faithful embeddings.

405 REPRODUCIBILITY STATEMENT

406 We provide our code as part of the supplementary materials. All dependencies are listed in the re-
 407 spective `pyproject.toml` file, and the `README` discusses how to install our package and run our
 408 experiments. Our implementation leverages multiple CPUs if available but has no specific hardware
 409 requirements otherwise.

410 REFERENCES

- 411 Hirokazu Anai, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hiroya Inakoshi, Raphaël Tinarrage, and
 412 Yuhei Umeda. DTM-based filtrations. In N. A. Baas, G. E. Carlsson, G. Quick, M. Szymik,
 413 and M. Thaulé (eds.), *Topological Data Analysis*, pp. 33–66. Springer, 2020. doi: 10.1007/
 414 978-3-030-43408-3_2.
- 415 Ulrich Bauer. Ripser: efficient computation of Vietoris–Rips persistence barcodes. *Journal of*
 416 *Applied and Computational Topology*, 5(3):391–423, 2021. doi: 10.1007/s41468-021-00071-5.
- 417 Paul Bendich. *Analyzing stratified spaces using persistent versions of intersection and local homol-*
 418 *ogy*. PhD thesis, Duke University, 2008.
- 419 Paul Bendich, David Cohen-Steiner, Herbert Edelsbrunner, John Harer, and Dmitriy Morozov. In-
 420 ferring local homology from sampled stratified spaces. In *IEEE Symposium on Foundations of*
 421 *Computer Science (FOCS)*, pp. 536–546, 2007. doi: 10.1109/FOCS.2007.45.
- 422 Tolga Birdal, Aaron Lou, Leonidas J. Guibas, and Umut Simsekli. Intrinsic dimension, persistent
 423 homology and generalization in neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin,
 424 P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*,
 425 volume 34, pp. 6776–6789, 2021.
- 426 Bradley C.A. Brown, Anthony L. Caterini, Brendan Leigh Ross, Jesse C. Cresswell, and Gabriel
 427 Loaiza-Ganem. The union of manifolds hypothesis and its implications for deep generative mod-
 428 elling. *arXiv:2207.02862*, 2022.
- 429 Francesco Camastra and Antonino Staiano. Intrinsic dimension estimation: Advances and open
 430 problems. *Information Sciences*, 328:26–41, 2016. doi: 10.1016/j.ins.2015.08.029.
- 431 Gunnar Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–
 432 308, 2009.
- 433 Frédéric Chazal, Vin De Silva, and Steve Oudot. Persistence stability for geometric complexes.
 434 *Geometriae Dedicata*, 173(1):193–214, 2014. doi: 10.1007/s10711-013-9937-z.
- 435 David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams.
 436 *Discrete & Computational Geometry*, 37(1):103–120, 2007. doi: 10.1007/s00454-006-1276-5.
- 437 Vin de Silva and Gunnar Carlsson. Topological estimation using witness complexes. In M. Gross,
 438 H. Pfister, M. Alexa, and S. Rusinkiewicz (eds.), *Symposium on Point-Based Graphics*. The Eu-
 439 rographics Association, 2004. doi: 10.2312/SPBG/SPBG04/157-166.
- 440 Tamal K. Dey, Fengtao Fan, and Yusu Wang. Dimension detection with local homology. In *Can-*
 441 *adian Conference on Computational Geometry (CCCG)*, 2014.
- 442 Herbert Edelsbrunner and John Harer. *Computational topology: An introduction*. American Math-
 443 ematical Society, Providence, RI, USA, 2010.
- 444 Brittany Terese Fasy and Bei Wang. Exploring persistent local homology in topological data analy-
 445 sis. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp.
 446 6430–6434, 2016.
- 447 Felix Hensel, Michael Moor, and Bastian Rieck. A survey of topological machine learning methods.
 448 *Frontiers in Artificial Intelligence*, 4, 2021. doi: 10.3389/frai.2021.681108.

