

Mapper, Manifolds, and More!

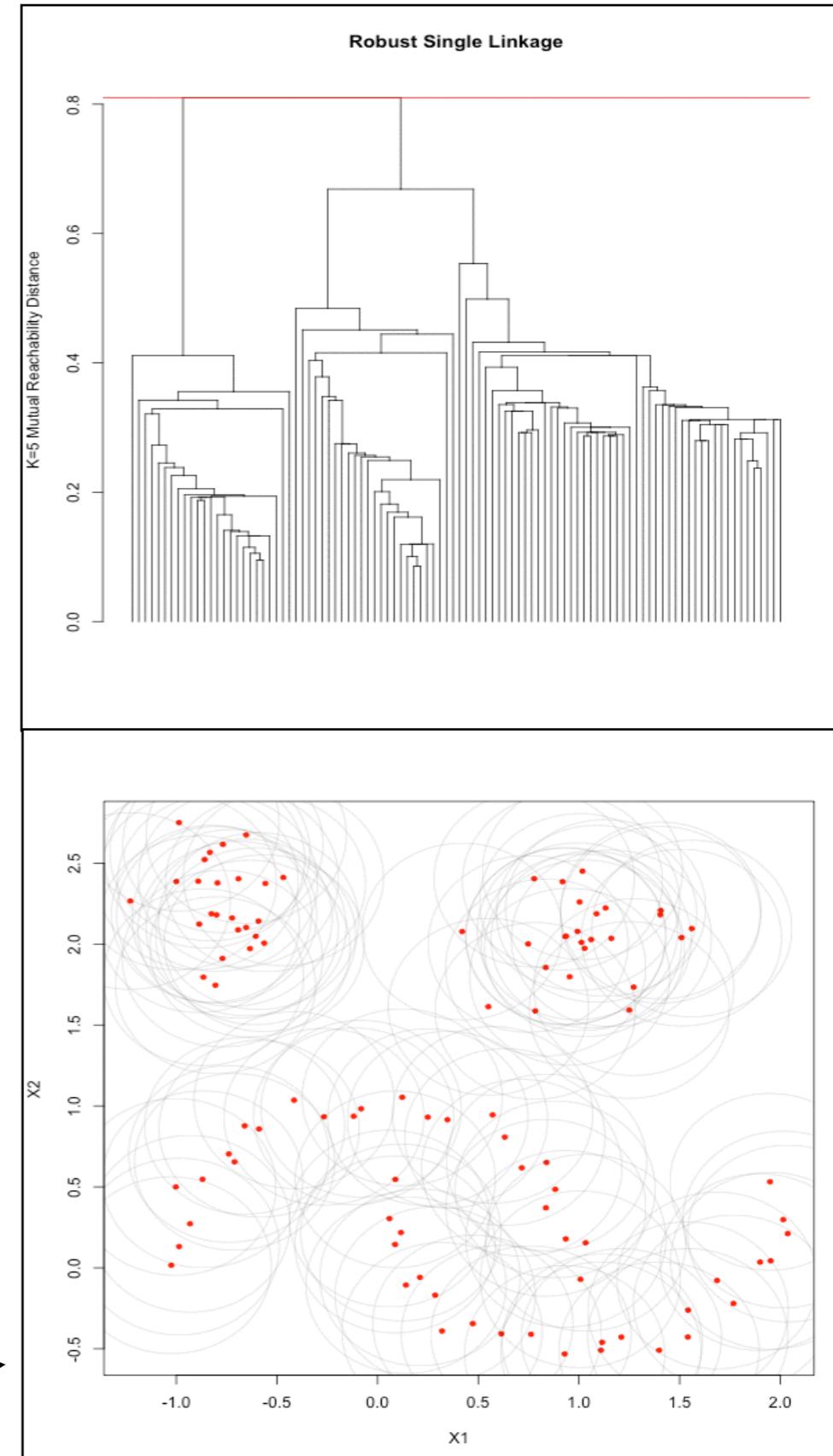
Topological Data Analysis and Mapper

Matt Piekenbrock

Data Science and Security Cluster (DSSC) Talk: March 28th, 2018

My Background

- My Background: Density-based clustering
 - What I've learned: **the best clustering solutions understand how a structure evolves across parameter ranges**
 - Effectively studied how [expanding] spheres relate to density level set estimation
 - Studied some things in theoretical statistics
- Not my Background: **Mathematics**
 - Only studied topology for about a year
 - *Never* taken a real analysis course
 - *Never* taken a theoretical mathematics course
- Summary:
 - Novice in Topology
 - Novice in Manifold learning
 - “Grain of salt” ...
- But I have done some cool things! 



Topology

- Topology is a branch of mathematics concerned with spaces and maps
 - Formalizes notions of *proximity* and *continuity*
 - What topology is useful for:
 - Feature characterization
 - e.g. suitable for classification
 - Formalize notions of robustness
 - “At what *noise level* will my estimate be off by X amount?”
 - Learning global representations from *local information*
 - Complexity / Feasibility
 - Proving *asymptotic behavior of functions* (e.g. Big-O)

Topological Data Analysis (TDA)

- TDA comprises “a collection of powerful tools that can quantify shape and structure in data in order to answer questions from the data’s domain.” [Munch]
- Is an *emerging* field for data analysis!
- Motivation for TDA:
 - Data is **huge**, often **high dimensional**, and **complex**
 - Traditional techniques have not “kept up”
 - i.e. Rely on *overly-simplistic* assumptions
- Basic Idea:
 - **Data has shape**
 - This shape can be rigorously quantified via **topological signatures**

Much of this is summarized from: Munch, Elizabeth. "A user's guide to topological data analysis." Journal of Learning Analytics 4.2 (2017): 47-61.

Whats a Topological Signature?

- Informally, a topology on a set is a *description of how elements in the set are spatially related*
 - Can be seen as a **formalization of clustering**
 - i.e. *the collection of all open sets in a space is called its topology*
- A topological signature is a **simplified representation** of the topology of a given space
- Often, a [discrete] representation used as a **topological signature** is a **simplicial complex**
 - * 0-simplex == vertex
 - * 1-simplex == edge
 - * 2-simplex == triangle
 - * 3-simplex == tetrahedron
 - * ... k -simplex == ...



Example of a Topological Signature

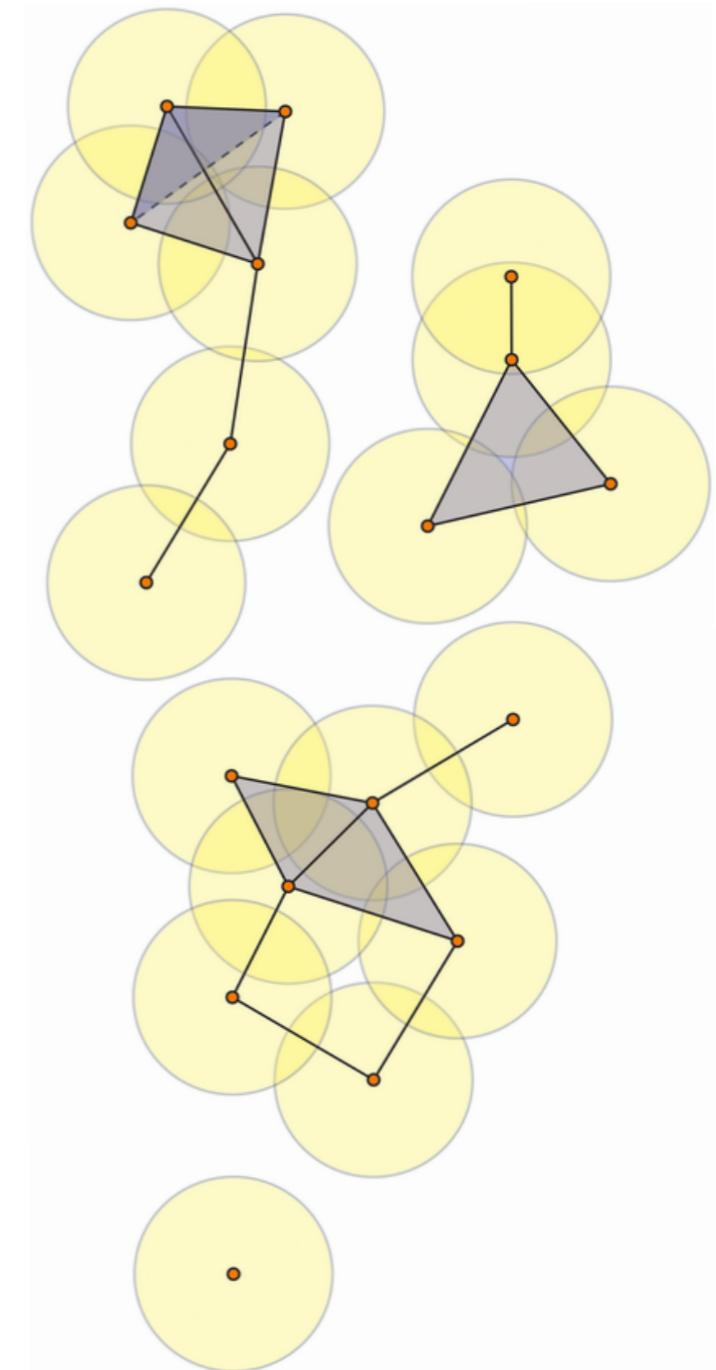
- Perhaps the most common topological signature is the so-called Rips Complex $VR_\epsilon(X)$ formed by forming an simplex between all points which have pairwise distances less than ϵ

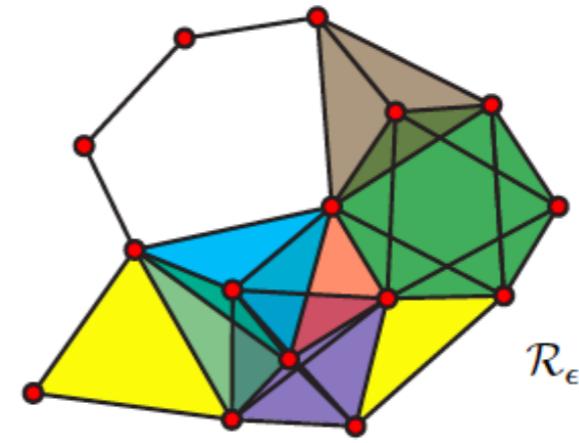
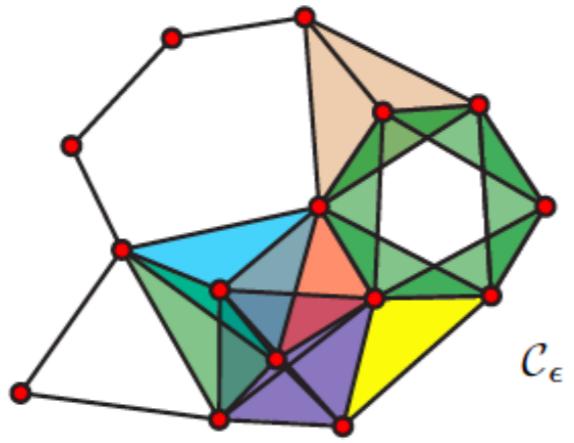
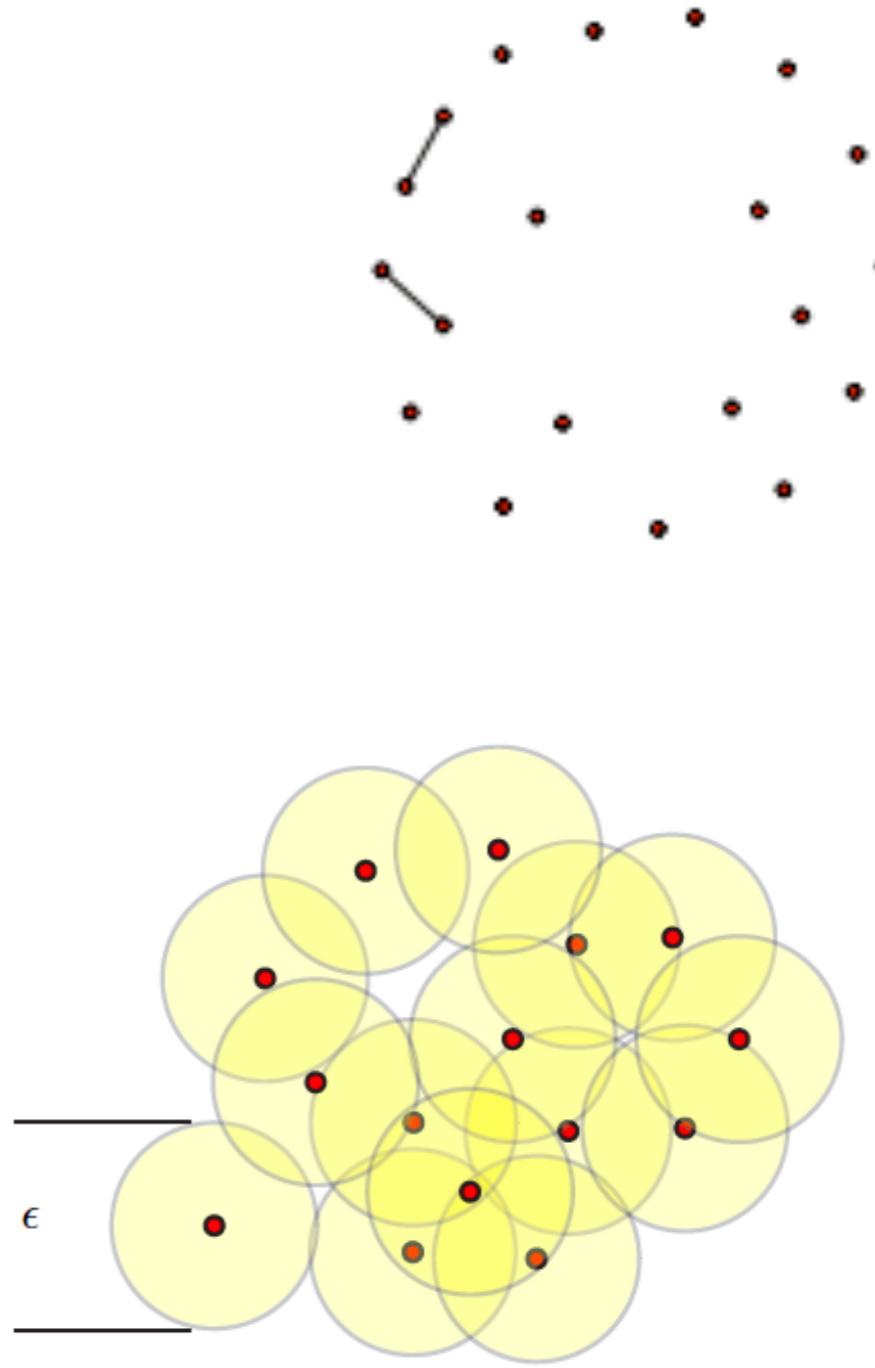
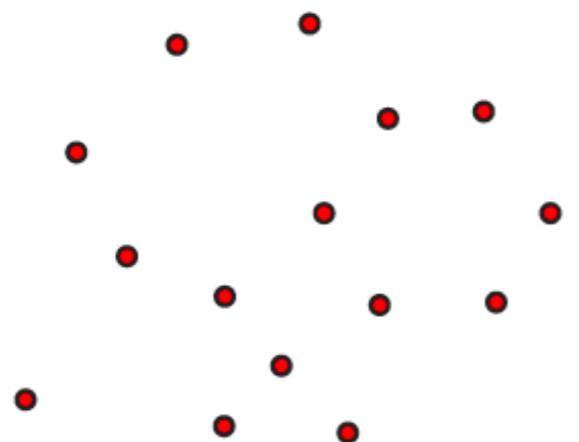
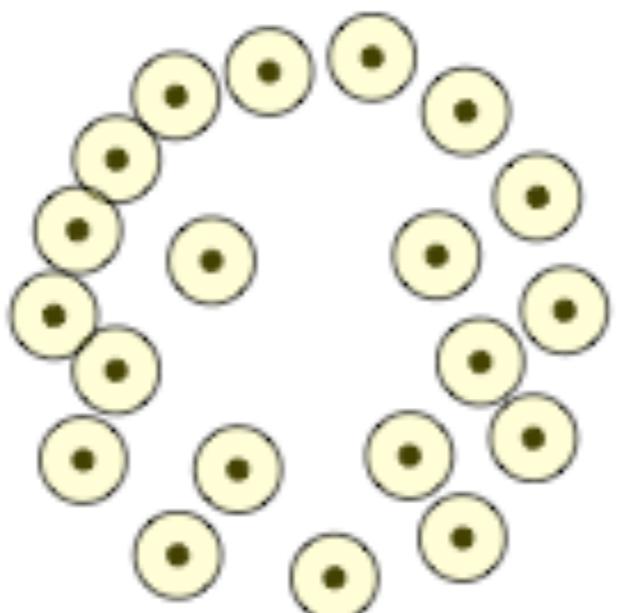
- The simplicial complex formed by non-empty intersections of

$$B(X_1, \epsilon) \cap B(X_2, \epsilon) \cap \dots \cap B(X_n, \epsilon)$$

where

$$B(X_i, \epsilon) = \{x \in X \mid d(X_i, x) \leq \epsilon\}$$





Recall: The Goal of TDA

- Basic Idea:
 - **Data has shape**
 - This shape can be **rigorously quantified** with **topology**
 - Quantify shape through **topological signatures**
 - Such signatures act as **summaries of the data**
- The Goal of TDA:
 - Use tools from topology to make meaningful signatures of the data
 - Topological signatures lead to topological **invariants**, and such invariants **enable greater understanding of the relationships in—and transformations of—real data**

