

Enabling Multi-Scale Simplicial Complex Generation for Mapper

Matt Piekenbrock* and Derek Doran†

Abstract. Mapper is perhaps the most widely used algorithm for practical topological data analysis. The primary result of the algorithm is a graph depicting topological relationships and structures usable for comparing, visualizing, and analyzing high dimensional data. As an unsupervised, exploratory tool that may be used for multimodal and heterogeneous forms data, Mapper naturally relies on a number of parameters which explicitly control the quality of the resulting construction; one such parameter controls the entire relational component of the output graph, and offers little guidance or intuition on what values may provide ‘better’ results. In this effort, we provide a closed-form solution that allows the efficient computation of all possible realizations of this crucial parameter, and we discuss further extensions and theoretical advancements that may be made as a result of this development.

Key words. Mapper, Topology Theory, Unsupervised Learning

AMS subject classifications. 54H99, 54C05, 54J05

Introduction. The notion of *nearness* plays an important role in studying and understanding relationships in data. Within the realm of unsupervised machine learning, there is now a well-developed theory of how measures of nearness may be exploited for the purpose of modeling, understanding, or visualizing data, even if the data is sufficiently complex or higher dimensional. For example, in manifold learning, distance interpolation methods along manifold surfaces enable realistic, perceptual changes between images to be more semantically clear [14]. In clustering, distance measures are used to express the (dis)similarity between objects towards creating a meaningful representation suitable for classification, based on characteristics of the data itself [8]. Self-organizing maps (SOMs), often referred to in the context of artificial neural networks for unsupervised learning, have been used extensively for mapping data to lower-dimensional representations which preserve some degree of topological properties of the input space, allowing for robust feature representation and visualization [9]. Indeed, it is not difficult to find unsupervised models which exploit notions of nearness to describe something useful or interesting about the space the data, or functions of the data, were observed in.

Although there are many approaches in unsupervised learning, perhaps the most foundational field that unifies these various concepts of “nearness” is the field of *topology theory*. Topology theory explores deformations of maps and spaces, and has been referred to as the natural evolution of the notions of proximity and continuity [6]. The development of topological methods and algorithms that bring this theory into practice are emerging, and specific applications of topology in data mining tasks have been referred to as *topological data analyses* (TDA). Perhaps the leading tool for applying TDA in practice is Mapper, developed by Singh

*Department of Computer Science, Wright State University, Dayton, OH (piekenbrock.5@wright.edu, <http://www.wright.com/~piekenbrock.5/>).

†Department of Computer Science, Wright State University, Dayton, OH (derek.doran@wright.edu, <http://knoesis.org/people/derek/>).

et al. [13]. **Mapper** is a coordinatization framework capable of summarizing high-dimensional data sets, or functions of them, and has proven its worth as a useful tool for data visualization [7], genetic recombination characterization [4], and dimensionality reduction [13]. **Mapper** has also had much success in commercial settings.¹ In its simplest interpretation, **Mapper** is a topologically-motivated framework for reducing high dimensional data into an interpretable graph $G(V, E)$, where vertices correspond to clusters of data, and edges correspond to interactions (non-empty intersections) between such clusters. The goal of **Mapper** is to provide a simplified representation of high-dimensional data that not only *preserves* certain topological structures of interest, but is amenable to qualitative analysis and visualization. The goal of **Mapper** is *not* to achieve a complete representation of the data, but to capture significant areas of interest in an explainable way, enabling fruitful exploratory data analysis.

As an unsupervised framework capable of generalizing to heterogeneous formats of data, **Mapper** naturally requires a number of parameters. Provided a reference map f , a covering over a metric space $\{\mathcal{U}\}$, a clustering algorithm \mathcal{C} using some distance metric d , and their corresponding parameters, **Mapper** deterministically generates G . **Mapper** admits a large parameter space, requiring several choices to be made by the user based on the goal of the analysis at hand. It's well known in both unsupervised and supervised learning that more often than not, a large parameter space implies a large solution space: changing the value of any particular parameter in may lead to a different set of solutions, which may in turn lead to inconsistent interpretations and analyses. This is true for the **Mapper** construction as well. Solutions which reduce this complexity would not only enable more robust analyses of the topological solutions produced by **Mapper**, but may indeed accelerate the adoption of TDA as a useful tool for data analysis compared to other forms of learning.

The focus of this effort is to analyze the *overlap* parameter of **Mapper**, which we describe further in Section 1, as it solely determines the connectivity, and thus the relational aspect, of the output graph G . That is, when the overlap parameter is set to 0, the graph is completely disconnected, and as the parameter increases, the graph generally becomes more connected. Unfortunately, there is little guidance on how to determine parameter values relevant to the data set being analyzed. If the overlap is too high, spurious connections may be shown in the output graph, conveying a relational structure that is misrepresentative of the topology. Conversely, if it is too low, subtle yet important connections may be missed.

Is there an optimal scale that best captures the topology of the data with respect to **Mapper**'s overlap parameter? To answer this, it's worth exploring how such scaling choices are performed for other types of TDA. Consider a well-known topological construction, such as the Vietoris-Rips Complex at scale ϵ , which is a simplicial complex obtained by adding n -dimensional simplices spanning the points of X if the pairwise intersection of $B(X_1, \epsilon) \cap B(X_2, \epsilon) \cap \dots \cap B(X_n, \epsilon)$ is non-empty, where $B(X_i, \epsilon) = \{x \in X \mid d(X_i, x) \leq \epsilon\}$. As the sole parameter which determines connectivity, the ϵ parameter is analogous to the overlap parameter for **Mapper**, and as Ghrist argues: "*Despite being both computable and insightful, the homology of a complex associated to a point cloud at a particular ϵ is insufficient: it is a mistake to ask which value of ϵ is optimal*" [5]. We extend this assertion to **Mapper**: a fixed

¹The **Mapper** framework actually lies at the core of the **Ayasdi platform** (<https://www.ayasdi.com/solutions/>)

scale parameter is insufficient if the task is study the structure in the context of *homology*. Rather, the solution is to look at the *continuum* of possible values towards understanding how the complex changes between ranges of parameter settings. Intuitively, substructures which seem to “persist” across parameter ranges forms a foundation for many approaches in Persistent Homology [12]. This solution extends to exploratory use-cases of Mapper as well—rather than choose an optimal scale to extract the “best” output of Mapper, a better solution may be to allow an analyst to explore the space of possible parameter values. By viewing multiple topological representations of the data at once, an analyst could compare and contrast solutions, derive insights, and extract richer knowledge on the topology of the data. The analyst could also further understand how the topological representation of the data (i.e. the output graph of Mapper) changes across other settings of interest.