- 449 Alexander Jakubowski, Milica Gašić, and Marcus Zibrowius. Topology of word embeddings: Sin-
450 gularities reflect polysemy. In *Proceedings of the Ninth Joint Conference on Lexical and Compu-*
451 *tational Semantics*, pp. 103–113. Association for Computational Linguistics, 2020.
- 452 Uyeong Jang, Susmit Jha, and Somesh Jha. On the need for topology-aware generative models for
453 manifold-based defenses. In *International Conference on Learning Representations*, 2020. URL
454 https://openreview.net/forum?id=r11F_CeYwS.
- 455 Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations*
456 *and Trends® in Machine Learning*, 12(4):307–392, 2019. doi: 10.1561/22000000056.
- 457 Leland McInnes, John Healy, and James Melville. UMAP: Uniform manifold approximation and
458 projection for dimension reduction. *arXiv:1802.03426*, 2018.
- 459 Kevin R. Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B. Burkhardt, William S. Chen,
460 Kristina Yim, Antonia van den Elzen, Matthew J. Hirn, Ronald R. Coifman, Natalia B. Ivanova,
461 Guy Wolf, and Smita Krishnaswamy. Visualizing structure and transitions in high-dimensional
462 biological data. *Nature Biotechnology*, 37(12):1482–1492, 2019.
- 463 Michael Moor, Max Horn, Bastian Rieck, and Karsten Borgwardt. Topological autoencoders. In
464 H. Daumé III and A. Singh (eds.), *Proceedings of the 37th International Conference on Machine*
465 *Learning (ICML)*, number 119 in Proceedings of Machine Learning Research, pp. 7045–7054.
466 PMLR, 2020.
- 467 Phillip Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic
468 dimension of images and its impact on learning. In *International Conference on Learning*
469 *Representations*, 2021.
- 470 Bastian Rieck and Heike Leitte. Persistent homology for the evaluation of dimensionality reduction
471 schemes. *Computer Graphics Forum*, 34(3):431–440, June 2015. doi: 10.1111/cgf.12655.
- 472 Bastian Rieck, Markus Banagl, Filip Sadlo, and Heike Leitte. Persistent intersection homology for
473 the analysis of discrete data. In H. Carr, I. Fujishiro, F. Sadlo, and S. Takahashi (eds.), *Topological*
474 *Methods in Data Analysis and Visualization V*, pp. 37–51. Springer, Cham, Switzerland, 2020.
475 doi: 10.1007/978-3-030-43036-8_3.
- 476 Luis Scoccola and Jose A. Perea. Fiberwise dimensionality reduction of topologically complex data
477 with vector bundles. *arXiv:2206.06513*, 2022.
- 478 Donald R. Sheehy. Linear-size approximations to the vietoris–rips filtration. *Discrete & Computa-*
479 *tional Geometry*, 49(4):778–796, 2013. doi: 10.1007/s00454-013-9513-1.
- 480 Bernadette J. Stolz, Jared Tanner, Heather A. Harrington, and Vidit Nanda. Geometric anomaly
481 detection in data. *Proceedings of the National Academy of Sciences*, 117(33):19664–19669, 2020.
482 doi: 10.1073/pnas.2001741117.
- 483 Guillaume Tuzin, Umberto Lupo, Lewis Tunstall, Julian Burella Pérez, Matteo Caorsi, Anibal
484 Medina-Mardones, Alberto Dassatti, and Kathryn Hess. *giotto-tda*: A topological data anal-
485 ysis toolkit for machine learning and data exploration. *arXiv:2004.02551*, 2020.
- 486 The GUDHI Project. *GUDHI User and Reference Manual*. GUDHI Editorial Board, 2015. URL
487 <http://gudhi.gforge.inria.fr/doc/latest/>.
- 488 John W. Tukey. Comparing individual means in the analysis of variance. *Biometrics*, 5(2):99–114,
489 1949.
- 490 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine*
491 *Learning Research*, 9(86):2579–2605, 2008.
- 492 Eli R. Zunder, Ernesto Lujan, Yury Goltsev, Marius Wernig, and Garry P. Nolan. A continuous
493 molecular roadmap to iPSC reprogramming through progression analysis of single-cell mass cy-
494 tometry. *Cell Stem Cell*, 16(3):323–337, 2015.

495 A APPENDIX

496 A.1 NOTATION

Symbol	Meaning
ϵ	local annulus scale parameter
\mathbb{R}	real numbers
H_i	i th (ordinary) homology functor (with $\mathbb{Z}/2\mathbb{Z}$ coefficients)
\tilde{H}_i	i th reduced homology functor (with $\mathbb{Z}/2\mathbb{Z}$ coefficients)
\inf	infimum
497 \sup	supremum
$ \cdot _\infty$	uniform (infinity) norm
n	intrinsic dimension of the space under consideration
N	ambient dimension of the space under consideration
\lim	(categorical) colimit
S^k	k -dimensional sphere
$c^\circ X := X \times (0, 1]/X \times \{1\}$	open cone of a topological space X

498 A.2 PROOFS OF THE MAIN STATEMENTS IN THE PAPER

499 We restate the theorems from the main paper for the convenience of readers, along with their proofs,
500 which were removed for space reasons. We first prove the stability theorem, first stated on p. 5 in the
501 main text, which shows that our method enjoys stability properties with respect to radius changes of
502 the intrinsic annuli.