Much of this is gleaned from either:

1. Munch, Elizabeth. "A user's guide to topological data analysis." *Journal of Learning Analytics* 4.2 (2017): 47-61.
2. Ghrist, Robert W. *Elementary applied topology*. Seattle: Createspace, 2014.

Mapper: a Topological Signature

- Mapper: *Perhaps* the most used signature in modern TDA *applications*
- Created by Singh, Mémoli, and Carlsson [Mapper]
- Simplest interpretation:
 - Interprets any set of data in \mathbb{R}^d as “point cloud data”, turns data into a *simplified topological graph*
 - “Mapper takes as input both a possibly **high dimensional dataset** and a **map defined on the data**, and produces a **summary of the data** by using a cover of the codomain of the map.”

▲	V1	▼	V2	▼	V2A	▼	V3	▼	V4	▼	V5	▼	V6	▼	V7	▼	V8	▼	V9	▼	V10	▼	V11	▼	V12	▼	V13	▼	V14	▼	V15	▼	V16	▼	V17	▼	V18	▼	V19	▼	V20	▼	V21	▼	V22	▼	V23	▼	V24	▼	V25	▼	V26	▼	V27	▼	V28	▼	V29	▼	V30	▼	V31	▼	V32	▼	V33	▼	V34	▼	V35	▼	V36	▼	V37	▼	V38	▼	V39	▼	V40	▼	V41	▼	V42	▼	V43	▼	V44	▼	V45	▼	V46	▼	V47	▼	V48	▼	V49	▼	V50	▼	V51	▼	V52	▼	V53	▼	V54	▼	V55	▼	V56	▼	V57	▼	V58	▼	V59	▼	V60	▼	V61	▼	V62	▼	V63	▼	V64	▼	V65	▼	V66	▼	V67	▼	V68	▼	V69	▼	V70	▼	V71	▼	V72	▼	V73	▼	V74	▼	V75	▼	V76	▼	V77	▼	V78	▼	V79	▼	V80	▼	V81	▼	V82	▼	V83	▼	V84	▼	V85	▼	V86	▼	V87	▼	V88	▼	V89	▼	V90	▼	V91	▼	V92	▼	V93	▼	V94	▼	V95	▼	V96	▼	V97	▼	V98	▼	V99	▼	V100	▼	V101	▼	V102	▼	V103	▼	V104	▼	V105	▼	V106	▼	V107	▼	V108	▼	V109	▼	V110	▼	V111	▼	V112	▼	V113	▼	V114	▼	V115	▼	V116	▼	V117	▼	V118	▼	V119	▼	V120	▼	V121	▼	V122	▼	V123	▼	V124	▼	V125	▼	V126	▼	V127	▼	V128	▼	V129	▼	V130	▼	V131	▼	V132	▼	V133	▼	V134	▼	V135	▼	V136	▼	V137	▼	V138	▼	V139	▼	V140	▼	V141	▼	V142	▼	V143	▼	V144	▼	V145	▼	V146	▼	V147	▼	V148	▼	V149	▼	V150	▼	V151	▼	V152	▼	V153	▼	V154	▼	V155	▼	V156	▼	V157	▼	V158	▼	V159	▼	V160	▼	V161	▼	V162	▼	V163	▼	V164	▼	V165	▼	V166	▼	V167	▼	V168	▼	V169	▼	V170	▼	V171	▼	V172	▼	V173	▼	V174	▼	V175	▼	V176	▼	V177	▼	V178	▼	V179	▼	V180	▼	V181	▼	V182	▼	V183	▼	V184	▼	V185	▼	V186	▼	V187	▼	V188	▼	V189	▼	V190	▼	V191	▼	V192	▼	V193	▼	V194	▼	V195	▼	V196	▼	V197	▼	V198	▼	V199	▼	V200	▼	V201	▼	V202	▼	V203	▼	V204	▼	V205	▼	V206	▼	V207	▼	V208	▼	V209	▼	V210	▼	V211	▼	V212	▼	V213	▼	V214	▼	V215	▼	V216	▼	V217	▼	V218	▼	V219	▼	V220	▼	V221	▼	V222	▼	V223	▼	V224	▼	V225	▼	V226	▼	V227	▼	V228	▼	V229	▼	V230	▼	V231	▼	V232	▼	V233	▼	V234	▼	V235	▼	V236	▼	V237	▼	V238	▼	V239	▼	V240	▼	V241	▼	V242	▼	V243	▼	V244	▼	V245	▼	V246	▼	V247	▼	V248	▼	V249	▼	V250	▼	V251	▼	V252	▼	V253	▼	V254	▼	V255	▼	V256	▼	V257	▼	V258	▼	V259	▼	V260	▼	V261	▼	V262	▼	V263	▼	V264	▼	V265	▼	V266	▼	V267	▼	V268	▼	V269	▼	V270	▼	V271	▼	V272	▼	V273	▼	V274	▼	V275	▼	V276	▼	V277	▼	V278	▼	V279	▼	V280	▼	V281	▼	V282	▼	V283	▼	V284	▼	V285	▼	V286	▼	V287	▼	V288	▼	V289	▼	V290	▼	V291	▼	V292	▼	V293	▼	V294	▼	V295	▼	V296	▼	V297	▼	V298	▼	V299	▼	V300	▼	V301	▼	V302	▼	V303	▼	V304	▼	V305	▼	V306	▼	V307	▼	V308	▼	V309	▼	V310	▼	V311	▼	V312	▼	V313	▼	V314	▼	V315	▼	V316	▼	V317	▼	V318	▼	V319	▼	V320	▼	V321	▼	V322	▼	V323	▼	V324	▼	V325	▼	V326	▼	V327	▼	V328	▼	V329	▼	V330	▼	V331	▼	V332	▼	V333	▼	V334	▼	V335	▼	V336	▼	V337	▼	V338	▼	V339	▼	V340	▼	V341	▼	V342	▼	V343	▼	V344	▼	V345	▼	V346	▼	V347	▼	V348	▼	V349	▼	V350	▼	V351	▼	V352	▼	V353	▼	V354	▼	V355	▼	V356	▼	V357	▼	V358	▼	V359	▼	V360	▼	V361	▼	V362	▼	V363	▼	V364	▼	V365	▼	V366	▼	V367	▼	V368	▼	V369	▼	V370	▼	V371	▼	V372	▼	V373	▼	V374	▼	V375	▼	V376	▼	V377	▼	V378	▼	V379	▼	V380	▼	V381	▼	V382	▼	V383	▼	V384	▼	V385	▼	V386	▼	V387	▼	V388	▼	V389	▼	V390	▼	V391	▼	V392	▼	V393	▼	V394	▼	V395	▼	V396	▼	V397	▼	V398	▼	V399	▼	V400	▼	V401	▼	V402	▼	V403	▼	V404	▼	V405	▼	V406	▼	V407	▼	V408	▼	V409	▼	V410	▼	V411	▼	V412	▼	V413	▼	V414	▼	V415	▼	V416	▼	V417	▼	V418	▼	V419	▼	V420	▼	V421	▼	V422	▼	V423	▼	V424	▼	V425	▼	V426	▼	V427	▼	V428	▼	V429	▼	V430	▼	V431	▼	V432	▼	V433	▼	V434	▼	V435	▼	V436	▼	V437	▼	V438	▼	V439	▼	V440	▼	V441	▼	V442	▼	V443	▼	V444	▼	V445	▼	V446	▼	V447	▼	V448	▼	V449	▼	V450	▼	V451	▼	V452	▼	V453	▼	V454	▼	V455	▼	V456	▼	V457	▼	V458	▼	V459	▼	V460	▼	V461	▼	V462	▼	V463	▼	V464	▼	V465	▼	V466	▼	V467	▼	V468	▼	V469	▼	V470	▼	V471	▼	V472	▼	V473	▼	V474	▼	V475	▼	V476	▼	V477	▼	V478	▼	V479	▼	V480	▼	V481	▼	V482	▼	V483	▼	V484	▼	V485	▼	V486	▼	V487	▼	V488	▼	V489	▼	V490	▼	V491	▼	V492	▼	V493	▼	V494	▼	V495	▼	V496	▼	V497	▼	V498	▼	V499	▼	V500	▼	V501	▼	V502	▼	V503	▼	V504	▼	V505	▼	V506	▼	V507	▼	V508	▼	V509	▼	V510	▼	V511	▼	V512	▼	V513	▼	V514	▼	V515	▼	V516	▼	V517	▼	V518	▼	V519	▼	V520	▼	V521	▼	V522	▼	V523	▼	V524	▼	V525	▼	V526	▼	V527	▼	V528	▼	V529	▼	V530	▼	V531	▼	V532	▼	V533	▼	V534	▼	V535	▼	V536	▼	V537	▼	V538	▼	V539	▼	V540	▼	V541	▼	V542	▼	V543	▼	V544	▼	V545	▼	V546	▼	V547	▼	V548	▼	V549	▼	V550	▼	V551	▼	V552	▼	V553	▼	V554	▼	V555	▼	V556	▼	V557	▼	V558	▼	V559	▼	V560	▼	V561	▼	V562	▼	V563	▼	V564	▼	V565	▼	V566	▼	V567	▼	V568	▼	V569	▼	V570	▼	V571	▼	V572	▼	V573	▼	V574	▼	V575	▼	V576	▼	V577	▼	V578	▼	V579	▼	V580	▼	V581	▼	V582	▼	V583	▼	V584	▼	V585	▼	V586	▼	V587	▼	V588	▼	V589	▼	V590	▼	V591	▼	V592	▼	V593	▼	V594	▼	V595	▼	V596	▼	V597	▼	V598	▼	V599	▼	V600	

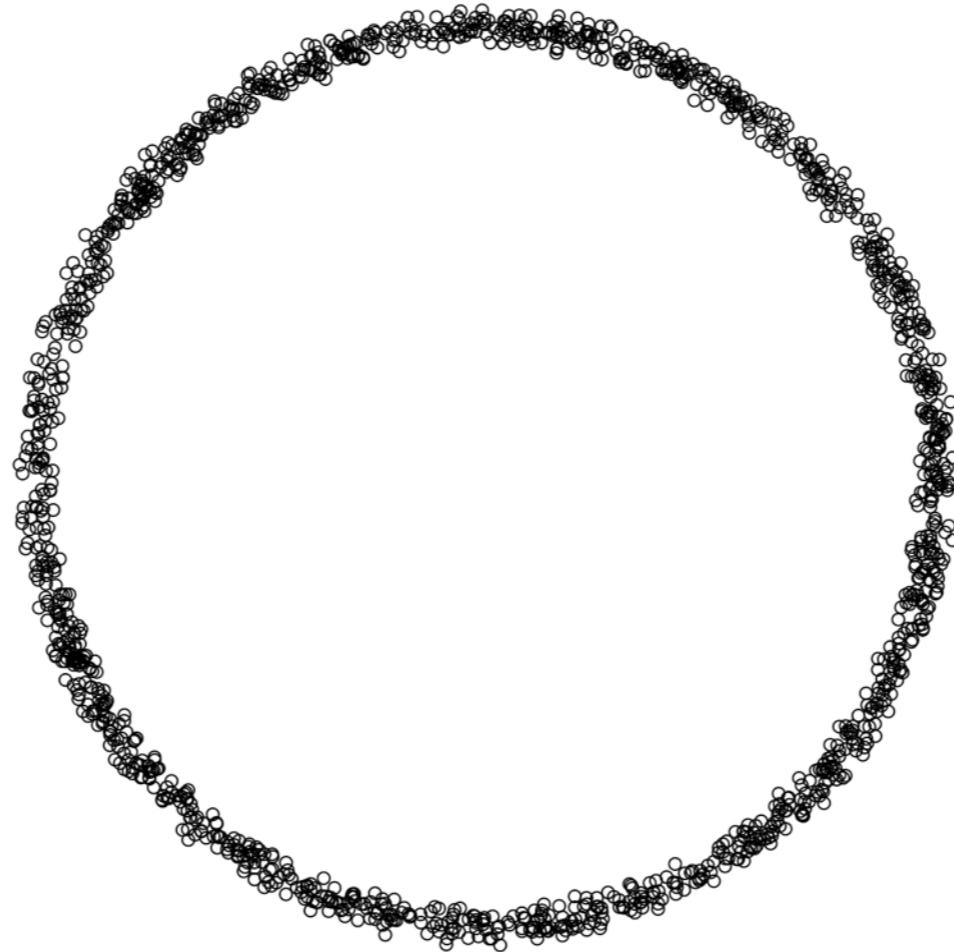
Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** \mathcal{C} to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**

$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

Mapper: Example

Consider the case where your
data looks like a circle



$$X = \{x_1, x_2, \dots, x_n\}$$

Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** \mathcal{C} to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**