This notion of studying topological constructions on a continuum of parameter ranges is not new. Methods which measure some notion of ‘stability’ or ‘persistence’ of a given real-valued function with respect to its domain have proven useful in both theoretical and practical instances of TDA [5]. In the context of the Mapper construction explicitly, theoretical extensions which analyze Mapper at multiple resolutions are emerging [2, 3]. Nonetheless, using such extensions to *performing TDA on empirical data* requires knowledge of how the Mapper construction changes with respect to ranges of parameter inputs, and what those parameter ranges are. Such types of analysis may only be enabled tractably, however, if the set of valid and useful values for the overlap parameter can be computed efficiently.

In this work, we derive a closed-form solution to compute the set of minimal parameter values that, at a given resolution, produce ‘distinct’ topological renderings under the Mapper algorithm, with the exact definition of ‘distinct’ to be defined later in Section 2. We further show that the solution can be computed very efficiently on a data set in $O(nd)$ time, where n is the size of the data set (i.e. the number of ‘points’ in the point cloud data), and d is the dimensionality of the space. We outline the benefits in tractability our approach enables in the context of 1-skeleton generation and in other application areas, and we note that our effort may also allow similar computational gains in generating the full simplicial complex.

In the next section, we discuss some prerequisite background material including preliminary notation that will be used throughout the paper. We then give a step-by-step overview of Mapper with examples. In Section 2, we derive a closed form solution for finding the minimal set of values for the overlap setting. In Section 3, examples are given and discussed in detail to further convey the motivation, significance, and utility of our solution. Finally, we discuss future applications and work related to Mapper in Section 5.

1. Background. This section provides a succinct summary of the Mapper algorithm, leaving the reader to refer to the original paper by Singh *et. al* [13] for deeper background material. An even greater treatment of the topological ideas underlying Mapper can be found in Carlsson’s exposition [1]. Still further, Merkulov *et. al* [10] provides much information on simplicial complexes, and Ghrist [6] may be consulted for a high-level overview of using topological constructs for applied topology.

We start by providing notation and some well-known topological constructs needed to describe Mapper and our corresponding derivation. Let S be a discrete set. An *abstract simplicial complex* is a collection X of finite subsets of S , closed under restriction. That is,

for each $\sigma \in X$, all subsets of σ are also in X . Each element $\sigma \in X$ is called a k -simplex whose faces are the simplices corresponding to all subsets of $\sigma \subset S$. The standard k -simplex is defined:

$$\Delta^k := \left\{ x \in [0, \infty)^{k+1} : \sum_{i=0}^k x_i = 1 \right\}$$

For simplicity, we may occasionally use the terms ‘vertex’, ‘edge’, ‘triangle’, and so on to describe the 0, 1, 2, and higher dimensional simplexes, respectively. The resulting topological construction returned by **Mapper** is a *geometric realization* of X , inductively defined by the k -skeletons of X , $k \in \mathbb{N}$, to be the quotient space:

$$X^{(k)} := \left(X^{(k-1)} \bigcup \prod_{\sigma: \dim \sigma = k} \Delta^k \right) / \sim$$

where \sim is the equivalence relation that identifies faces of Δ^k with the corresponding combinatorial faces of σ in $X^{(j)}$ for $j < k$. Thus, the 0-skeleton $X^{(0)} = S$ is a discrete set. The 1-skeleton $X^{(1)}$ is a space homeomorphic to a collection of vertices (S) and edges connecting vertices, i.e. a topological graph. Define a *cover* of a set X as a collection of sets U_α whose union contains X as a subset, i.e. $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ is a cover of X if $X \subseteq \bigcup_{\alpha \in A} U_\alpha$ for some index set A . In this paper, we will always assume that each U_α is path-connected and a cover means a finite open cover. Define the *nerve* of a covering \mathcal{U} as the simplicial complex $N(\mathcal{U})$ whose k -simplices correspond to non-empty intersections of $k + 1$ distinct elements of \mathcal{U} , $N(\mathcal{U}) = \{\sigma \subseteq A \mid \bigcap_{\alpha \in \sigma} U_\alpha \neq \emptyset\}$. Given a continuous map $f : X \rightarrow Z$ where Z is equipped with a cover $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$, we will write $f^*(\mathcal{U})$ as the cover of X obtained by considering the path connected components of $f^{-1}(U_\alpha)$ for each α . Given this map f , **Mapper** produces a *topological construction* M defined as the nerve of the $f^*(\mathcal{U})$, $M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$.

Unless otherwise stated, we will use lowercase symbols to represent scalars, capital symbols to represent sets, bold lowercase symbols to represent vectors ($1 \times m$), and bold capital symbols to represent matrices ($n \times m$). With a slight abuse of notation, we will occasionally express vectors with a scalar in **typewriter** font when the vector contains identical elements along all of its dimensions, e.g. $\mathbf{k} = k$ implies $\mathbf{k} = [k_1, k_2, \dots, k_m]$ where each $k_i = k \forall i \in [1, m]$. This is used to simplify the mathematical expressions that follow, and their corresponding examples, although we note the derivation still holds when the components are not identical. This simplification also aligns our analysis more closely with several practical applications of **Mapper**—many of the algorithmic implementations of **Mapper** allow the user to input scalar values out of simplicity, which then convert to their vector-valued equivalents internally. This is elaborated on in the next section.

1.1. The Mapper Algorithm. Mapper constructs a simplicial complex from a given set of ‘point cloud data’ (PCD) $X \subset \mathbb{R}^d$ by the following (informally presented) process:

1. Define a reference map $f : X \rightarrow Z$ where Z is some reference metric space. Note that because Z is a metric space, it’s assumed that it is possible to compute inter-point distances between the points in the data of Z .
2. Select a finite covering $\{U_\alpha\}_{\alpha \in A}$ of Z .
3. Construct the subsets $X_\alpha = f^{-1}U_\alpha$. Since f is continuous, these sets also form an open covering of X .