503 **Theorem 1.** Given a finite metric space \mathbb{X} and $x \in \mathbb{X}$, let $B_r^s(x)$ and $B_{r'}^{s'}(x)$ denote two
504 intrinsic annuli with $|r - r'| \leq \epsilon_1$ and $|s - s'| \leq \epsilon_2$. Furthermore, let $\mathcal{D}, \mathcal{D}'$ denote the
505 persistence diagrams corresponding to $\text{PH}_i(B_r^s(x); \mathcal{V})$ and $\text{PH}_i(B_{r'}^{s'}(x); \mathcal{V})$, respectively. Then
506 $\frac{1}{2} d_B(\mathcal{D}, \mathcal{D}') \leq \max\{\epsilon_1, \epsilon_2\}$.

507 *Proof.* The Hausdorff distance of two non-empty subsets $A, B \subset \mathbb{X}$ is $d_H(A, B) := \inf\{\epsilon \geq$
508 $0 \mid A \subset B_\epsilon, B \subset A_\epsilon\}$, where $A_\epsilon = \cup_{a \in A} \{x \in \mathbb{X}; d(x, a) \leq \epsilon\}$ denotes the ϵ -thickening of
509 A in X . Set $\epsilon := \max\{\epsilon_1, \epsilon_2\}$. By assumption, $B_r^s(x) \subset B_{r'}^{s'}(x)_\epsilon$ and $B_{r'}^{s'}(x) \subset B_r^s(x)_\epsilon$, i.e.
510 $d_H(B_r^s(x), B_{r'}^{s'}(x)) \leq \epsilon$. Using the geometric stability theorem of persistence diagrams (Chazal
511 et al., 2014), we have $\frac{1}{2} d_B(\mathcal{D}, \mathcal{D}') \leq d_H(B_r^s(x), B_{r'}^{s'}(x))$, which proves the claim. \square

512 Next, we prove that our *persistent intrinsic dimension (PID)* measure is capable of capturing the
513 dimension of manifolds correctly, provided sufficiently many samples are present. This theorem
514 was first stated on p. 5.

515 **Theorem 2.** Let $M \subset \mathbb{R}^N$ be an n -dimensional compact smooth manifold and let $\mathbb{X} :=$
516 $\{x_1, \dots, x_S\}$ be a collection of uniform samples from M . For a sufficiently large S , there exist
517 constants $\epsilon_1, \epsilon_2 > 0$ such that $i_x(\epsilon) = n$ for all $\epsilon_1 < \epsilon < \epsilon_2$ and any point $x \in \mathbb{X}$. Moreover, ϵ_1 can
518 be chosen arbitrarily small by increasing S .

519 *Proof.* Let $x \in \mathbb{X}$ be an arbitrary point. Since M is a manifold, x admits a Euclidean neighbour-
520 hood U . Since M is smooth, we can assume U to be arbitrarily close to being flat by shrinking it.
521 Thus, we can find $\epsilon_2 > 0$ with $B_r^s(x) \subset U$ for all $r, s < \epsilon_2$ such that $H_i(\mathcal{V}(B_r^s(x), t)) = 0$ for all
522 $i \geq n$, and all t . Hence, $\text{PH}_i(B_r^s(x)) = 0$ for all $i \geq n$, and therefore $i_x(\epsilon_2) \leq n$. By contrast, for
523 S sufficiently large, and r, s as before, there exists a parameter t such that $\mathcal{V}(B_r^s(x), t)$ is homotopy-
524 equivalent to an $(n - 1)$ -sphere, and so $H_{n-1}(\mathcal{V}(B_r^s(x), t))$ admits a generator, i.e. it is non-trivial.
525 Consequently, $\text{PH}_{n-1}(B_r^s(x)) \neq 0$, and $i_x(\epsilon_2) = n$. By further increasing S , we can ensure that
526 the statement still holds when we decrease ϵ_2 , which proves the two remaining claims. \square

527 A.3 ADDITIONAL PROOFS

528 To make this paper self-contained, we provide a brief proof of Eq. (1). By the excision axiom for
529 homology, we have

$$H_i(X, X \setminus x) \cong H_i(\text{St}(x), \text{St}(x) \setminus x). \quad (5)$$

530 Since $\text{St}(x)$ is *contractible*, the long exact reduced homology sequence of the pair $(\text{St}(x), \text{St}(x) \setminus x)$
531 records exactness of

$$0 = \tilde{H}_i(\text{St}(x)) \rightarrow H_i(\text{St}(x), \text{St}(x) \setminus x) \rightarrow \tilde{H}_{i-1}(\text{St}(x) \setminus x) \rightarrow \tilde{H}_{i-1}(\text{St}(x)) = 0$$

532 for all i , and therefore $H_i(\text{St}(x), \text{St}(x) \setminus x) \cong \tilde{H}_{i-1}(\text{St}(x) \setminus x)$. Eq. (1) now follows from the
533 observation that $\text{St}(x) \setminus x$ deformation retracts to $\text{Lk}(x)$.