$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

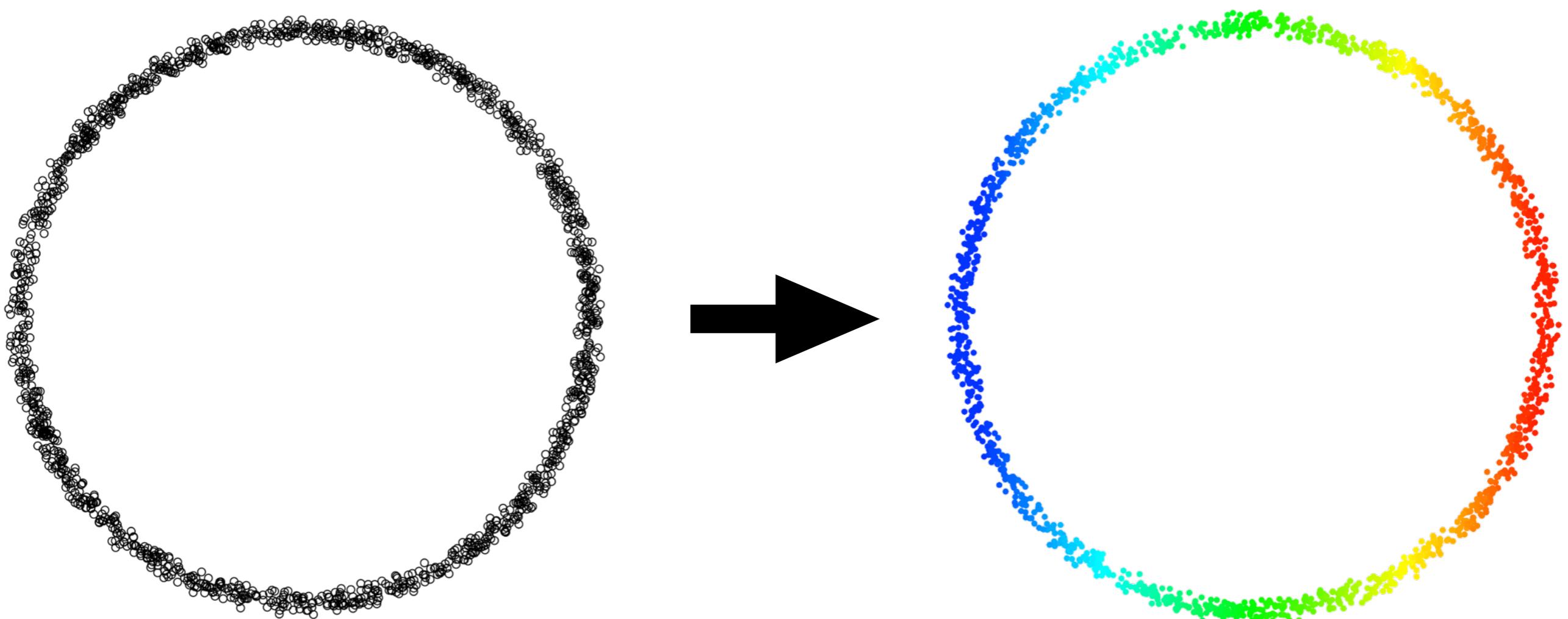
Mapper: Example — Step 1

... and you evaluate a function f

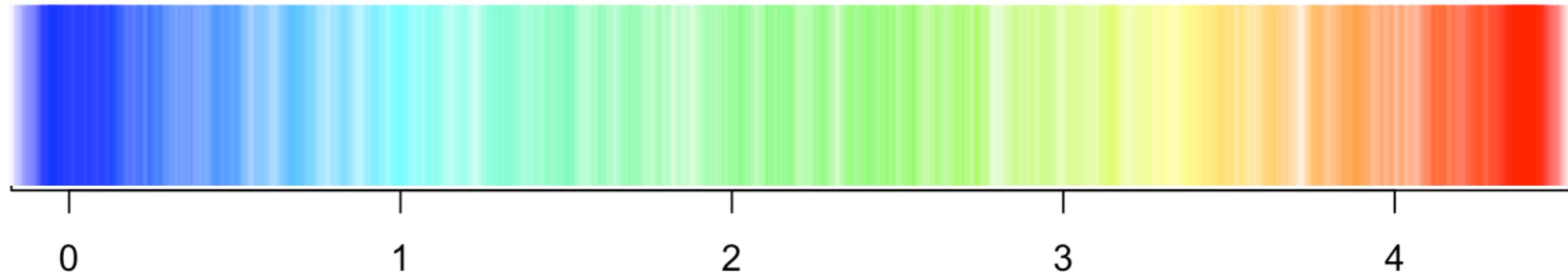
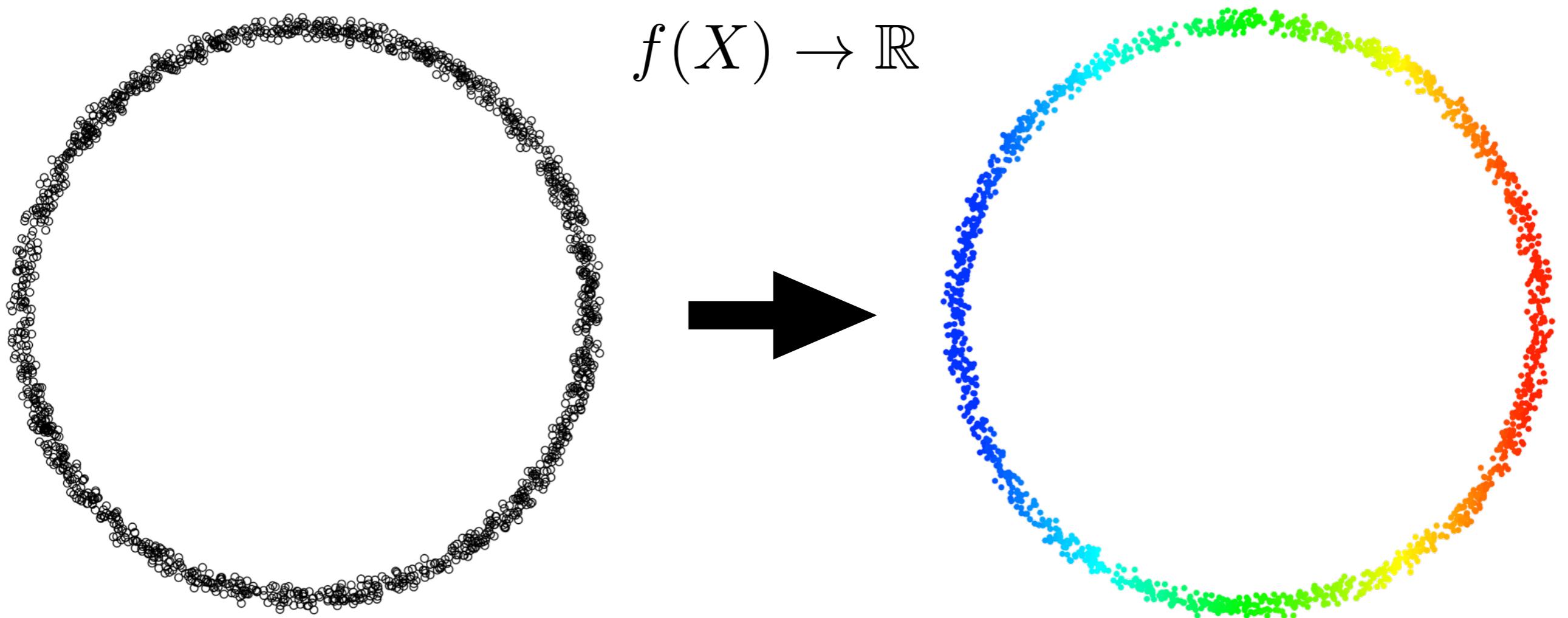
on that circle $f(x) = \|x - p\|_2$

where p is the left-most point in the circle

(blue == low distance, red == high distance)



Mapper: Example — Step 1



Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** \mathcal{C} to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**
$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

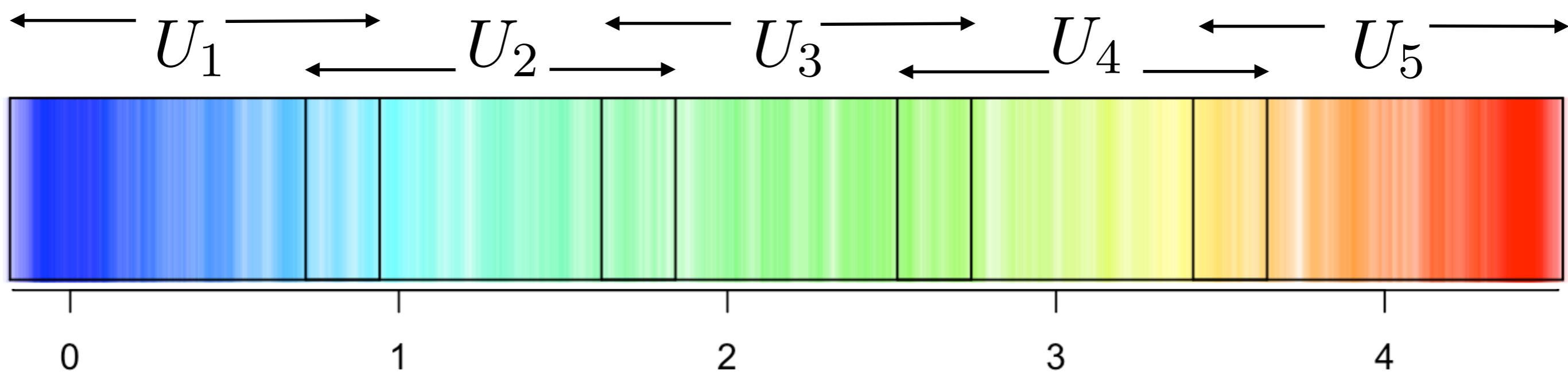
Mapper: Example — Step 2

Define a cover of the filter space Z

Recall the definition of a cover is a collection of sets whose union contains some space as a subset

$$C = \{U_\alpha : a \in A\} \quad X \subseteq \bigcup_{\alpha \in A} U_\alpha$$

In the case above, C covers X

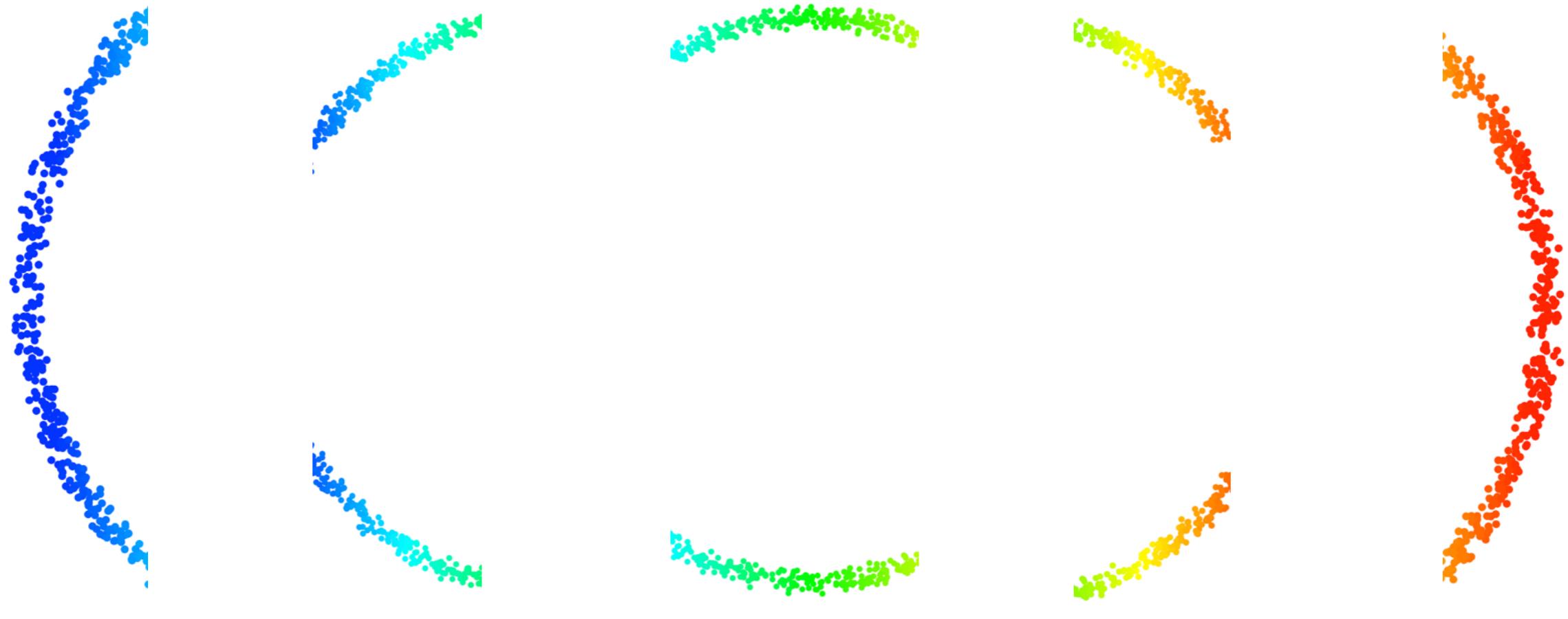


Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** C to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**