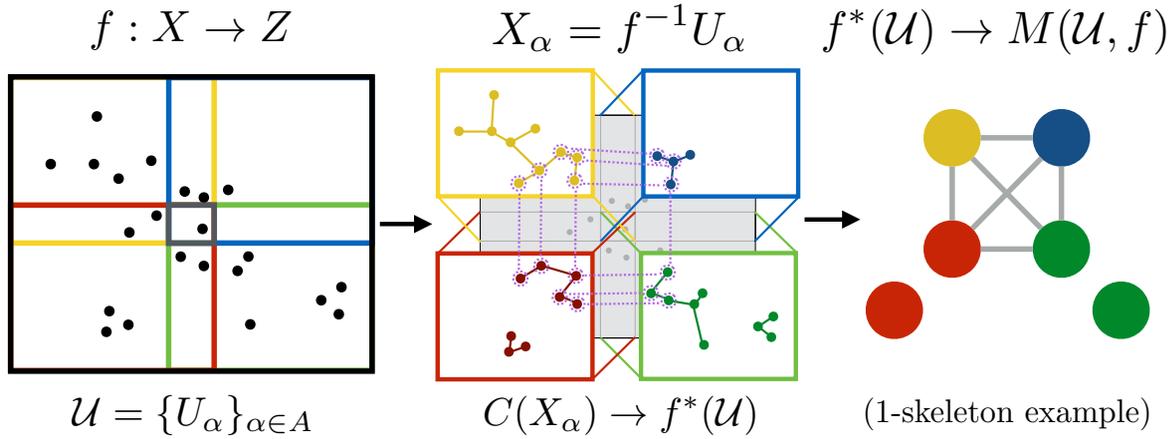


Figure 1. Mapper step by step example. The left panel shows the mapping of a given data set $X \rightarrow Z$ via f (where $Z \in \mathbb{R}^2$) and four overlapping sets that construct the covering \mathcal{U} of Z . The center figure shows the isolation the points in each of these sets that is used for partial clustering and the connected components represent the results of this clustering. Points which appear in multiple subsets are represented by the dashed lines. The final figure on the right represents the 1-skeleton realization of the simplicial complex constructed from the nerve of the covering.

4. Given a clustering algorithm \mathcal{C} , construct the set of clusters by applying \mathcal{C} to the sets X_α . This induces a covering of X parametrized by pairs (α, c) where $\alpha \in A$ and c is one of the clusters of X_α .
5. Construct the simplicial complex $N(\mathcal{U})$ whose vertex set is the set of all possible such pairs (α, c) and where a family $\{(\alpha_0, c_0), (\alpha_1, c_1), \dots, (\alpha_k, c_k)\}$ spans a k -simplex if and only if the corresponding clusters have a point in common.

Proceeding item-wise through these steps, one begins a TDA with **Mapper** by defining a reference map, sometimes referred to as a *filter* function, which takes as input a set of ‘point cloud data’ defined in the *data space* and maps to an appropriate metric space, which we’ll refer to as the *filter space* or Z -space. The outcome produced by **Mapper** is highly dependent on this function. The choice of this function is inevitably application-specific, however details on common functions that may of interest are discussed in [13]. In step 2, a parameterized covering is constructed on the Z -space, partitioning subsets of the data into subsets in the filter space based on which set of the covering they lie within. We will occasionally refer to these sets which compose the cover as *level sets*. Because level sets generally overlap, any given data point may be within multiple level sets. Since f is continuous, the constructed cover of Z also forms a covering of X . **Mapper** relies on the idea of *partial clustering*, wherein subsets of X are clustered independently of other subsets as a means of simplifying the data into topological prototypes. These subsets are determined by the parameterization of the covering in the Z -space. The prototypes formed within these level sets are represented as 0-simplexes, the interaction of these prototypes with other 0-simplexes form the higher dimensional simplexes in the resulting simplicial complex. Intuitively, a simplicial complex is an expression of a given space as a union of points, intervals, triangles, and higher dimensional

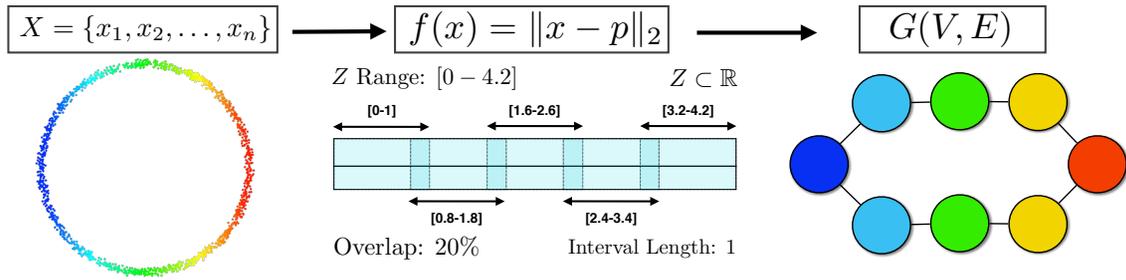


Figure 2. Mapper ring data example discussed in Example 1.1. In this example, $Z \in \mathbb{R}$ represents the distance from each point $x_i \in X$ to the point with the smallest coordinate in $X^{(1)}$, i.e. the “left-most” point. Note how the three inner intervals resulted in two clusters each, which connect at each end point to a single cluster in the first and last intervals, forming a ring.

analogues [1]. Figure 1 illustrates, at a high level of abstraction, how the Mapper algorithm create a simplicial complex from a given set of data.

Example 1.1: Mapper Ring Data Example

To further illustrate Mapper, consider the example given by Singh *et. al* in [13], shown in Figure 2. The left figure displays the ring data X in \mathbb{R}^2 . The filter function f used computes the distance of every point $x \in X$ to the left-most point p , resulting in a filter space in \mathbb{R} . This filter space is shown in the center figure, from which a covering is constructed of 5 sets, each with 20% overlap (the points lie in the filter space are not shown). The right figure shows the resulting 1-skeleton of the constructed simplicial complex, represented as a graph G .

1.2. Cover Parameterization. Many of the parameters of Mapper are user supplied. This includes the reference map (step 1), all parameters related to forming the cover (steps 2-3), and the clustering algorithm (as well as its corresponding hyper-parameters) (step 4). Because steps 1 and 4 are generally data set specific, we will focus on steps 2-3: constructing a cover within the Z -space.