534 A.4 PSEUDOCODE

535 We provide brief pseudocode implementations of the algorithms discussed in Section 4. In the fol-
536 lowing, we use $\# \text{Bar}_i(\mathbb{X})$ to denote the number of i -dimensional persistent barcodes of \mathbb{X} (w.r.t.
537 the Vietoris–Rips filtration, but any other choice of filtration affords the same description). Algo-
538 rithm 1 explains the calculation of *persistent intrinsic dimension* (see Section 4.1 in the main paper
539 for details). For the subsequent algorithms, we assume that the estimated dimension of the intrinsic
540 dimension of the data is n . We impose no additional requirements on this number; it can, in fact,
541 be obtained by any choice of intrinsic dimension estimation method. As a short-hand notation, for
542 $p_i = \text{PH}_{n-1}(\mathcal{V}(\text{EB}_r^s(x), \bullet))$ w.r.t. some sample of $\{y \in \mathbb{R}^n \mid r \leq d(x, y) \leq s\}$, we denote by
543 $p_i^{r,s} = \text{PH}_{n-1}(\mathcal{V}(\text{EB}_r^s(x), \bullet))$ the respective fibred persistent local homology barcode (calculated
544 w.r.t. the same sample). Algorithm 2 then shows how to calculate the *Euclidicity* values, following
545 Eq. (3) and one of its potential implementations, given in Eq. (4).

Algorithm 1 An algorithm for calculating the *persistent intrinsic dimension* (PID)

Require: $x \in \mathbb{X}$, s_{\max} , ℓ .

```

1: for  $s \in \Gamma$  do                                     ▷ Iterate over the parameter grid
2:    $i_x(s) \leftarrow 0$ 
3:   for  $r < s \in \Gamma$  do
4:     for  $i = 1, \dots, N - 1$  do
5:       Calculate  $\# \text{Bar}_i(B_r^s(x))$ 
6:       if  $\# \text{Bar}_i(B_r^s(x)) > 0$  then
7:          $i_x(s) \leftarrow i + 1$ 
8:       end if
9:     end for
10:  end for
11:  return  $i_x(s)$ 
12: end for

```

Algorithm 2 An algorithm for calculating the *Euclidicity values* δ_{jk}

Require: $x \in \mathbb{X}$, s_{\max} , ℓ , n , $\{p_1, \dots, p_m\}$.

```

1: for  $j = 1, \dots, m$  do
2:   for  $k = j + 1, \dots, m$  do
3:     for  $s \in \Gamma$  do
4:       for  $r \in \Gamma, r < s$  do
5:         Calculate  $d_B(p_j^{r,s}, p_k^{r,s})$                                      ▷ Calculate bottleneck distance
6:         return  $d_B(p_j^{r,s}, p_k^{r,s})$ 
7:       end for
8:     end for
9:     Calculate  $D(p_j, p_k)$                                              ▷ Evaluate Eq. (4)
10:    return  $D(p_j, p_k)$ 
11:  end for
12: end for

```

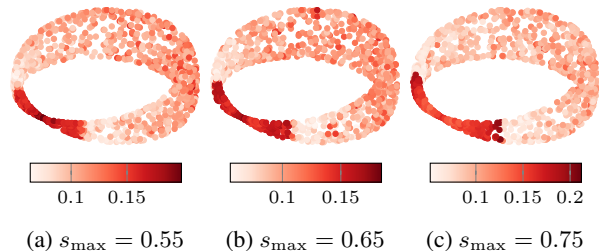


Figure 8: Modifying the outer radius s_{\max} still enables us to detect the singularity of the ‘pinched torus.’ Larger radii, however, progressively increase the field of influence of our method, thus starting to assign high Euclidicity values to larger regions of the point cloud.

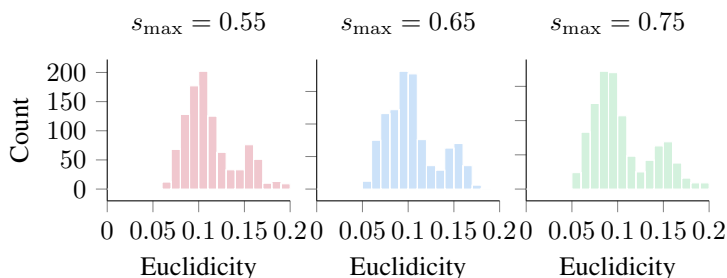


Figure 9: Histograms of the Euclidicity values for the point clouds shown in Fig. 8. Larger radii result in the distribution accumulating more probability mass at higher Euclidicity values, making the singularity detection procedure less local (but still succeeding in detecting the singularity and its environs).