$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

Mapper: Example — Step 3



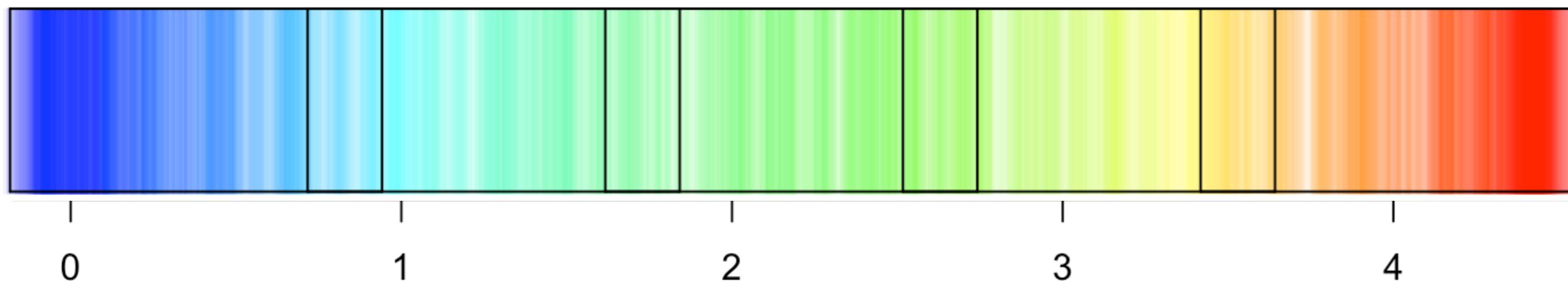
U_1

U_2

U_3

U_4

U_5

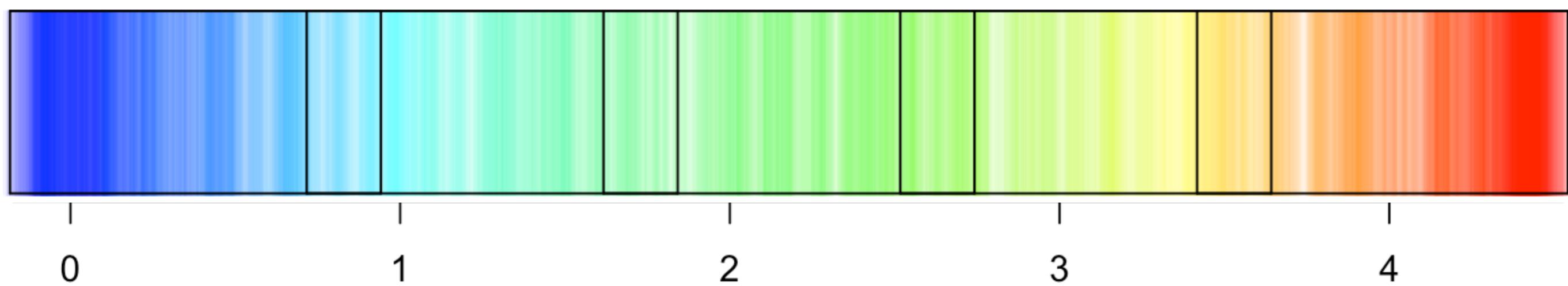
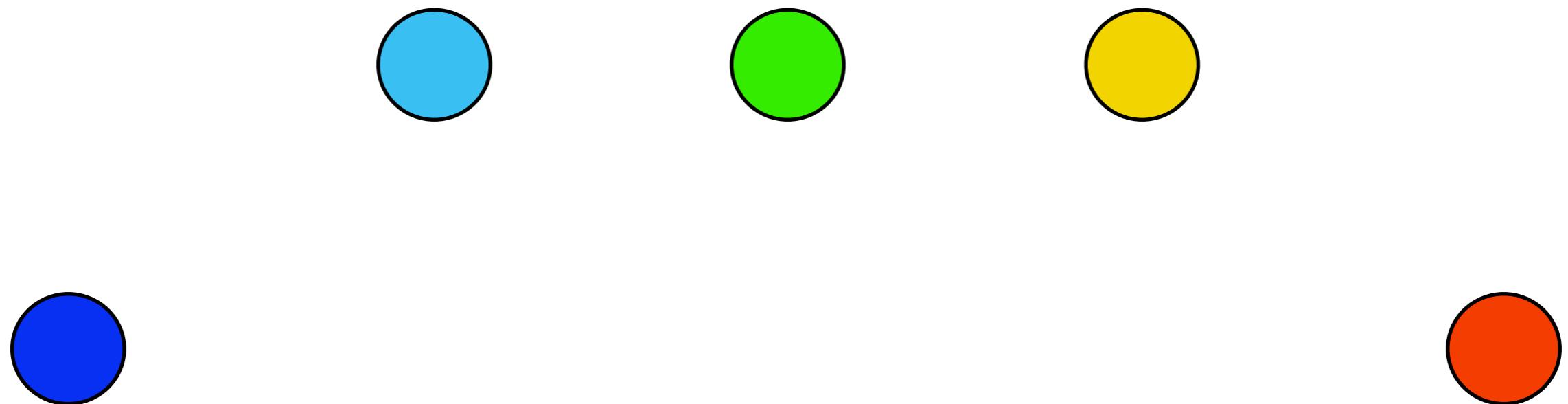


Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** \mathcal{C} to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X_α by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**

$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

Mapper: Example — Step 4

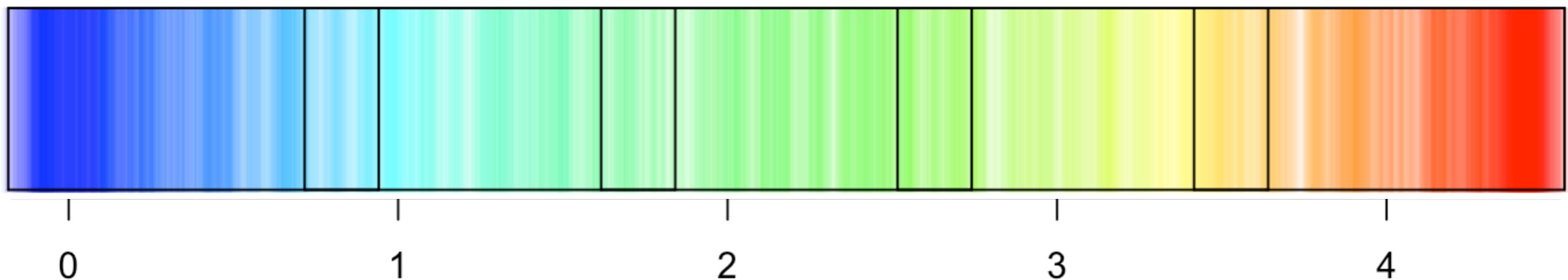
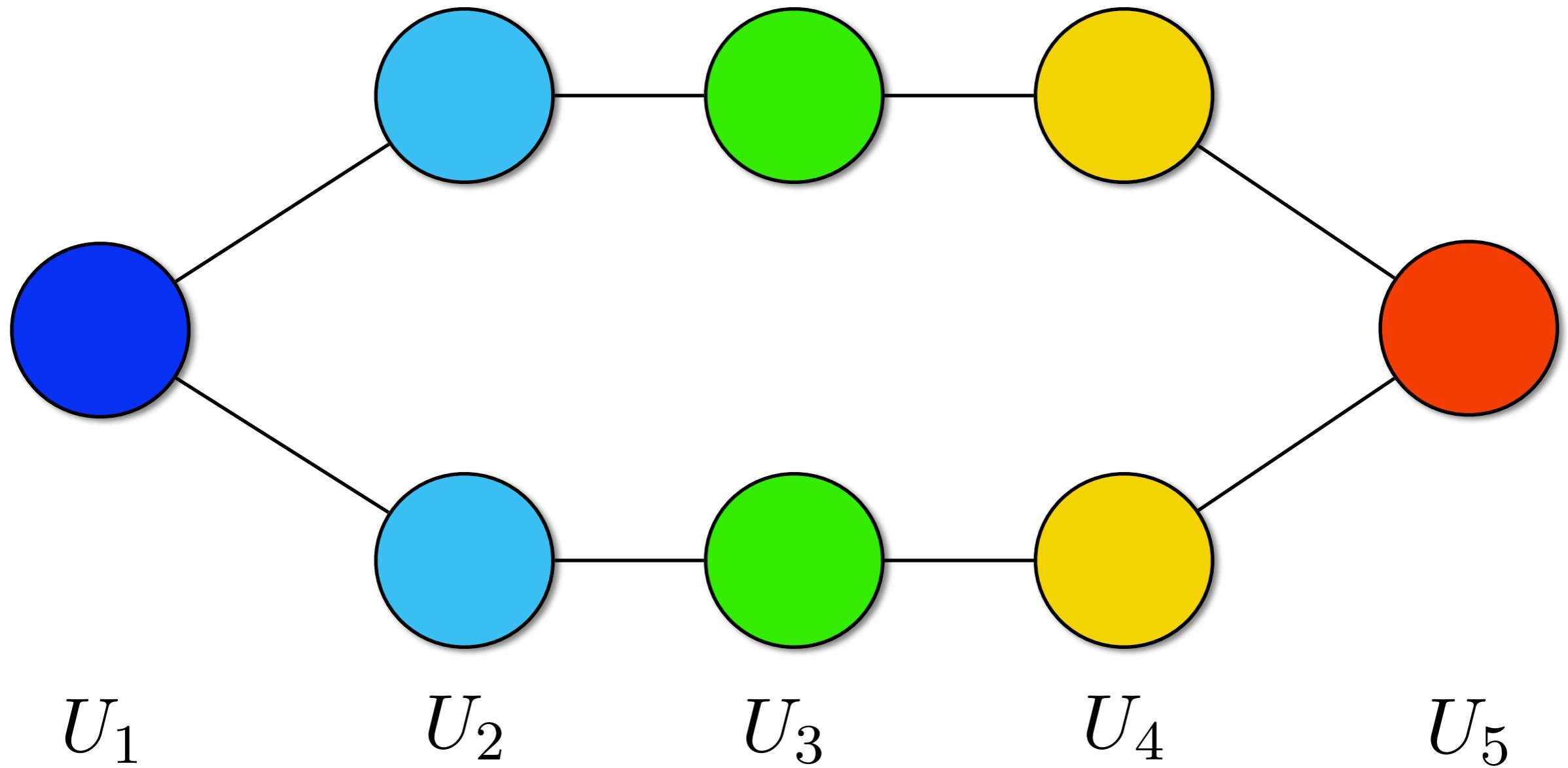


Mapper: Background

1. Define a **reference map** $f : X \rightarrow Z$
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ of Z
 - A is called **the index set**
3. Construct the subsets X_α via $f^{-1}(U_\alpha)$
4. Apply a **clustering algorithm** \mathcal{C} to the sets X_α
5. Obtain a cover $f^*(\mathcal{U})$ of X_α by considering the path-connected components of X
 - Clusters form “**nodes**” / 0-simplexes
 - Non-empty intersections form “**edges**” / 1-simplexes
6. The Mapper construction is **the nerve of this cover**

$$M(\mathcal{U}, f) := N(f^*(\mathcal{U}))$$

Mapper: Example — Step 5-6

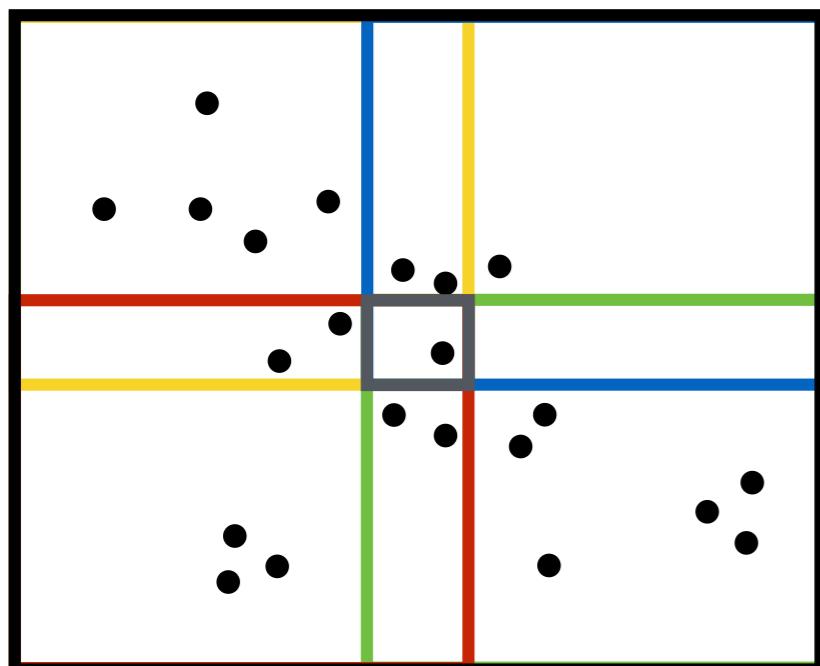


Mapper: In one picture

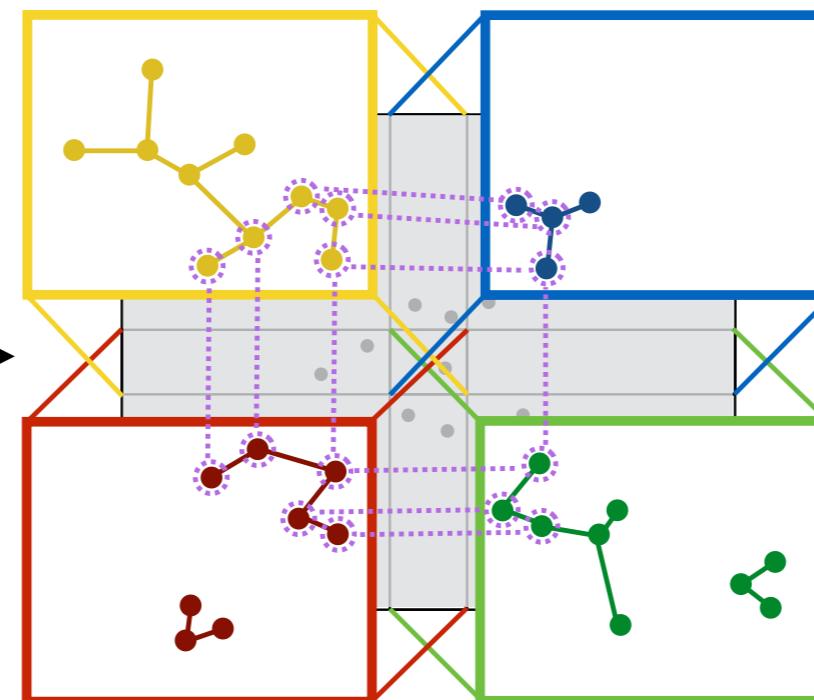
$$f : X \rightarrow Z$$

$$X_\alpha = f^{-1}U_\alpha$$

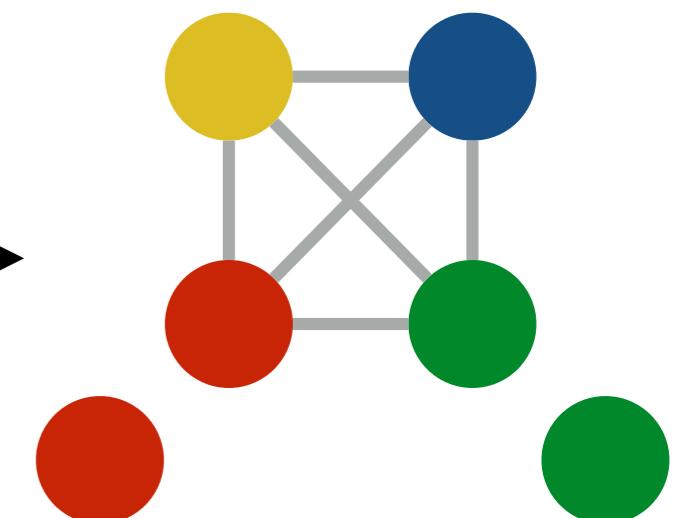
$$f^*(\mathcal{U}) \rightarrow M(\mathcal{U}, f)$$



$$\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$$



$$C(X_\alpha) \rightarrow f^*(\mathcal{U})$$



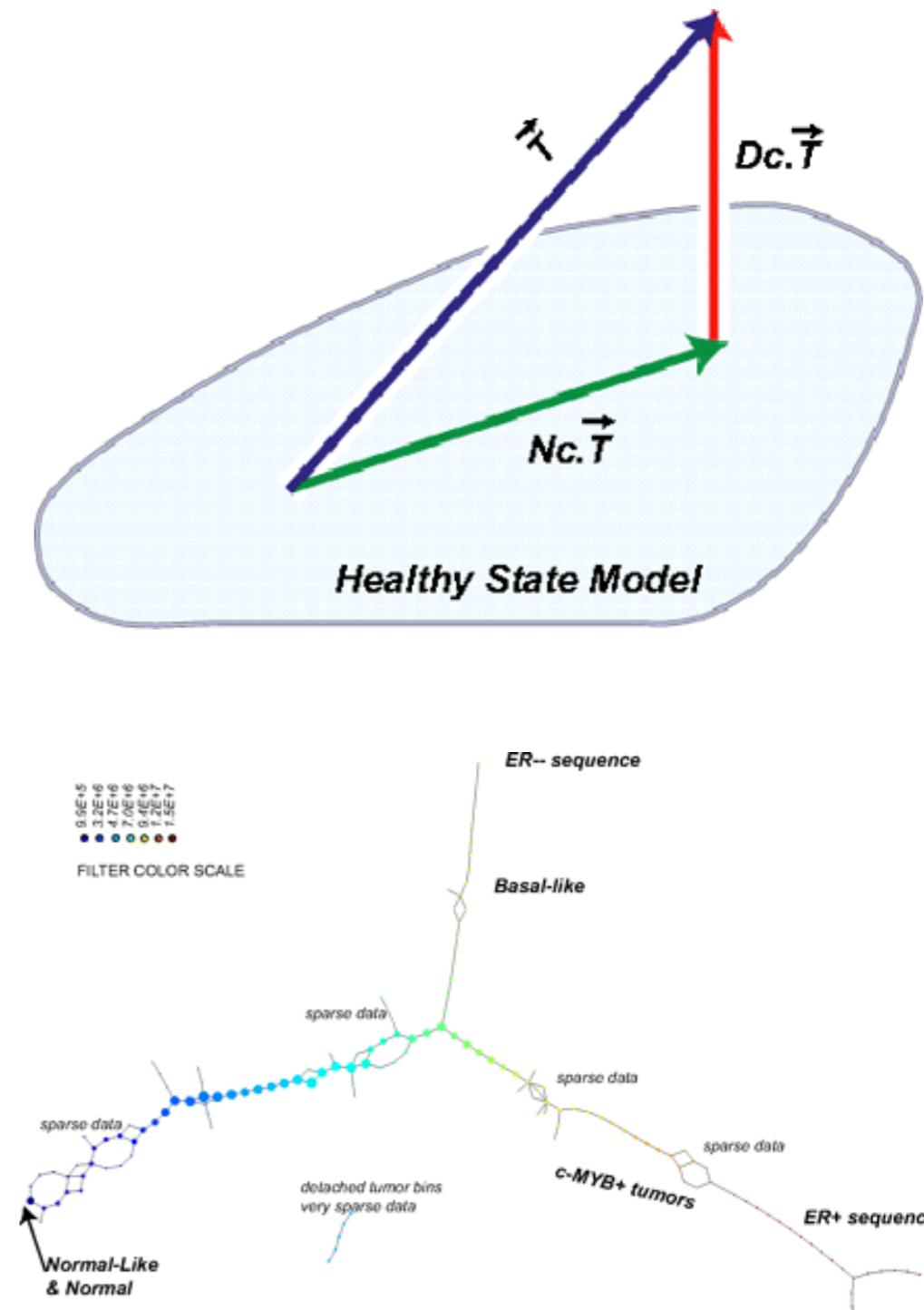
(1-skeleton example)