There are numerous ways to construct a suitable covering. Although it’s assumed each level set U_α is path-connected, and that the cover \mathcal{U} is finite, Mapper makes no assertions restricting the shape of the level sets used to construct the covering. For example, Figure 2 of [13] demonstrates two potential coverings of the parameter space in \mathbb{R}^2 , one with rectangles and another with hexagons. In practice, perhaps due to its simplicity or its computational tractability, a rectangular covering is often used. As a 2-parameter family of coverings, the rectangular cover generally requires the user to supply 1) a tuple of positive integers representing the number of intervals to construct (per dimension) and 2) a tuple of positive real values representing the percent overlap to allow between adjacent intervals (per dimension). Examples of this cover are shown in Figures 1 and 2.

Consider steps 2 and 3 of the Mapper framework: constructing a finite covering $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ of Z , given a user-supplied continuous map $f : X \rightarrow Z$. The hyper-rectangular covering method requires as input the number of rectangles $\mathbf{k} = [k_1, k_2, \dots, k_m]$ to distribute

along each dimension (where $\mathbf{k} \in \mathbb{Z}^m$), and an overlap percentage $\mathbf{g} = [g_1, g_2, \dots, g_m]$ neighboring rectangles should intersect (where $\mathbf{g} \in [0, 1]^m$). One proceeds constructing the cover by first finding the range of the function f restricted to the set of given points Z . That is, denote the i^{th} point of Z as \mathbf{z}_i , and its j^{th} component of this point as $z_i^{(j)}$. Similarly, denote the j^{th} component Z as $\mathbf{z}^{(j)}$. The extrema vectors of Z are defined as follows:

$$\begin{aligned} \mathbf{z}_{max} &= [\max(\mathbf{z}^{(1)}), \max(\mathbf{z}^{(2)}), \dots, \max(\mathbf{z}^{(m)})] \\ \mathbf{z}_{min} &= [\min(\mathbf{z}^{(1)}), \min(\mathbf{z}^{(2)}), \dots, \min(\mathbf{z}^{(m)})] \end{aligned}$$

And the range vector is the component-wise difference between them:

$$\bar{\mathbf{z}} = \mathbf{z}_{max} - \mathbf{z}_{min}$$

Given a user-supplied input for \mathbf{g} and \mathbf{k} , a cover $\mathcal{U}[\mathbf{r}, \mathbf{e}]$ is constructed where:

$$(1) \quad \mathbf{r} = \frac{\bar{\mathbf{z}}}{(\mathbf{k} - \mathbf{g}(\mathbf{k} - 1))}, \quad \mathbf{e} = \mathbf{r} \circ (1 - \mathbf{g})$$

We'll refer to \mathbf{r} as the *interval length* of a given covering, and \mathbf{e} as the *step size* (which is inversely proportional to the overlap \mathbf{g} with respect to \mathbf{r}). When $\mathbf{k} > 1$, this cover implies Z is effectively divided up into a number of overlapping, rectangular level sets. Each point in Z lies within at least one level set of $\mathcal{U}[\mathbf{r}, \mathbf{e}]$. To formalize operations on these sets, it's necessary to express this cover as an *indexed family* of sets. That is, the cover $\mathcal{U}[\mathbf{r}, \mathbf{e}] = \{U_\alpha\}_{\alpha \in A}$ of the space Z is an indexed family of level sets U_α , each of length \mathbf{r} , indexed by a given index set A represented by the mapping function $A_\alpha : A \rightarrow U_\alpha$.

From a theoretical perspective, there is no restriction on how the index set is represented. But consider the case where $\mathbf{k} = \mathbf{k}$ and A is given by the m -fold Cartesian product of \mathbf{k} , i.e.

$$A = \kappa^m = \underbrace{\kappa \times \dots \times \kappa}_m = \{(l_1, l_2, \dots, l_m) \mid l_i \in \kappa \text{ for all } i = 1, \dots, m\}$$

where $\kappa = [0, 1, \dots, \mathbf{k} - 1]$. This indexing scheme offers the practical advantage of not only being easy to compute, but it also explicitly records locational information with regards to the position of each level set in the cover. For example, when this indexing scheme is used, the level set indexed by the tuple $(0, 0)$ contains any point in Z between $[\mathbf{z}_{min}, \mathbf{r})$, and similarly the level set indexed by $(\mathbf{k} - 1, \mathbf{k} - 1)$ will contain any point in Z between $[\mathbf{r} \circ (\mathbf{k} - 1), \mathbf{z}_{max}]$. This approach has additional advantages, which we elaborate on in the next section.

2. Cover Analysis. We now derive the main result of the paper: a closed-form solution to find the minimal set of overlap values which may be used to compute distinct simplicial complexes. As discussed above, we specifically consider the case where a (possibly multivariate) hyper-rectangular covering method is used, and we focus on the case where the number of intervals \mathbf{k} is fixed.² Additionally, we only concern ourselves with the 1-skeleton, which is often sufficient for most types of analysis [].

²This is the only covering method supported by the *Python Mapper* library for \mathbb{R}^d when $d > 1$ and the R package *TDAmapper*.

The goal of the following analysis is to create a method for computing the minimal set of input values \mathbf{g} which realize all possible *distinct* simplicial complexes in a closed-form solution. To formalize this, it's necessary to first define what is meant by a 'distinct' solution. Recall that the **Mapper** construction provides a discrete, combinatorial summary of the data used to construct it, $M(\mathcal{U}_r, f)$, given by the nerve of $f^*(\mathcal{U}_r)$. Although the overlap parameter \mathbf{g} lies on a continuum, there may be ranges of values which produce identical simplicial complexes. For example, consider a covering constructed from $\mathbf{g} = \mathbf{g}$ and another from $\mathbf{g}' = \mathbf{g} + \epsilon$ for a sufficiently small ϵ . Denote the corresponding interval lengths as \mathbf{r} and \mathbf{r}' . Since \mathbf{k} is fixed, \mathbf{r}' will be slightly larger than \mathbf{r} in the corresponding covers, \mathcal{U}_r and $\mathcal{U}_{r'}$. However, if the point membership of each level set does not change between the two covers, then the corresponding dissimilarities between all points within each level set will also not change; there will be no new vertices or edges introduced between the two **Mapper** constructions. This motivates a definition of *distinction* between possible realized values of \mathbf{g} .