546 A.5 STABILITY OF EUCLIDICITY ESTIMATES

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

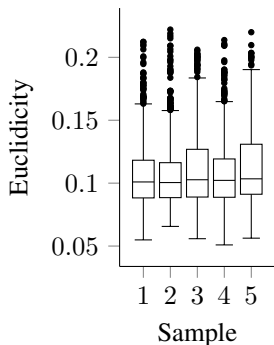


Figure 7: Boxplots of the Euclidicity values of different random samples of the ‘pinched torus’ data set. While each sample invariably exhibits some degree of geometric variation, we are able to reliably identify the singularity and its neighbourhood.

that we lose a clear distinction between singular values and non-singular values. Our data-driven parameter selection procedure is thus to be preferred in practice since it incorporates data density.

Fig. 7 shows that Euclidicity is robust under sampling; repeating the calculations for smaller batches of the ‘pinched torus’ data set (500 points each) still lets us detect the singularity and its neighbours reliably. This robustness is an important property in practice where we are dealing with samples from an unknown data set whose shape properties we want to capture. Euclidicity enables us to perform these calculations in a robust manner. Following the brief discussion in Section 5.1, we show the results of varying s_{\max} , the outer radius of the local annulus, for the ‘pinched torus’ data set. Fig. 8 depicts point clouds of 1000 samples; we observe that the singularity, i.e. the ‘pinch point,’ is always detected. For larger radii, however, this detection becomes progressively more *global*, incorporating larger parts of the point cloud. Fig. 9 depicts the corresponding histograms; we observe the same shift in probability mass towards the tail end of the distribution. For extremely large annuli, we estimate

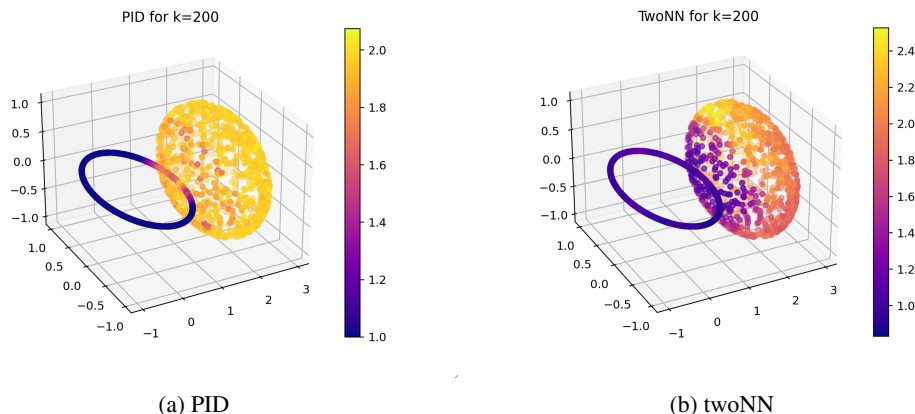


Figure 10: Even for large values of k , PID still does not overestimate the local dimensionality of the data, exhibiting a clear distinction between the circle and the sphere, respectively.

570 A.6 COMPARISON OF PID WITH OTHER DIMENSION ESTIMATES

571 In order to assess the quality of PID, we decided to test its performance on a space that is both singular
 572 and has non-constant dimension. The data space we chose consists of 2000 samples of $S^1 \vee S^2$,
 573 i.e. a 1-sphere glued together with a 2-sphere at a certain concatenation point. We then applied the
 574 PID procedure for a maximum locality scale that was given by the k nearest neighbour distances, for
 575 $k \in \{25, 50, 75, 100, 125, 150, 175, 200\}$. We assigned to each point the average of the PID scores
 576 at the respective scales that are less than or equal to the k nearest neighbour bound. Subsequently,
 577 we compared the results with other local dimension estimates for the respective number of neigh-
 578 bours. The methods that were chosen for comparison include `lpca`, `twoNN`, `KNN`, and `DANC`; we
 579 used the respective implementation from the `scikit-dimension` Python package.⁷

580 Fig. 10a shows the PID results for a maximum locality scale of 200 neighbours, with colours show-
 581 ing the estimated dimension values for each point. Overall, the correct intrinsic dimension is de-
 582 tected for most of the points. However, points that lie close to the singular point show a PID value
 583 between 1 and 2. Similarly to what we already discussed for Euclidicity, PID should therefore also
 584 be interpreted as a measure that incorporates the intrinsic dimension of a point on *multiple scales*
 585 of locality. For real-world data, the dimension will generally change when changing the locality
 586 scale. However, since there is no canonical choice of scale, we believe that any such scale provides
 587 valuable information about the intrinsic dimension that is worth being measured. We therefore argue
 588 that a multi-scale approach like ours is appropriate in practice, especially in a regime that is agnostic
 589 with respect to the underlying intrinsic dimension. By contrast, Fig. 10b shows the corresponding
 590 dimension estimates for `twoNN`, where we observe less stable and reliable results across the dataset.