Why Mapper is useful

- Mapper provides a **succinct summary** of **the shape of a data set** (expressed via the codomain of the mapping)
- Mappers utility lies in its **generality**:
 - *Any mapping function* can be used
 - **Cover** may be constructed *arbitrarily*
 - *Any clustering algorithm* may be used
- The resulting graph is often much easier to interpret than, e.g. individual scatter plots of pairwise relationships
- Mapper is often paired with **high-dimensional data**, and is generally used to see the ‘true’ shape or **structure of the data**
- Mapper is the core algorithm behind the AI Company, [Ayasdi Inc.](#)
 - Anti-Money Laundering
 - Detecting Payment Fraud
 - Assessing health risks

How to choose a mapping function?

- Mapper is **highly dependent on the choice of filter function**
- Ideally, a **domain-specific** map that is well-understood may be appropriate
 - Ex. Biologists created a Healthy State Model (HSM) which encodes **tumor cell tissues** into orthogonal **“disease”** and **“normal”** components
 - Using the disease component allows for disease-specific analysis of their data
 - What if we want a **general** filter function that can be used for many kinds of data sets?



Manifold Learning

- What is a Manifold?
 - A topological space that **locally resembles Euclidean space**
 - A manifold [generally] is “smooth” if it permits the use of partial differentiation*
 - Ex. of Manifolds: The Earth, a Torus, a swiss roll
- The Manifold Hypothesis —————
 - The manifold hypothesis is the idea that data tend to lie *on or near a low dimensional manifold*
 - Alt. Def: **The dimensionality of data is only arbitrarily high**; rather, data may exist in some “**intrinsic**” dimensionality

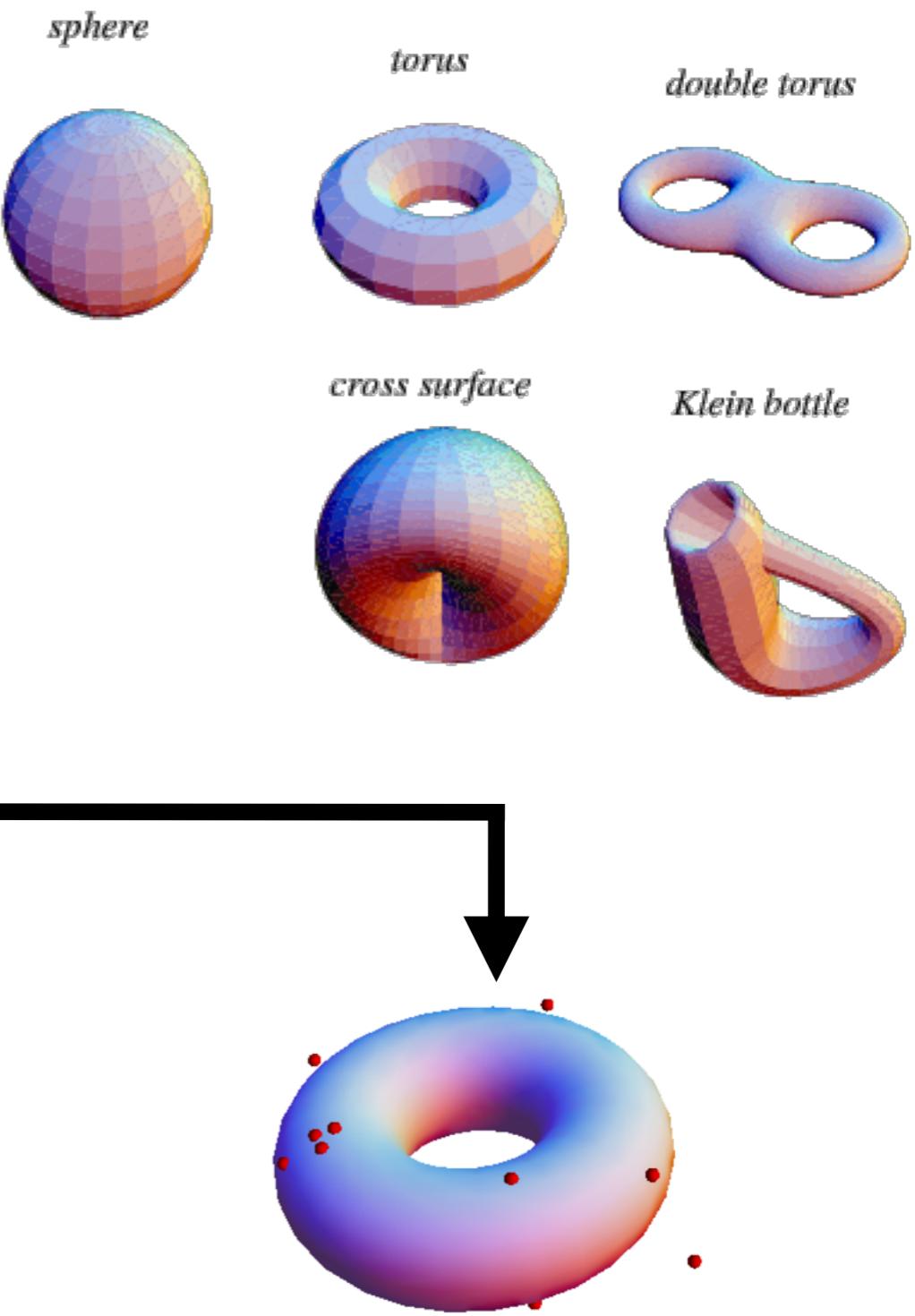


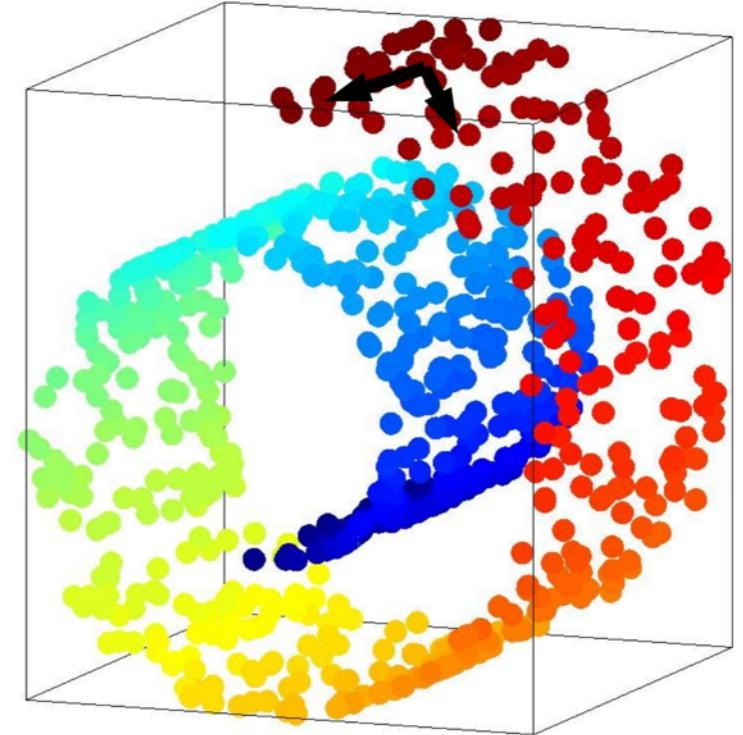
FIGURE 1. Data lying in the vicinity of a two dimensional torus.

*Technically, a Manifold is *smooth* when a Hausdorff space is furnished with an atlas of charts, each of which are infinitely differentiable whenever they intersect. Image from: Fefferman, Charles, Sanjoy Mitter, and Hariharan Narayanan. "Testing the manifold hypothesis." Journal of the American Mathematical Society 29.4 (2016): 983-1049.

IsoMap Example

1. Make a **neighborhood graph**,
connected points that are either:

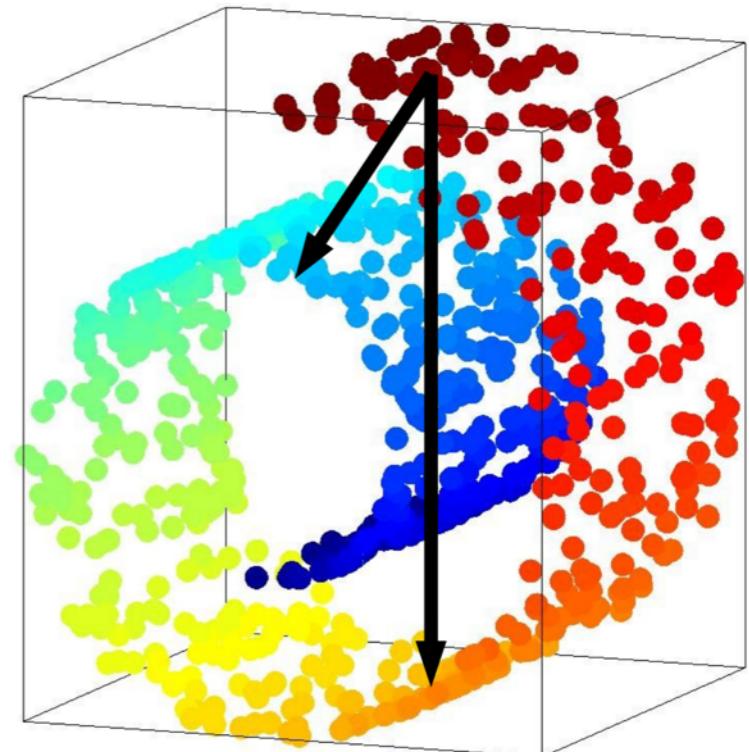
- Within some distance ϵ to each other
- Is the k -th nearest neighbor of another



2. Compute **Shortest path** (Dijkstra's)

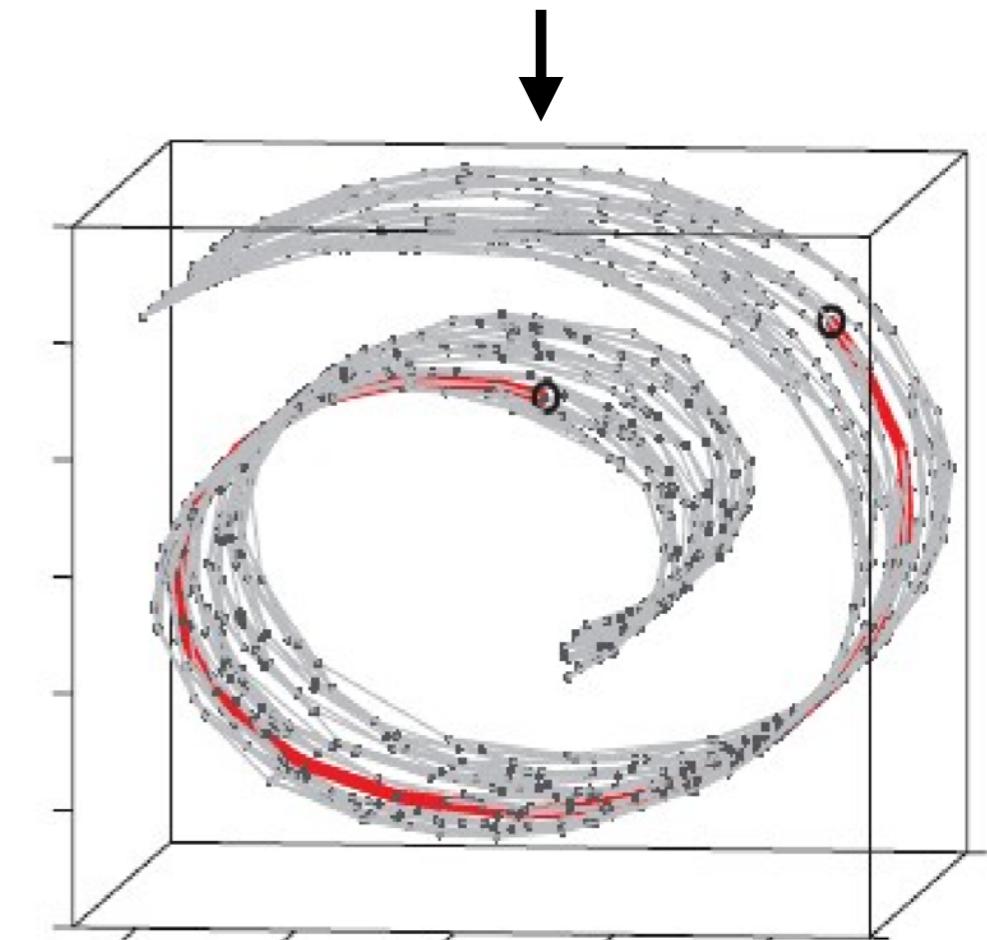
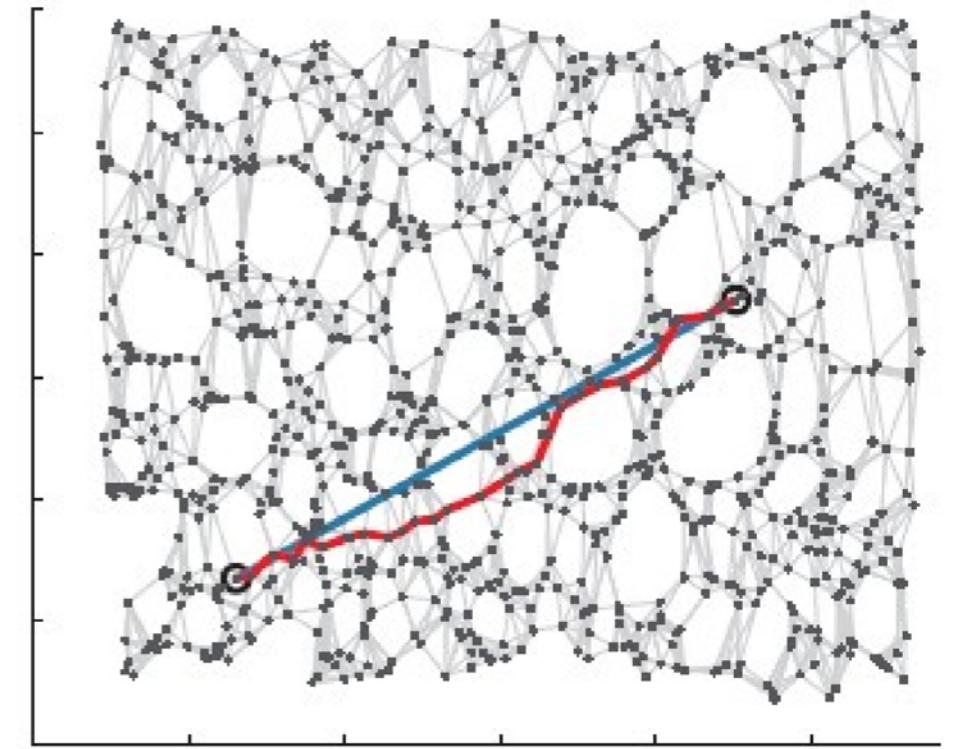
3. Compute **lower-dimensional embedding** (MDS)