Definition 2.1 (Distinct Simplicial Complex). *Given a finite set of point cloud data $\mathbf{Z}_n \in Z$ where $Z \subset \mathbb{R}^m$ and two covers $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ and $\mathcal{V} = \{V_\alpha\}_{\alpha \in A}$ of Z which share the same index set (but may index different families), denote the corresponding membership sets of $\mathbf{Z}_n \cap U_\alpha$ and $\mathbf{Z}_n \cap V_\alpha$ as follows:*

$$(2) \quad \begin{aligned} \mathcal{U}_\alpha(\mathbf{Z}) &= \{z : z \in \mathbf{Z}_n \text{ and } z \in U_\alpha\} \\ \mathcal{V}_\alpha(\mathbf{Z}) &= \{z : z \in \mathbf{Z}_n \text{ and } z \in V_\alpha\} \end{aligned}$$

We define the simplicial complex $N(\mathcal{U})$ as being distinct from the simplicial complex $N(\mathcal{V})$ if:

$$(3) \quad \text{Any } \mathcal{U}_\alpha(\mathbf{Z}) \neq \mathcal{V}_\alpha(\mathbf{Z}) \quad \forall \alpha \in A$$

Thus, given two coverings with the same index set, if there is at least one membership set that is not equal for some index, then the resulting simplicial complexes are considered distinct.

It's worth noting that under this definition of distinctness, it is possible to observe two **Mapper** constructions whose 1-skeletons share the same simplex representation. That is, it is *not* guaranteed that . However, under this definition, the *mapping* of simplices to points in \mathbf{Z}_n will be different.

With an established notion of a distinct simplicial complex, we begin our derivation by assuming we are given a finite point cloud $\mathbf{Z}_n \in Z$, and consider the special case where the topological space Z is zero-dimensional with respect to its Lebesgue covering dimension. This implies that each open set of the cover of Z is disjoint, or equivalently $\mathbf{g} = 0$. For brevity, we'll occasionally refer to this case as the *base cover*. Under this case, every point $z_i \in \mathbf{Z}_n$ is contained in exactly one set of the cover, and additionally, is mapped by exactly one index $\alpha \in A$. Let this initial mapping be expressed as $\psi : \mathbf{Z}_n \rightarrow \alpha$, where $\alpha \in A$, and denote \mathbf{A} as the matrix of mappings by ψ over a set of n data points (i.e. where row \mathbf{A}_i corresponds to mapping $\psi(z_i)$).

$$(4) \quad \mathbf{A} = [\alpha_0 \quad \alpha_1 \quad \dots \quad \alpha_n]^\top = \begin{bmatrix} \psi(z_0) \\ \psi(z_1) \\ \vdots \\ \psi(z_n) \end{bmatrix} \quad (\text{where } \mathbf{A} \in \mathbb{Z}^{n \times m})$$

Given a set of *fixed* parameters, the parameters of the cover \mathbf{r} and \mathbf{e} are traditionally defined as in (1). We first need to establish the parameterization of the cover with the defined mapping \mathbf{A} , i.e. when $\mathbf{g} = 0$. In doing so, we define the *base interval length* as follows:

$$(5) \quad \bar{\mathbf{r}} = \bar{\mathbf{z}} \circ \mathbf{k}^{-1}$$

Since \mathbf{e} is inversely proportional to \mathbf{g} , the corresponding step size $e = 0$. This establishes a lower bound on the interval length (any interval length below this value is not a cover of Z), and provides a starting point for our solution. Given a base interval length, each point \mathbf{z}_i is contained in exactly one level set α_i , given by the map $\psi(\mathbf{z}_i)$, which we'll refer to the points' *originating set*. A simple illustrations of the realizations of these definitions for a small set of points is given in Example 1.1.

Example 1.1

Consider a set of points $\mathbf{Z}_n \subset \mathbb{R}^2$ with the defined range vector $\bar{\mathbf{z}} = [2, 4]$. Suppose $\mathbf{k} = 2$, i.e. $\mathbf{k} = [2, 2]$. Assume a base cover is be constructed $\mathcal{U}[\bar{\mathbf{r}}, 0]$ where $\bar{\mathbf{r}} = [1, 2]$. The index set A is given by the cartesian product $\kappa \times \kappa = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$, and it's corresponding family of level sets is given by $\{\{[0, 1], [0, 2]\}, \{(1, 2], [0, 2]\}, \{[0, 1], (2, 4]\}, \{(1, 2], (2, 4]\}\}$. Observe that when \mathbf{k} is fixed, the family of coverings of $\mathcal{U}[\bar{\mathbf{r}}, \mathbf{e}]$ for all values of \mathbf{e} have the same indexing set.

When $\mathbf{g} = 0$, there are no simplices of dimension ≥ 1 in the resulting Mapper construction, i.e. M is a fully disconnected graph. Consider the case when $\mathbf{g}' = 0 + \epsilon$ for a sufficiently small ϵ . What is the smallest ϵ that will produce a distinct simplicial complex? According to our definition of a distinct simplicial complex 2.1, this question is equivalent to asking: what is the smallest overlap percentage $\mathbf{g}' > 0$ that induces a set membership change, where $\mathcal{U}_\alpha(\mathbf{Z})$ corresponds to the cover $\mathbf{g} = 0$, and $\mathcal{V}_\alpha(\mathbf{Z})$ corresponds to the cover $\mathbf{g}' > 0$? Clearly, we must consider overlap values $\mathbf{g}' > \mathbf{g}$, and thus interval lengths larger than $\hat{\mathbf{r}} > \bar{\mathbf{r}}$. For the membership sets between two parameterizations to be different, a point $\mathbf{z}_i \in \mathbf{Z}_n$ must intersect another level set, outside of its originating set. To compute this efficiently, it's necessary know the *relative* position of each point within it's originating set, such that the metric distance between any pair of points from the same originating level set is unchanged. Consider a projection of a point \mathbf{z}_i onto a single index set, such that $0 \leq z_i^{(d)} \leq \bar{r}^{(d)} \forall d \in \{1, 2, \dots, m\}$, with projected positions $\tilde{\mathbf{Z}}$ computed as follows:

$$(6) \quad \tilde{\mathbf{Z}} = (\mathbf{Z}_n - \mathbf{z}_{min}) - \mathbf{A} \circ \bar{\mathbf{r}} \quad (\text{where } \tilde{\mathbf{Z}} \in \mathbb{R}^{n \times m})$$