591 Fig. 11a shows boxplots of the distributions of the dimension estimates, for all points that lie on
 592 the 1D-sphere. We see that for PID, the mass is concentrated at a value of 1. Although there are
 593 outliers present, these correspond to points that are close to the singularity, as it was expected. We
 594 note that other methods like `KNN` and `lpca` might highly overestimate the dimension, and that the
 595 interquartile range is significantly higher for `twoNN` and `KNN`. Fig. 11b shows the same distributions
 596 for the points that lie on the 2D-sphere. Again, `lpca` highly overestimates the dimension since the
 597 median lies at a value of 3. Again, the interquartile range of PID is the tightest, and the estimates
 598 are closest to the ground truth. Moreover, the lower-value outliers again correspond to points that
 599 are close to the singular gluing point.

600 Fig. 12a and Fig. 12b show average dimension estimate scores of all investigated methods for vary-
 601 ing values of k , both for points on the 1-sphere and the 2-sphere. We note that on average, only
 602 `twoNN` and `DANC` lead to results which are comparable with the reliability of our method. How-
 603 ever, as we already saw in Fig. 11a and Fig. 11b, the variance of the scores of our method is signifi-
 604 cantly lower, leading to more reliable outputs for each of the points.

⁷<https://scikit-dimension.readthedocs.io/en/latest/>

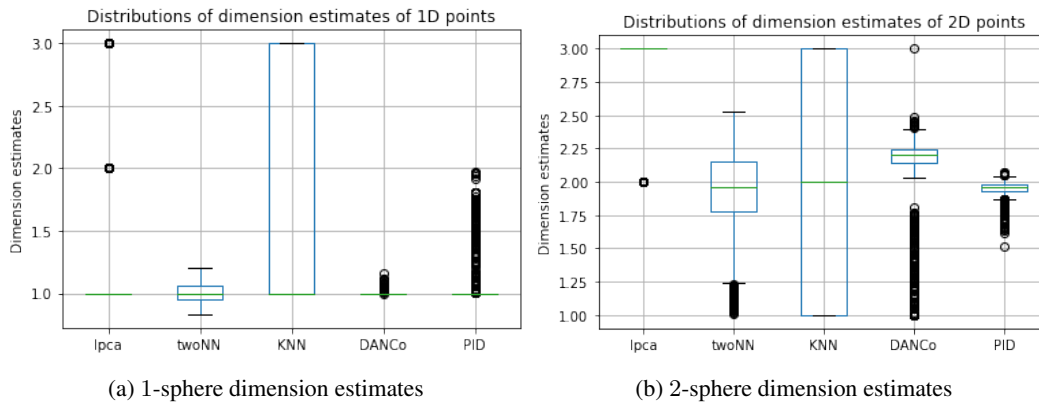


Figure 11: Estimates of the local intrinsic dimension for points that are close to the 1D-sphere, i.e. the circle, or the 2D-sphere, respectively.

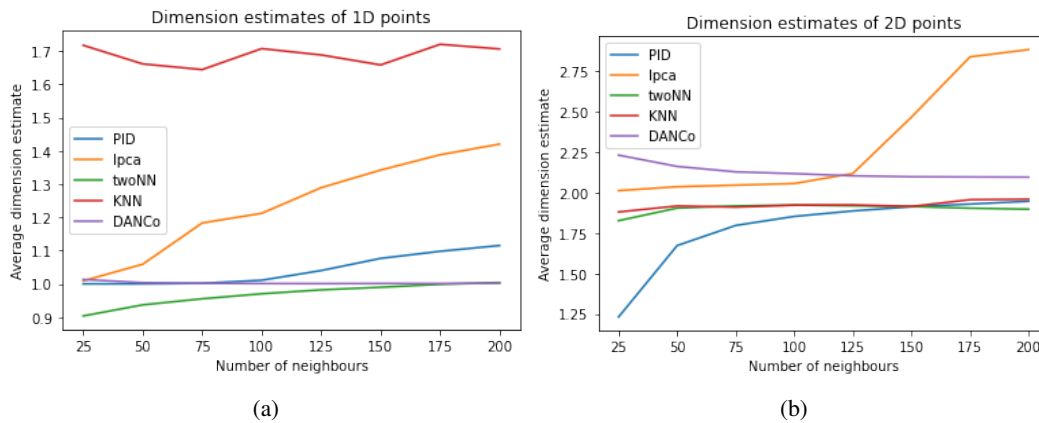


Figure 12: Dimension estimates of the 1D-sphere and the 2D-sphere for different methods, plotted as a function of the number of neighbours k .



Figure 13: From left to right: more examples of low Euclidicity values, median Euclidicity values, and high Euclidicity values for the MNIST data set.