- Goal: Given distance matrix $D = (d_{ij})$, compute the set:
 $x_1, \dots, x_n \in \mathbb{R}^m$
Such that: $d_{ij} \approx \|x_i - x_j\|_2$



IsoMap Example

1. Make a **neighborhood graph**, connected points that are either:
 - Within some distance ϵ to each other
 - Is the k -th nearest neighbor of another
2. Compute **Shortest path** (Dijkstra's)
3. Compute **lower-dimensional embedding** (MDS)
 - Goal: Given distance matrix $D = (d_{ij})$, compute the set:
$$x_1, \dots, x_n \in \mathbb{R}^m$$
Such that: $d_{ij} \approx \|x_i - x_j\|_2$

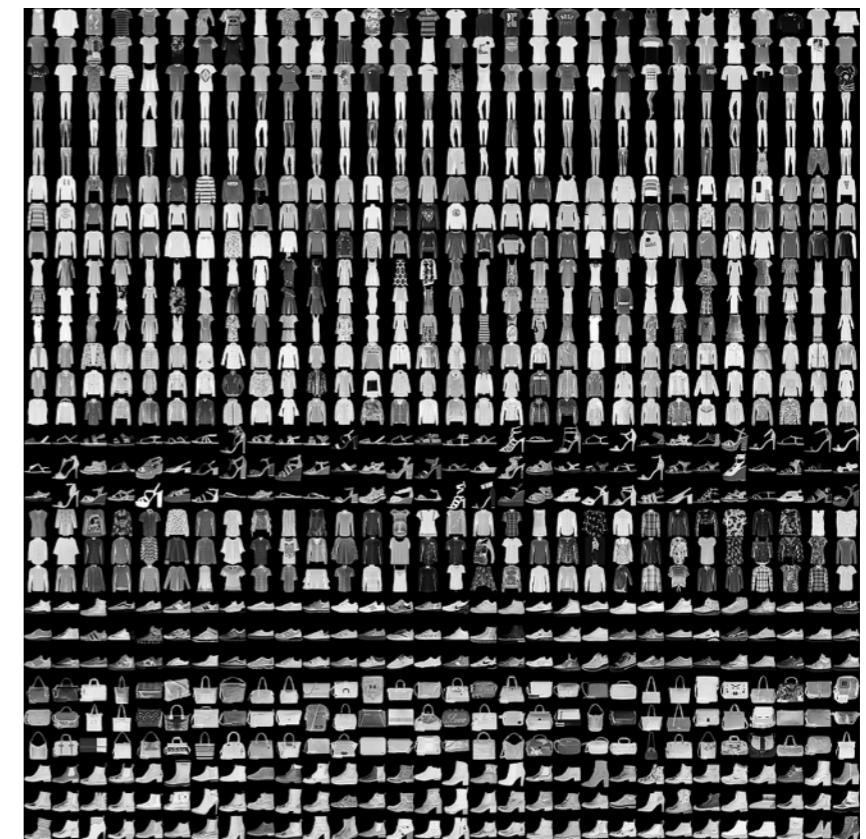


Why care about the Manifold Hypothesis?

- There's a lot of theoretical [3, 4], experimental [2], and empirical [1, 3] evidence supporting the Manifold Hypothesis
- Related: Non-linear Dimensionality Reduction
 - SNE, t-SNE, IsoMap, Locally Linear Embedding
- There's *several ways* to get a different perspective of the manifold; there is (to my knowledge) no “**best**” manifold approximation technique
- The hypothesis has actually motivated developments in **Generative Adversarial Networks** [3]



t-SNE
applied to
MNIST



1. <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/#fn3>
2. Carlsson, Gunnar, et al. "On the local behavior of spaces of natural images." *International journal of computer vision* 76.1 (2008): 1-12.
3. Lui, Kry Yik Chau, et al. "Implicit Manifold Learning on Generative Adversarial Networks." *arXiv preprint arXiv:1710.11260* (2017).
4. Fefferman, Charles, Sanjoy Mitter, and Hariharan Narayanan. "Testing the manifold hypothesis." *Journal of the American Mathematical Society* 29.4 (2016): 983-1049.

Background: Reeb Graphs

- Level set definition:

$$L_c(f) = \{(x_1, \dots, x_n) \mid f(x_1, \dots, x_n) = c\}$$

- A Reeb graph is a mathematical object reflecting **the evolution of the level sets** of a real-valued function **on a manifold**.

- Points are part of the same ‘edge’ if they belong in the same connected component in $f^{-1}(c)$

- Reeb space == multivariate generalization of Reeb graph

- “**...compresses the components of the level sets** of a multivariate and obtains a summary representation of their relationships”

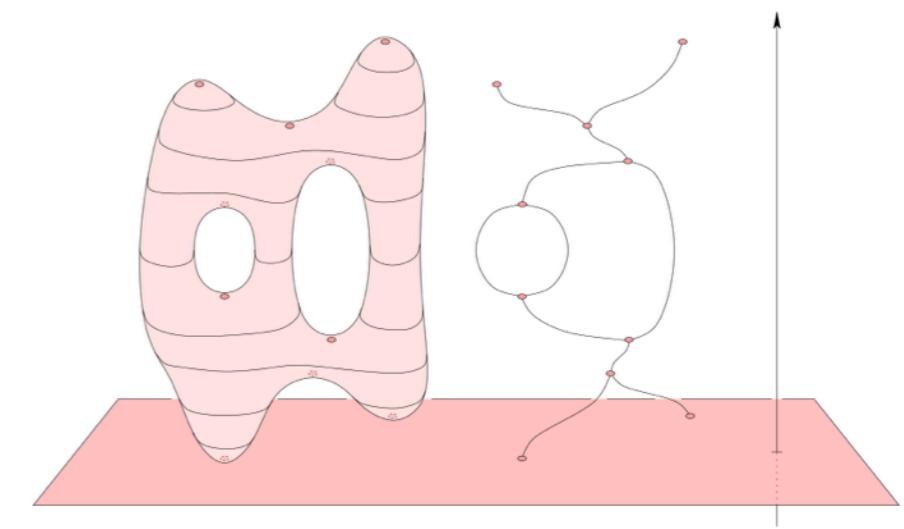
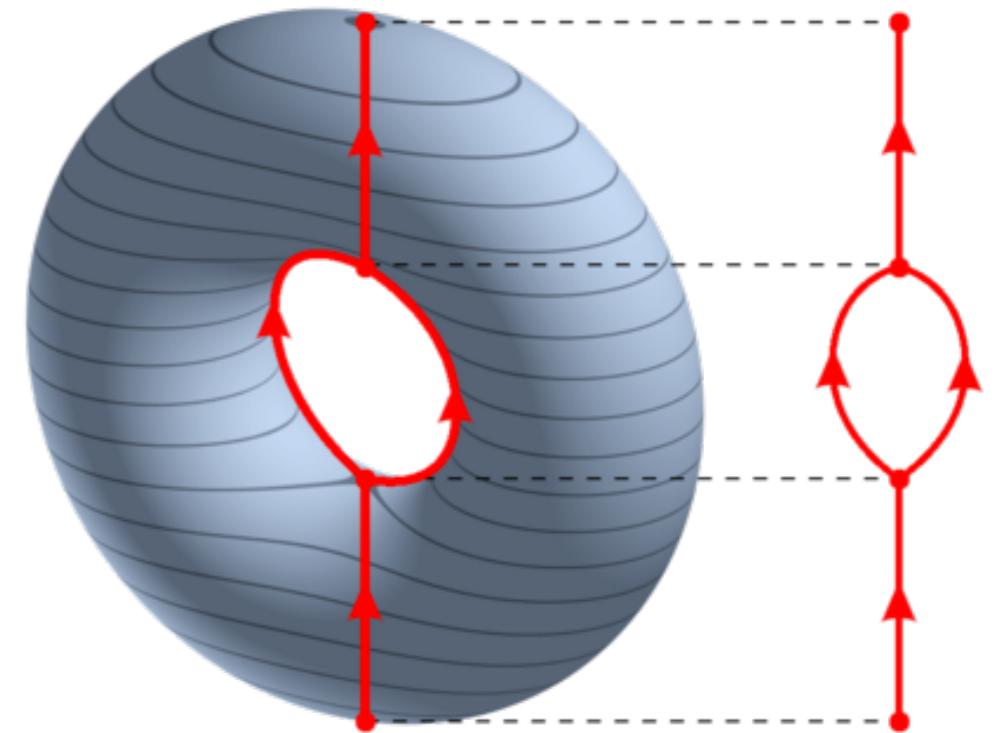


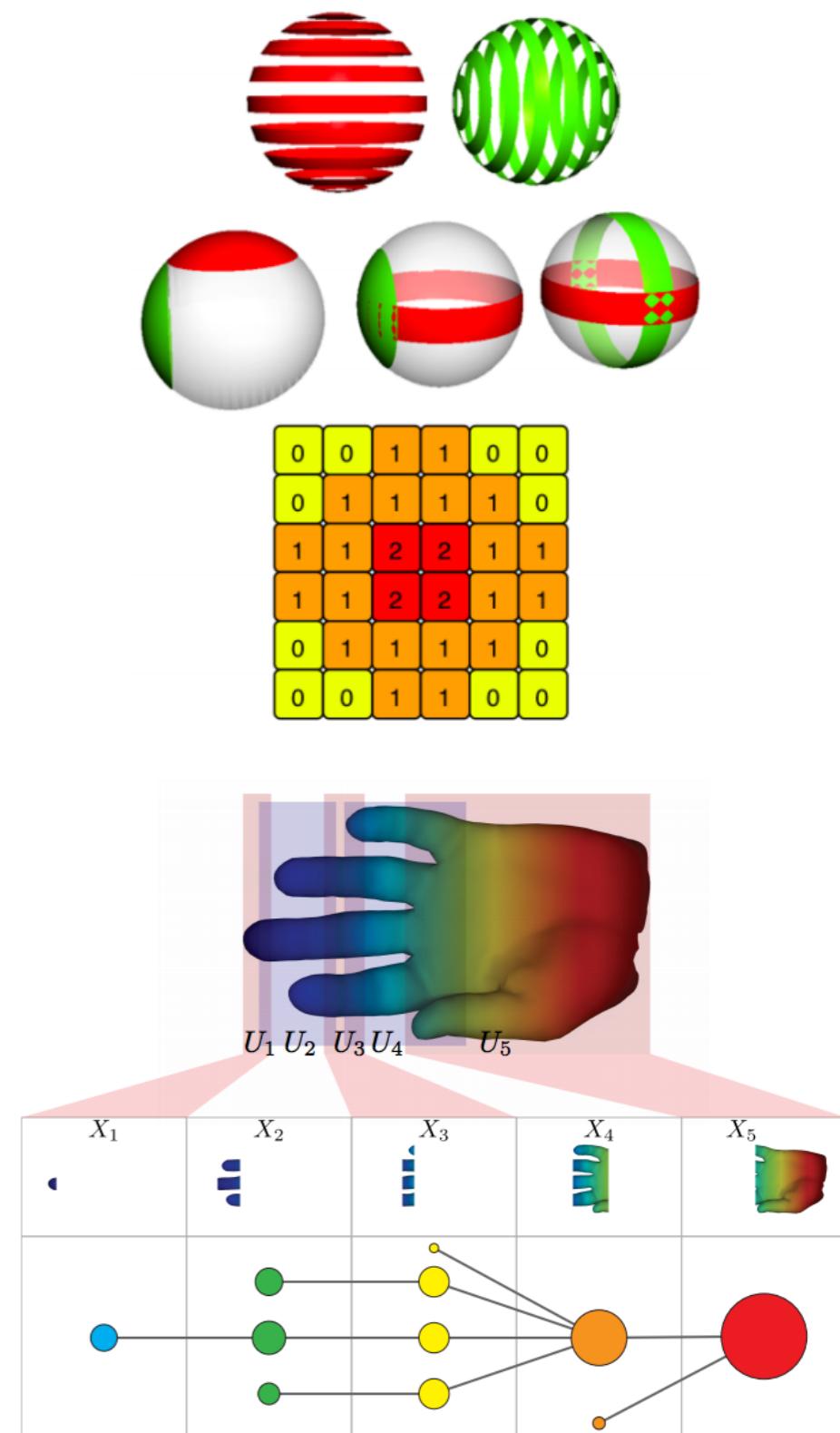
Figure V.13: Level sets of the 2-manifold map to points on the real line and components of the level sets map to points of the Reeb graph.

Top from: https://en.wikipedia.org/wiki/Reeb_graph

Bottom image from: <https://www2.cs.duke.edu/courses/fall06/cps296.1/Lectures/sec-V-4.pdf>

What does Mapper represent?

- What does mapper *actually* do?
 - "...our construction amounts to a stochastic version of the **Reeb graph** associated with the filter function. If the covering of R is too coarse, **we will be constructing an image of the Reeb graph of the function**, while if it is fine enough **we will recover the Reeb graph precisely.**"
 - Munch et. al [1] proved that the categorical representations of the Reeb space and Mapper **converge in terms of interleaving distance**
- So we may *think* of Mapper as an approximation of the Reeb space



1. Munch, Elizabeth, and Bei Wang. "Convergence between categorical representations of Reeb space and Mapper." arXiv preprint arXiv:1512.04108 (2015).

Top From: Singh, Gurjeet, Facundo Mémoli, and Gunnar E. Carlsson. "Topological methods for the analysis of high dimensional data sets and 3d object recognition." SPBG. 2007.

Bottom from: <http://web.cse.ohio-state.edu/~wang.1016/courses/5559/Lecs/mapper-lec5559.pdf>

Demo of World Values Survey



► LiveSlides web content

To view

Download the add-in.

liveslides.com/download

Start the presentation.