Equation 6 describes $\tilde{\mathbf{Z}}$ as a translation of Z to the origin, followed by a projection of every point to the level set $[0, \bar{r}]$. Note that this projection is only possible with the original surjective mapping of $\psi(\mathbf{z}_i) \rightarrow \alpha_i$, where $\alpha_i \in \mathbb{Z}^m$ is a vector of integers representing the indexed location of each point in the base cover. A visual illustration of the projection is shown in Figure 3. Given $\tilde{\mathbf{Z}}$, we can collectively compute the distance of every point in \mathbf{Z}_n to its *nearest* adjacent set, per dimension, by considering which halfspace each individual point lies within. If a given point $\tilde{\mathbf{z}}_i$ lies in the lower halfspace along dimension d , then the nearest level set along dimension d that \mathbf{z}_i is simply the d^{th} component of $\tilde{\mathbf{z}}_i$. If $\tilde{\mathbf{z}}_i$ lies within the upper halfspace

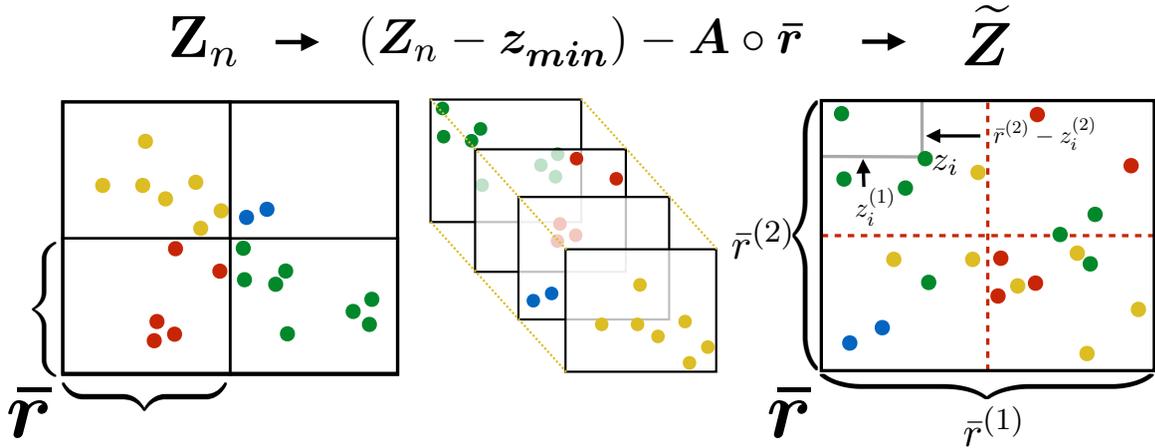


Figure 3. The left panel shows the given \mathbf{Z}_n in the Z -space (where $Z \in \mathbb{R}^2$) and four overlapping sets that construct the covering \mathcal{U} of Z . The center figure shows conceptually how the projection is aligning each level set into one level set. The final figure on the right illustrates how each points corresponding $\tilde{z}_i^{(d)}$ value depends on which halfspace (delimited with red dashed lines) it lies in. Note that the left panel is showing four level sets, and the right panel is showing only one defined on the interval $[0, \mathbf{r}]$.

along dimension d , then the distance to the nearest adjacent level set is given by $\bar{r}^{(d)} - \tilde{z}_i^{(d)}$. These simple observations express all the information needed to compute the entire finite set of \mathbf{g} that produce distinct simplicial complexes, and in the next section, we consider two use-cases of where these observations become useful for 1-skeleton generation.

2.1. Restricted Cover. In practice, it's often sufficient to assume reasonable restrictions on parameter ranges of the cover. For example, if we consider using a rectangular cover over \mathbb{R}^2 under the case where $e < \frac{r}{2}$, there will exist regions where at most four rectangles will intersect, and thus the dimension of the simplices used to construct the complex will be 3 or less. We first analyze this situation in the context of our end-goal.

Consider the projection previously discussed in equation 6 and the corresponding component wise differences, $\tilde{z}_i^{(d)}$ if \tilde{z}_i lies in the lower halfspace of $\tilde{\mathbf{Z}}$ along dimension d , and $\bar{r}^{(d)} - \tilde{z}_i^{(d)}$ if it lies in the upper halfspace. Since these distances are complements to each other with respect to \mathbf{r} , regardless of which halfspace $\tilde{z}_i^{(d)}$ lies within, the distance to the closest adjacent level set may be simplified by minimum of the two values: $\delta(z_i) = \min(\tilde{z}_i^{(d)}, \bar{r}^{(d)} - \tilde{z}_i^{(d)})$. Denote this mapping as follows:

$$(7) \quad \mathbf{Z}_\Delta = \begin{bmatrix} \delta(\mathbf{z}_1^{(1)}) & \delta(\mathbf{z}_1^{(2)}) & \dots & \delta(\mathbf{z}_1^{(m)}) \\ \delta(\mathbf{z}_2^{(1)}) & \delta(\mathbf{z}_2^{(2)}) & \dots & \delta(\mathbf{z}_2^{(m)}) \\ \vdots & \vdots & \ddots & \vdots \\ \delta(\mathbf{z}_n^{(1)}) & \delta(\mathbf{z}_n^{(2)}) & \dots & \delta(\mathbf{z}_n^{(m)}) \end{bmatrix} \quad (\text{where } \mathbf{Z}_\Delta \in \mathbb{R}^{n \times 1})$$

These minimal distances, when added to the base interval length, define the *smallest* interval length (for a given point \mathbf{z}_i) $\hat{\mathbf{r}}$ which induces a change in the membership set of the

corresponding cover \mathcal{U} . Furthermore, they can be computed simply:

$$(8) \quad \hat{\mathbf{R}} = \mathbf{Z}_\Delta + \bar{\mathbf{r}} \quad (\text{where } \hat{\mathbf{R}} \in \mathbb{R}^{n \times m})$$