605 A.7 EUCLIDICITY OF MNIST AND FASHIONMNIST

606 Fig. 13 and Fig. 14 show the Euclidicity results for the 4 additional runs on both the MNIST and
 607 FASHIONMNIST data sets. Again, we depicted the 9 images with lowest (left), medium (middle),
 608 and highest (right) Euclidicity scores for the two datasets. Moving from left to right, the images
 609 exhibit increases in the complexity of the local geometry, giving evidence for the reproducibility of
 610 the observation we remarked in Section 5.5.

611 Finally, as Fig. 15 shows, the empirical distributions of the calculated Euclidicity scores differ signif-
 612 icantly for the MNIST and FASHIONMNIST data sets, with the distribution for MNIST exhibiting
 613 a bimodal behaviour, whereas the FASHIONMNIST Euclidicity value distribution is unimodal. We
 614 hypothesise that this corresponds to regions of simple complexity—and locally linear structures—in
 615 the MNIST data set, which are absent in the FASHIONMNIST data set.

616 A.8 ONE-PARAMETER VERSUS MULTI-PARAMETER EUCLIDICITY FOR WEDGED SPHERES

617 Fig. 16 shows the empirical distributions of Euclidicity scores for fixed locality parameters (left) and
 618 for our proposed multi-scale locality approach (right). We see that the variance is *significantly lower*
 619 in the multi-scale regime, indicating more stable and robust results. Moreover, the ratio of maximum
 620 and mean is higher in the multi-parameter setting, where high Euclidicity scores correspond to data
 621 points that lie close to the singularity, resulting in more reliable outcomes.

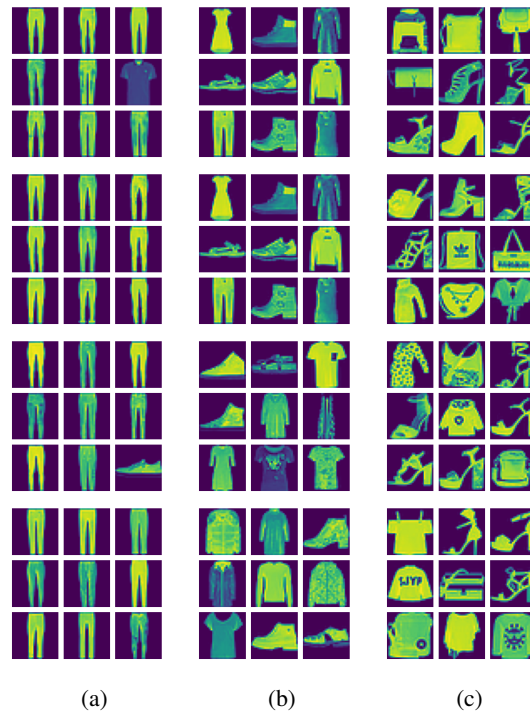


Figure 14: From left to right: more examples of low Euclidity values, median Euclidity values, and high Euclidity values for the FASHIONMNIST data set.

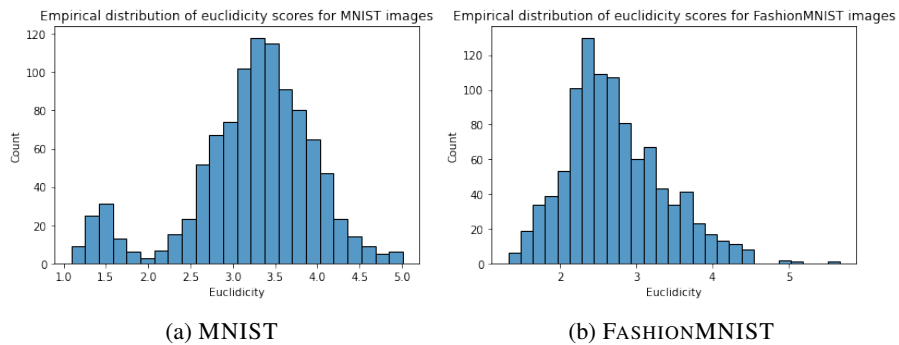


Figure 15: Both MNIST and FASHIONMNIST exhibit markedly different distributions in terms of Euclidity: MNIST Euclidity values are bimodal, whereas FASHIONMNIST Euclidity values are unimodal.