Time Permitting

Persistent Homology

- Recurring theme in applied topological data analysis
- “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.” - Ghrist

- “The motivation is that, for a parameterized family of spaces (i.e. VR complexes) modeling a point-cloud data set, **qualitative features which persist over a large parameter range have greater statistical significance**”

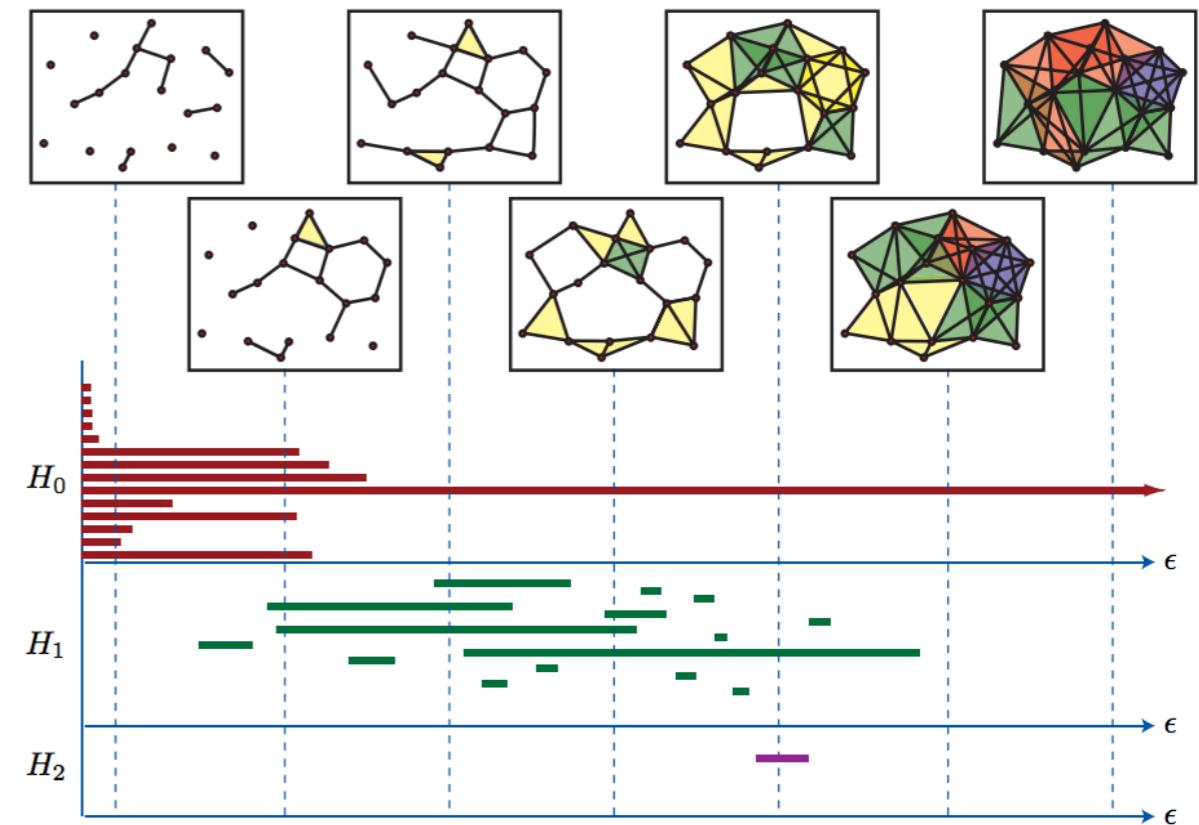


FIGURE 4. [bottom] An example of the barcodes for $H_*(R)$ in the example of Figure 3. [top] The rank of $H_k(\mathcal{R}_{\epsilon_i})$ equals the number of intervals in the barcode for $H_k(R)$ intersecting the (dashed) line $\epsilon = \epsilon_i$.

Studying Mapper: In the context of Persistent Homology

- Is it possible to study Mapper **in the context of Persistent Homology?**
- “The icon of persistence is a **monotone sequence**

$$\cdots \rightarrow \cdot \rightarrow \cdot \rightarrow \cdots$$

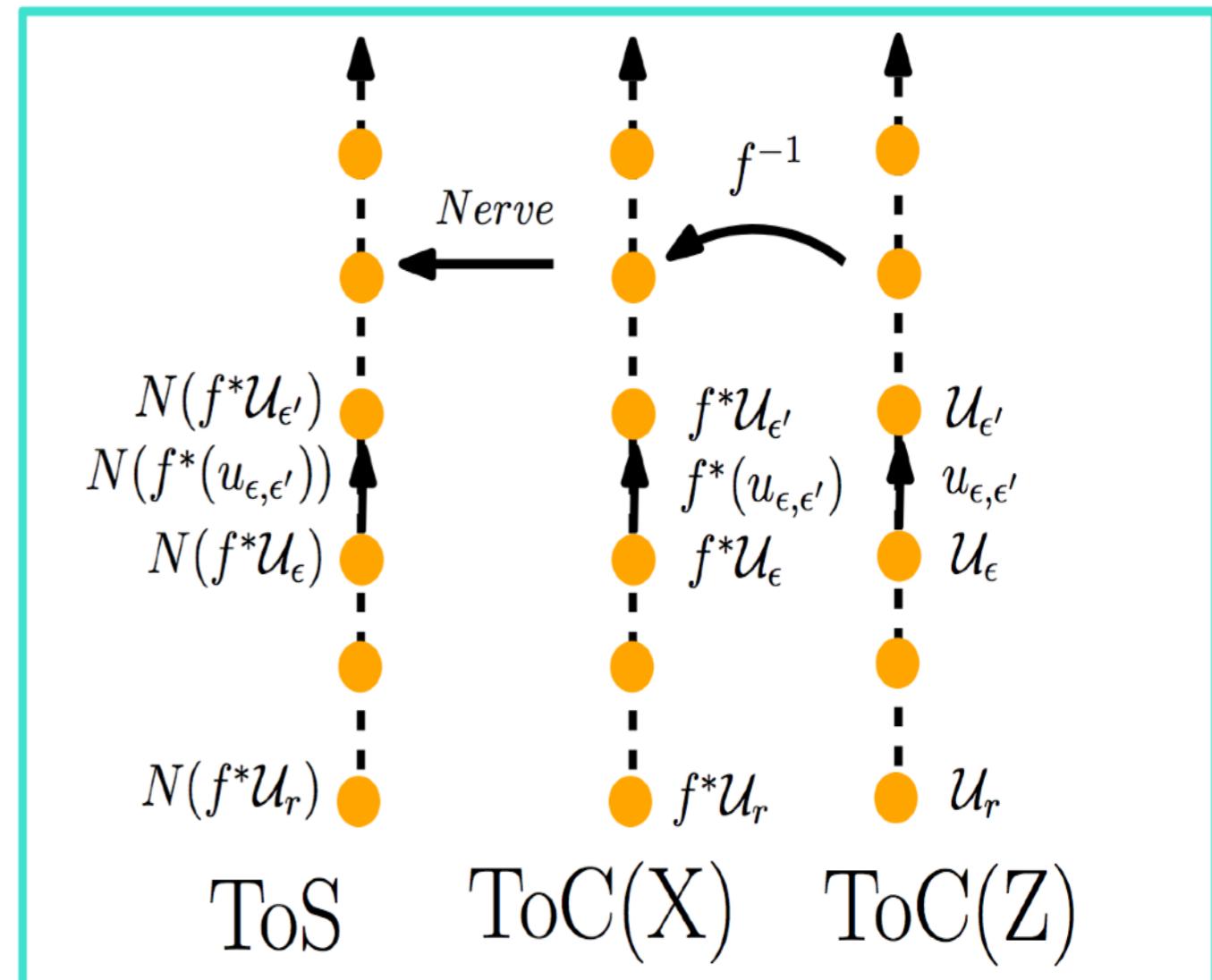
where arrows connote maps of spaces or **chains** or the **induced homomorphisms** on homology." - Ghrist

- So, we just need some kind of monotone sequence between Mapper constructions?
- As it turns out, someone* already has!
 - Published analysis of so-called “Multi-Scale Mapper” [1]
 - *Happens to be one of the original authors of Mapper

Multi-scale Mapper

- “The resulting view of the data [produced by Mapper] through a cover of the codomain offers flexibility in analyzing the data. However, *it offers only a view at a fixed scale at which the cover is constructed.*”

- Multiscale mapper:* a “tower” of simplicial complexes, which is a chain of simplicial complexes connected by [induced] simplicial maps
- Nice benefit:** if the map is a real-valued PL function, the *exact* persistence diagram *from only the 1-skeleton* (!)



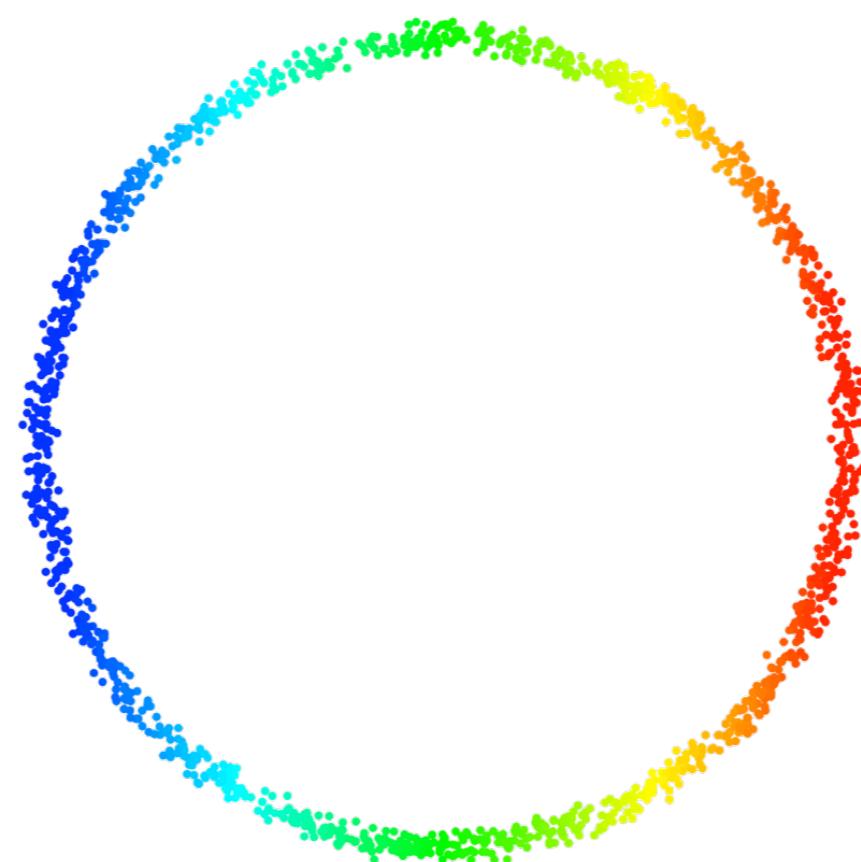
Consider Mappers Complexity

1. Define a **reference map** $f : X \rightarrow Z$ (e.g. $O(n^2)$)
2. Construct a covering $\{U_\alpha\}_{\alpha \in A}$ $+ O(nd)$
3. Apply a **clustering algorithm** \mathcal{C} to the sets X_α $+ O(n^2) + O(n^2) + O(n\alpha)$
4. Obtain a cover $f^*(\mathcal{U})$ of X by considering the path-connected components of $f^{-1}(U_\alpha)$ $+ O(n\alpha)$
5. The Mapper construction is **the nerve of this cover**, either the $O(n^3)$
 1. 1-skeleton
 2. n -skeletonor $O(3^{n/3} \times n^2)$

Problem: How to compute?

- While such theoretical results are amazing, how would one actually compute Mapper for large data sets?
- Consider the following:

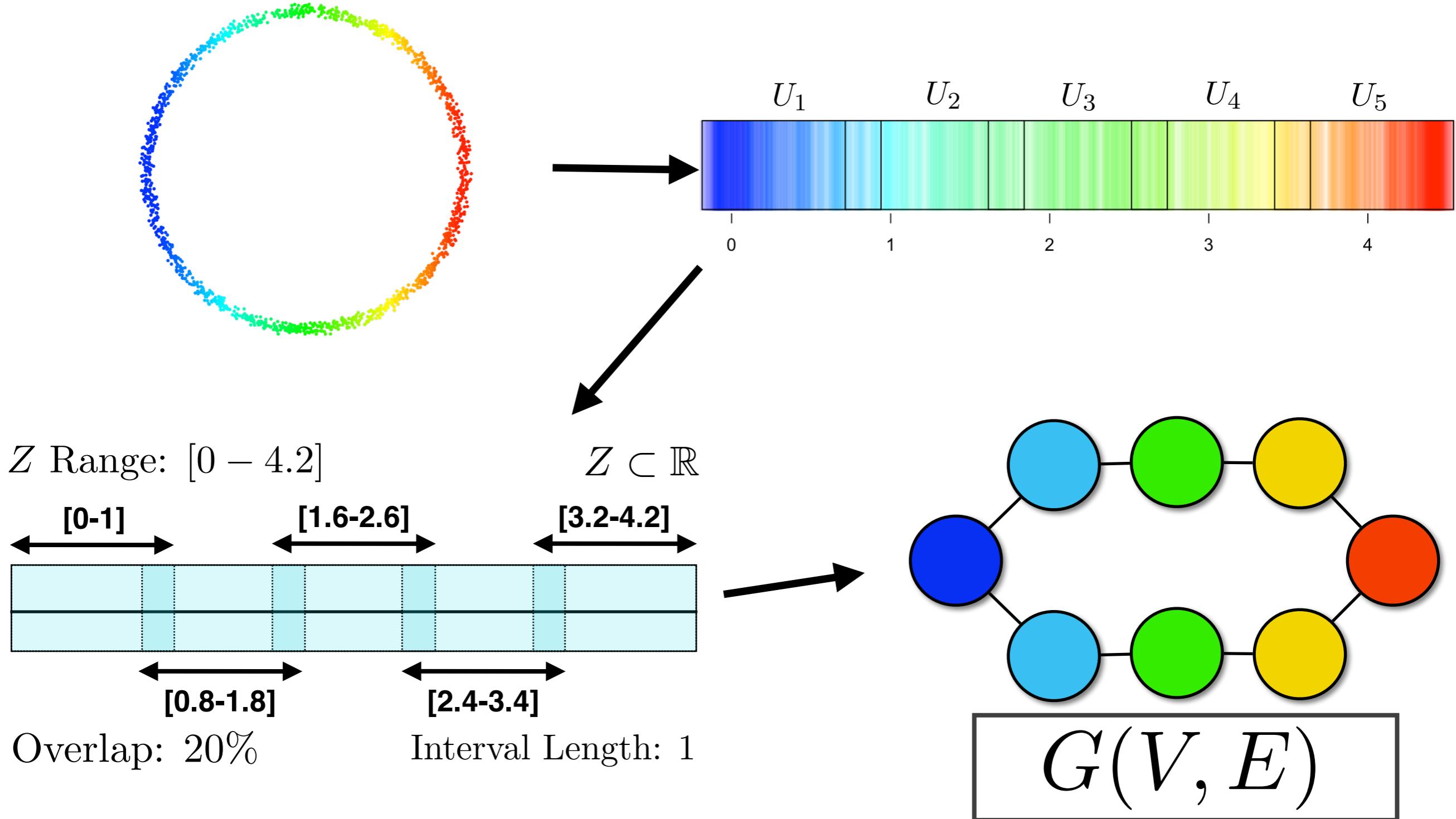
$$X = \{x_1, x_2, \dots, x_n\}$$



Problem: How to compute?

$$X = \{x_1, x_2, \dots, x_n\}$$

$$f(x) = \|x - p\|_2$$

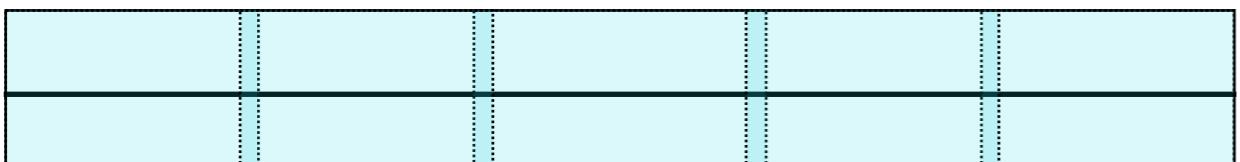


Problem: How to compute?