Note that because these distances are all orthogonal, we have all of the information needed to compute when *any* set of the covering will intersect the originating set of a given point \mathbf{z}_i . For example, consider any subset in the cover which is not along one of the orthogonal directions in $\hat{\mathbf{r}}$. The smallest interval length wherein a given point \mathbf{z}_i outside of this set will intersect the subset is achieved when $\mathbf{r} = \mathbf{r}_{max}$ where \mathbf{r} is an m -length vector of values all set to $\max(\hat{\mathbf{r}})$. Since the step size e is a function of \mathbf{g} and \mathbf{r} , and because \mathbf{k} is fixed, knowledge of the smallest interval length(s) which cause a level set membership change enables \mathbf{g} to be computed in closed-form. Recall that fixed Mapper traditionally takes as input the number of intervals to divide the filter space along each dimension, \mathbf{k} , and the overlap (percentage) between these sets \mathbf{g} , as shown in equation 1. We take the inverse of this function, supplying the set of $\hat{\mathbf{R}}$ values computed from the \mathbf{Z}_n point cloud as input with the fixed \mathbf{k} to recover \mathbf{g} . The derivation, given in the appendix 5, yields:

$$(9) \quad \mathbf{G} = \frac{(\bar{\mathbf{z}} - \hat{\mathbf{R}}\mathbf{k})}{\hat{\mathbf{R}}(-\mathbf{k} + 1)} \quad (\text{where } \mathbf{G} \in \mathbb{R}^{n \times m})$$

Given \mathbf{G} , the step size can be computed analogously:

$$(10) \quad \hat{\mathbf{E}} = \hat{\mathbf{R}}(1 - \mathbf{G})$$

Each entry $\mathbf{g}_i^{(j)} = \mathbf{G}_{ij}$ expresses the smallest overlap percentage necessary for the corresponding point \mathbf{z}_i to intersect it's nearest adjacent level set along the j^{th} dimension. Since we only consider the restricted cover parameterization in this section, these values encompass all of the possible values to produce distinct simplicial complexes. More specific uses cases for when \mathbf{G} are available is discussed more in Section 3. We next discuss the case where $e \geq \frac{r}{2}$, or equivalently, where the overlap is less than 50%.

2.2. Unrestricted Cover. In this situation, \mathbf{g} is not bounded above by 0.50. As a result, the value of $\mathbf{z}_i^{(d)}$ may intersect more than just its nearest and originating level sets. For example, consider a point $z = 2.1$ in 2, but instead consider the case where the cover was constructed with an overlap of 80%. When $k = 5$, the corresponding interval length comes out to ≈ 2.41 . In this situation, z intersects every level set of the cover. This implies that a more general formulation of equation 9 is needed. Note that this is only needed when $\mathbf{g} \geq 0.50$.

We use the same strategy as before, but now replace equation 7 with two variables, \mathbf{Z}_{near} and \mathbf{Z}_{far} . Set $\mathbf{Z}_{near} = \mathbf{Z}_\Delta$, and compute \mathbf{Z}_{far} analogously but instead by looking at the distance to the *farther* level set, i.e. $\delta(\mathbf{z}_i) = \max(\tilde{\mathbf{z}}_i^{(d)}, \bar{\mathbf{r}}^{(d)} - \tilde{\mathbf{z}}_i^{(d)})$. Additionally, we add natural multiples of $\bar{\mathbf{R}}$ to account for level sets which are not immediately adjacent to the originating level sets. The \mathbf{p}^{th} nearest level sets correspond to:

$$(11) \quad \begin{aligned} \hat{\mathbf{R}}_{near} &= \mathbf{Z}_{near} + (\mathbf{p} - 1)\bar{\mathbf{R}} \quad (\text{where } \hat{\mathbf{R}}_{near} \in \mathbb{R}^{1 \times m}) \\ \hat{\mathbf{R}}_{far} &= \mathbf{Z}_{far} + (\mathbf{p} - 1)\bar{\mathbf{R}} \quad (\text{where } \hat{\mathbf{R}}_{far} \in \mathbb{R}^{1 \times m}) \end{aligned}$$

where \mathbf{p} is an m -dimensional vector where each component $\mathbf{p}^{(d)} \in [1, 2, \dots, \mathbf{k}^{(d)}]$.

Note that the magnitude of each component of \mathbf{p} dramatically increases the total number of overlap values. Consider the case where $\mathbf{p} = \mathbf{k}$.

2.3. Practical considerations. In the previous sections, we outlined methods to compute the full finite set of overlap values under both the cases where $0 \leq g < 0.50$ and $0.50 \leq g < 1$. The analysis performed so far has generally assumed that all level sets within distance $\hat{\mathbf{R}}$ were *valid* level sets in the cover. However, it may be the case that for a given point z_i , the level set(s) that were considered in computing $\hat{\mathbf{R}}$ lie outside of the cover. For example, consider a point $z_i = 0.2$ positioned in the filter space of the ring data set from example 2. When $0 \leq g < 0.50$, z_i never intersects a level set beyond its originating set. Because z_i lies at the extrema of the filter space, the halfspace distance computed for z_i does not correspond to a level set of the cover, and therefore its corresponding value in G must be discarded. When $0.50 \leq g < 1$, the corresponding overlap values for z_i may be computed using multiples of $\hat{\mathbf{R}}_{far}$, as given in equation 11, for all $p = \{1, 2, \dots, 4\}$.

We note that this situation is simply an exception case, does not detract from the validity of the solutions given above, and can be readily handled by simply ensuring the index of the level set whose distance is being computed is within the index set of the cover. Furthermore, these indexes may be readily computed. Consider the use of the m -fold Cartesian product index set described in section 1. If $e < \frac{r}{2}$, these are simply the sets adjacent to \mathbf{A} in the direction of the halfspace each point in \mathbf{Z}_n lies. Conceptually, this index set $\hat{\mathbf{A}}$ represents the mapping as follows:

$$(12) \quad \hat{\mathbf{A}} = \begin{bmatrix} \psi(\mathbf{z}_0 + \hat{\mathbf{R}}_0) \\ \psi(\mathbf{z}_1 + \hat{\mathbf{R}}_1) \\ \vdots \\ \psi(\mathbf{z}_n + \hat{\mathbf{R}}_n) \end{bmatrix} \quad (\text{where } \mathbf{A} \in \mathbb{Z}^{n \times m})$$

Since we directly account for this offset in equation 11, this mapping can be simplified for both the restricted and unrestricted covers as follows:

$$(13) \quad \hat{\mathbf{A}} = \mathbf{A} + \gamma \circ \mathbf{p}$$

$$\text{where } \gamma = \begin{cases} 1 & \text{if } \bar{r} - \tilde{\mathbf{Z}} < \frac{\hat{\mathbf{R}}}{2} \\ -1 & \text{if } \bar{r} - \tilde{\mathbf{Z}} \geq \frac{\hat{\mathbf{R}}}{2} \end{cases}$$

Any index $\hat{\mathbf{A}} \notin A$ corresponds to invalid overlap parameter, and can be excluded through a simple post-processing step.