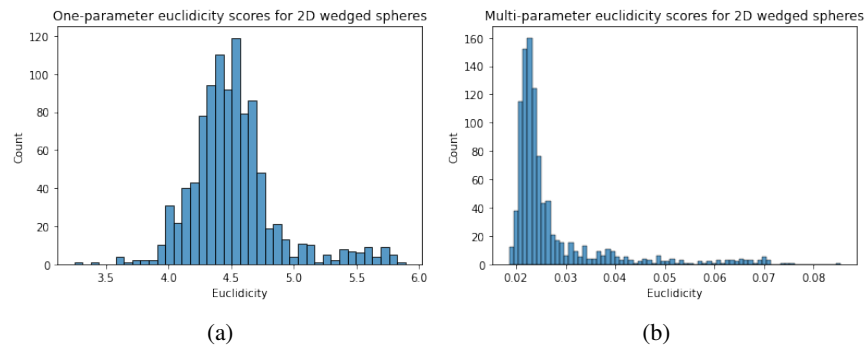


Figure 16: A comparison of Euclidity values of a one-parameter approach (left) and our multi-parameter approach (right) demonstrates that multiple scales are necessary to adequately capture singularities.

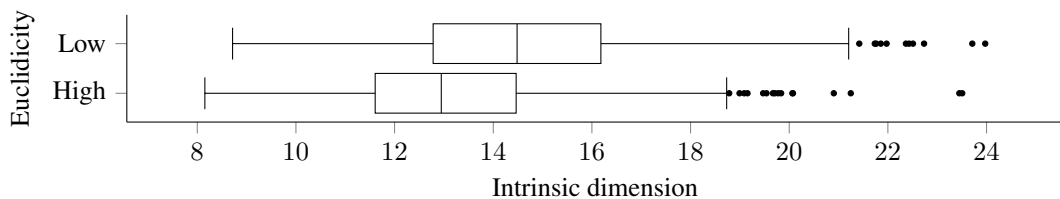


Figure 17: A comparison of intrinsic dimension estimates computed for points in the iPSC dataset that admit high (left) and low (right) Euclidity scores. The $t_{\text{wO}}\text{NN}$ dimensionality estimator was used for this example.

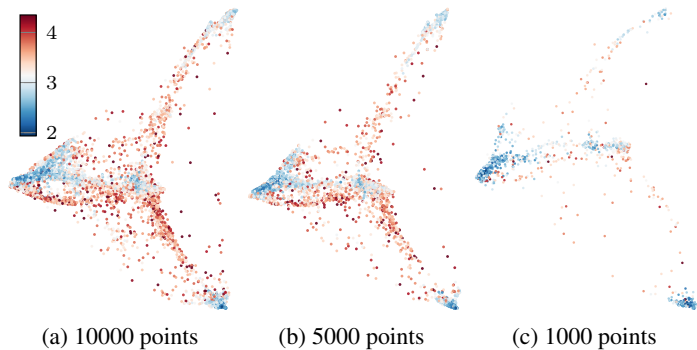


Figure 18: Euclidity remains stable under subsampling the iPSC data set. Minor variations in the point cloud shape are due to the PHATE embedding algorithm; Euclidity was calculated on the raw data.

622 A.9 EUCLIDICITY OF IPSC DATA

623 The iPSC data set Zunder et al. (2015) consists of 33 variables and around 220k samples. It is
 624 known to contain branching structures that can best be extracted using PHATE (Moon et al., 2019),
 625 a non-linear dimensionality reduction algorithm. We only employ this algorithm for *visualisation*
 626 *purposes*; all Euclidity calculations are performed on the original data. Using $t_{\text{wO}}\text{NN}$ for dimension-
 627 *ality* estimation, we obtained a mean intrinsic dimension of 16; as outlined above, other dimension-
 628 *ality* estimators may be employed as well—we consider this analysis to be a proof of concept
 629 first and foremost. We selected parameters as described in Section 5.5, and computed Euclidity
 630 for 10000 samples.

631 We observe that high-Euclidity scores correspond to points that exhibit a *lower density* in the
 632 PHATE embedding,⁸ and according to the $t_{\text{wO}}\text{NN}$ estimates we see that such points are in fact of
 633 *lower intrinsic dimension*; see Fig. 17 for details. More specifically, we calculated the intrinsic
 634 dimension for the subsample, observing that the interquartile range for the 1000 points with *highest*
 635 *Euclidity* is around 12–14, whereas the interquartile range of the 1000 *lowest Euclidity* points
 636 ranges between around 13–16. Again, we used the $t_{\text{wO}}\text{NN}$ algorithm for intrinsic dimensionality
 637 estimates (using $k = 50$ nearest neighbours). Since lower-dimensional points in a space can be
 638 regarded as being singular in the sense of stratified spaces, we see further evidence for Euclidity
 639 as a useful tool for the detection of non-manifold regions in the data. Finally, we remark that
 640 our analyses remain valid under *subsampling*. Fig. 18 depicts subsamples of different sizes for
 641 which we calculated Euclidity (on the raw data, respectively, using PHATE to obtain embeddings).
 642 Euclidity distributions remain stable and the same phenomena are highlighted for each subsample.

⁸However, notice that low-density regions in the PHATE visualisation need not necessarily correspond to low-density regions in the original dataset.