3. Implementation. In 2013, Müllner et. al introduced the *Python Mapper* software tool [11] as a limited yet extensible realization of the **Mapper** framework for data analysis. The *Python Mapper* library supports three types of covers for \mathbb{R} , and one type of cover for \mathbb{R}^d (multidimensional ‘patch’ covering, the one studied in this effort). Similarly, Pearson [] has an R package *TDAmapper* which relies on the hyper-rectangular covering method as well. All of the derived equations as part of this effort are available in both Python and R scripts, available

open-source online, as unofficial extensions to these libraries.³ In the following subsections, we demonstrate the utility of this result with a few example data sets previously analyzed with Mapper.

4. Applications and Discussion. There are several applications of Mapper enabled from this work. We outline only a few here.

4.1. Improvement in Efficiency. TODO

4.2. Multiscale Mapper: Expanding n -orthotopes. TODO

4.3. Sensitivity and Stability. TODO

5. Conclusion and Future Work. In this paper, we've derived a closed-form solution to computing all of the possible parameter values to the overlap parameter in Mapper, enabling efficient computation of not only the 1-skeleton of the simplicial complex traditionally produced by Mapper, but also of the full simplicial complex. Theoretical evidence was shown indicating that this solution may also be used to compare output constructions from Mapper very quickly across a discrete set of parameter settings, which in turn may provide a useful foundation for both comparative and exploratory analysis of how the topological structure of G changes as a function of such parameter changes. Finally, we included a small discussion of how the solution provided enables further development in the application of topological data analysis to various theoretical and application areas.

Acknowledgement. This work was supported by an appointment to the Internship/Research Participation Program at the U.S. Air Force Research Laboratory (AFRL), administered by the Oak Ridge Institute for Science and Education through an interagency agreement between the U.S. Department of Energy and Wright-Patterson Air Force Base.

Appendix: Overlap Inverse. Traditionally, given the number of intervals k and the percentage overlap between consecutive intervals g , Z may be equipped with a hyper-rectangular covering $\mathcal{U}[\mathbf{r}, \mathbf{e}]$ given one has a finite point cloud $\mathbf{Z}_n \subset Z$, where \mathbf{r} and \mathbf{e} are computed as follows:

$$(14) \quad \mathbf{r} = \frac{z}{(\mathbf{k} - g(\mathbf{k} - 1))}, \quad \mathbf{e} = \mathbf{r}(1 - g)$$

Since \mathbf{e} is dependent on g , and assuming \mathbf{k} is fixed, it's clear that all one needs to compute g is a supplied interval length \mathbf{r} . This is shown below:

$$(15) \quad \begin{aligned} \mathbf{r}(\mathbf{k} - g(\mathbf{k} - 1)) &= \bar{z} \\ \mathbf{r}\mathbf{k} - \mathbf{r}g\mathbf{k} + \mathbf{r}g &= \bar{z} \\ g(-\mathbf{r}\mathbf{k} + \mathbf{r}) &= \bar{z} - \mathbf{r}\mathbf{k} \\ g &= \frac{\bar{z} - \mathbf{r}\mathbf{k}}{\mathbf{r}(-\mathbf{k} + 1)} \end{aligned}$$

It's possible that a user of Mapper could parameterize the cover using \mathbf{r} and k directly instead of g , however determining what values of \mathbf{r} are relevant is difficult since the range of Z will

³TODO:link github

inevitably differ in value and interpretation with respect to the data set and reference map used. This is perhaps the reason its preferred by the open-source implementations of `Mapper` to parameterize the cover by g , which always carries the same interpretation as a percentage. We consider an alternative approach which explicitly computes r from a finite set of data, thus reducing the possible values of g (for a given set of data) to a finite number of parameter settings.

REFERENCES

- [1] G. CARLSSON, *Topology and data*, Bulletin of the American Mathematical Society, 46 (2009), pp. 255–308.
- [2] T. K. DEY, F. MÉMOLI, AND Y. WANG, *Multiscale mapper: Topological summarization via codomain covers*, in Proceedings of the twenty-seventh annual acm-siam symposium on discrete algorithms, SIAM, 2016, pp. 997–1013.
- [3] T. K. DEY, F. MÉMOLI, AND Y. WANG, *Topological analysis of nerves, reeb spaces, mappers, and multiscale mappers*, arXiv preprint arXiv:1703.07387, (2017).
- [4] K. J. EMMETT AND R. RABADAN, *Characterizing scales of genetic recombination and antibiotic resistance in pathogenic bacteria using topological data analysis*, in International Conference on Brain Informatics and Health, Springer, 2014, pp. 540–551.
- [5] R. GHRIST, *Barcodes: the persistent topology of data*, Bulletin of the American Mathematical Society, 45 (2008), pp. 61–75.
- [6] R. W. GHRIST, *Elementary applied topology*, Createspace, 2014.
- [7] C. HEINE, H. LEITTE, M. HLAWITSCHKA, F. IURICICH, L. DE FLORIANI, G. SCHEUERMANN, H. HAGEN, AND C. GARTH, *A survey of topology-based methods in visualization*, in Computer Graphics Forum, vol. 35, Wiley Online Library, 2016, pp. 643–667.
- [8] A. K. JAIN, *Data clustering: 50 years beyond k-means*, Pattern recognition letters, 31 (2010), pp. 651–666.
- [9] T. KOHONEN, *Self-organized formation of topologically correct feature maps*, Biological cybernetics, 43 (1982), pp. 59–69.
- [10] S. MERKULOV, *Hatcher, a. algebraic topology (cambridge university press, 2002)*, Proceedings of the Edinburgh Mathematical Society, 46 (2003), pp. 511–512.
- [11] D. MÜLLNER AND A. BABU, *Python mapper: An open-source toolchain for data exploration, analysis, and visualization*, URL <http://math.stanford.edu/muellner/mapper>, (2013).
- [12] E. MUNCH, *A users guide to topological data analysis*, Journal of Learning Analytics, 4 (2017), pp. 47–61.
- [13] G. SINGH, F. MÉMOLI, AND G. E. CARLSSON, *Topological methods for the analysis of high dimensional data sets and 3d object recognition.*, in SPBG, 2007, pp. 91–100.
- [14] J. B. TENENBAUM, *Mapping a manifold of perceptual observations*, in Advances in neural information processing systems, 1998, pp. 682–688.