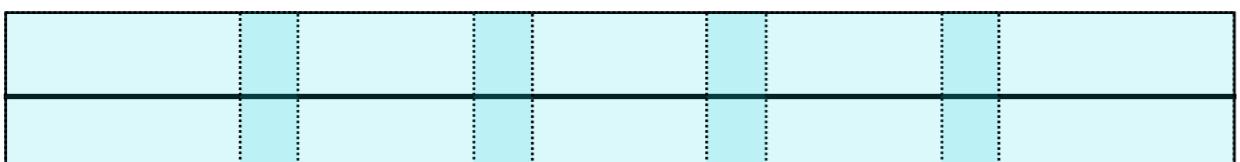
$$f(x) = \|x - p\|_2$$

$$G(V, E)$$

Z Range: $[0 - 4.2]$ $Z \subset \mathbb{R}$



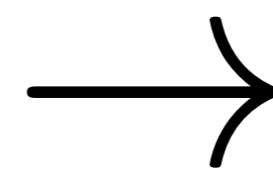
Overlap : 5%



Overlap: 20%



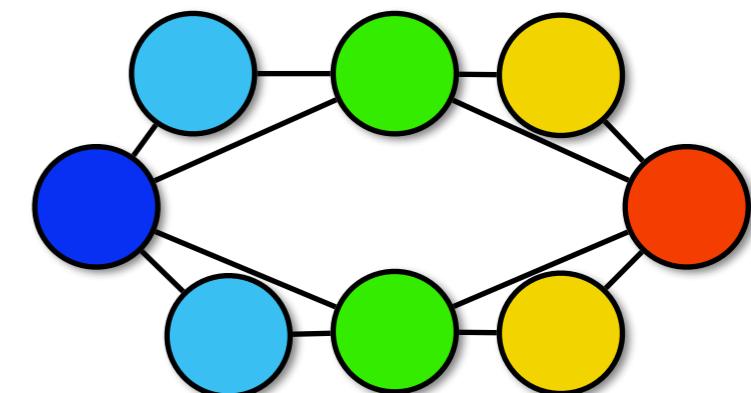
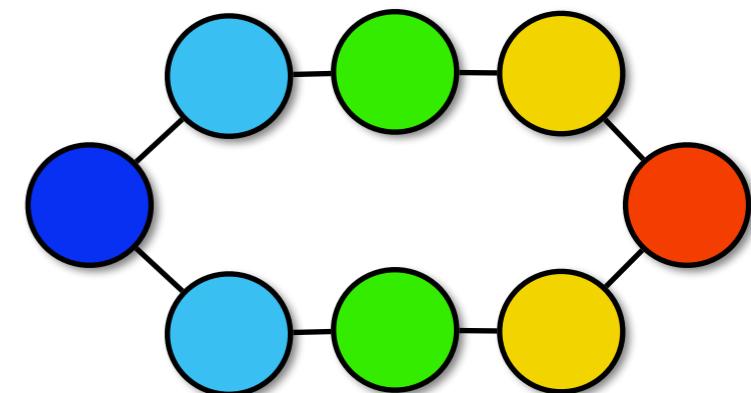
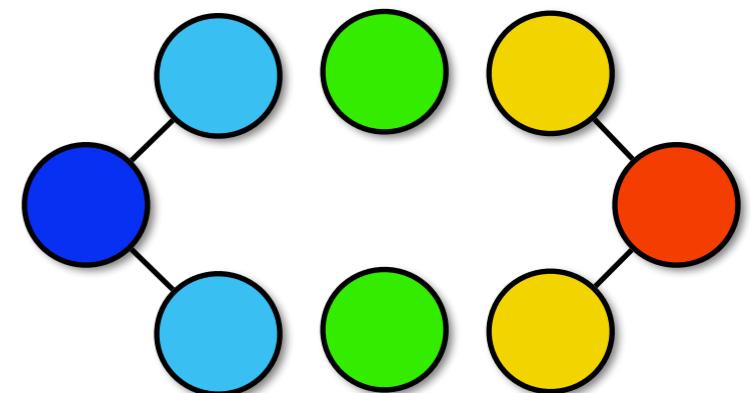
Overlap : 40%



$$\epsilon \rightarrow \infty$$



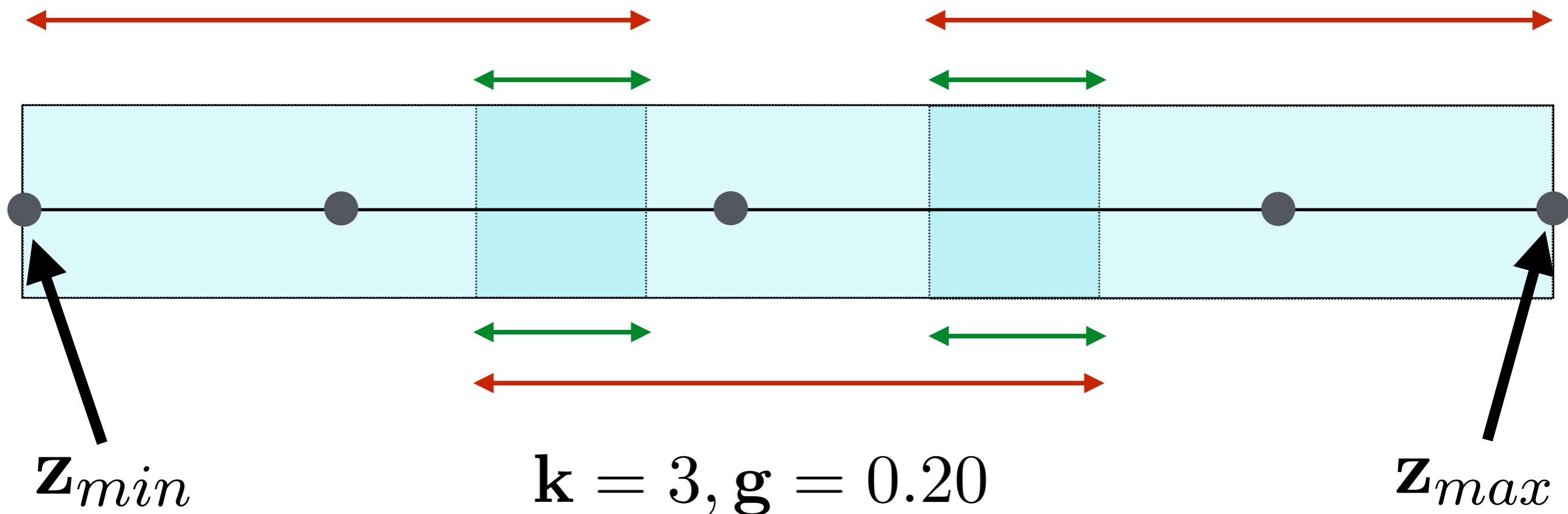
$$\epsilon \rightarrow \infty$$



Consider the problem

$$f : X \rightarrow Z$$

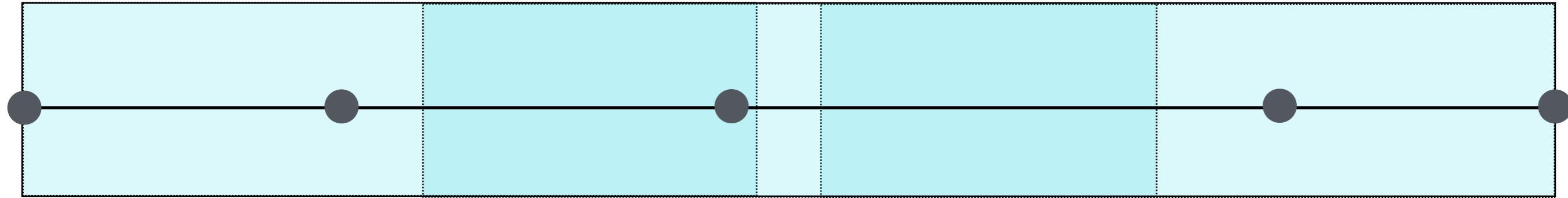
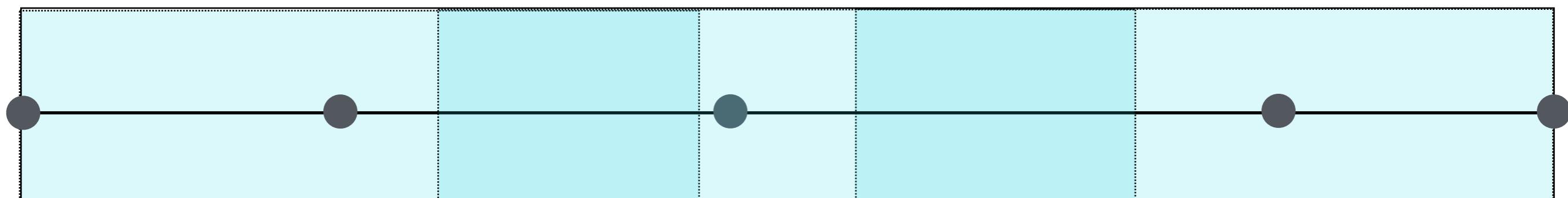
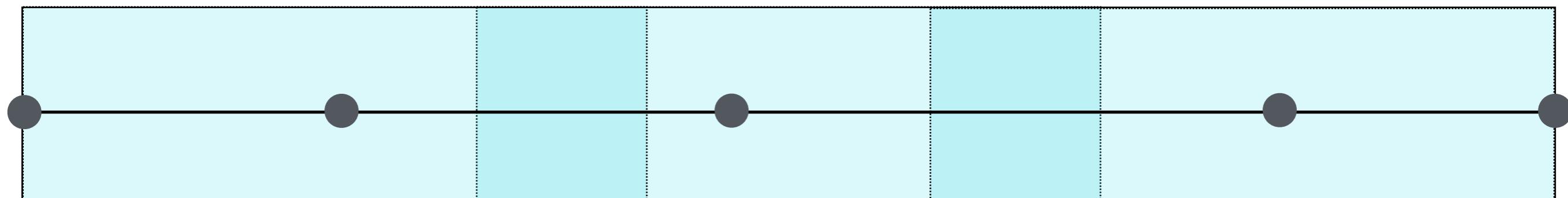
$$\bar{z} = z_{max} - z_{min}$$



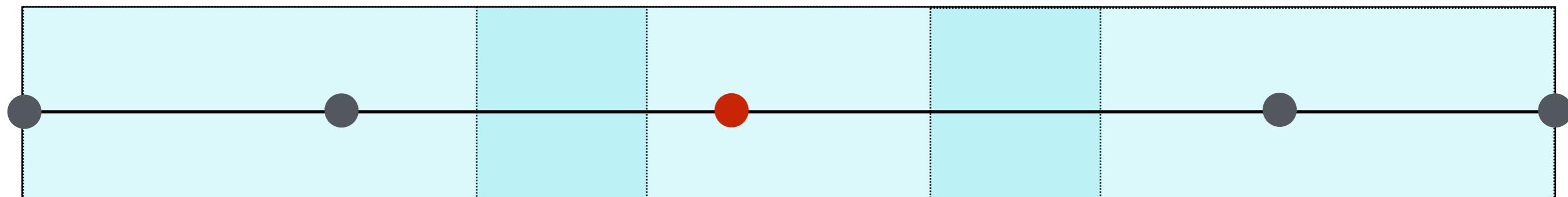
$$r = \frac{\bar{z}}{(k - g(k - 1))} \quad e = r \circ (1 - g)$$

Consider the problem

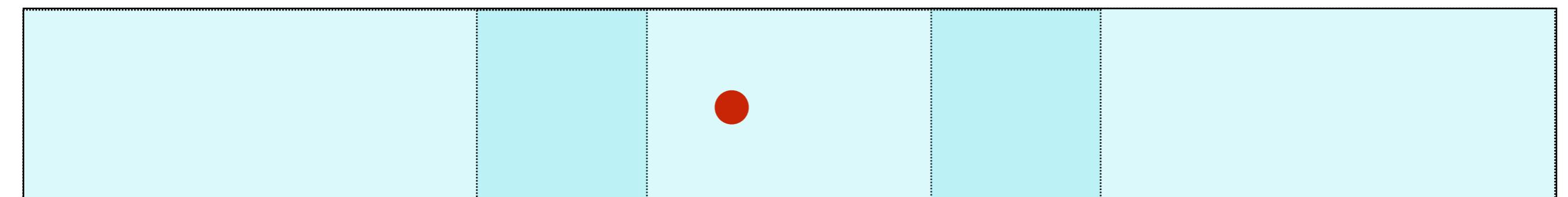
What is the smallest overlap value that could induce a new simplicial complex?



Consider the problem

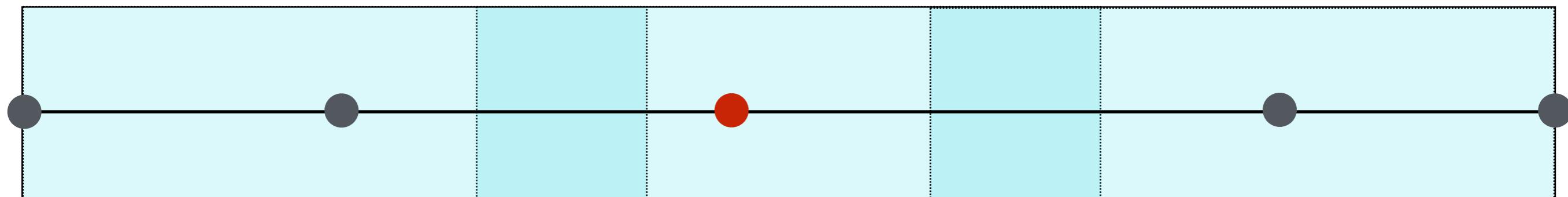


What is the smallest overlap values that *this point* could potentially induce a new complex?

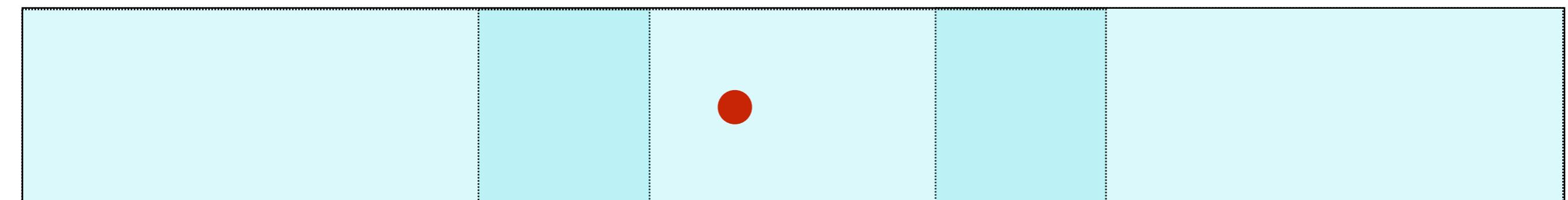


What interval length will this point intersect a new level set?

Consider the problem



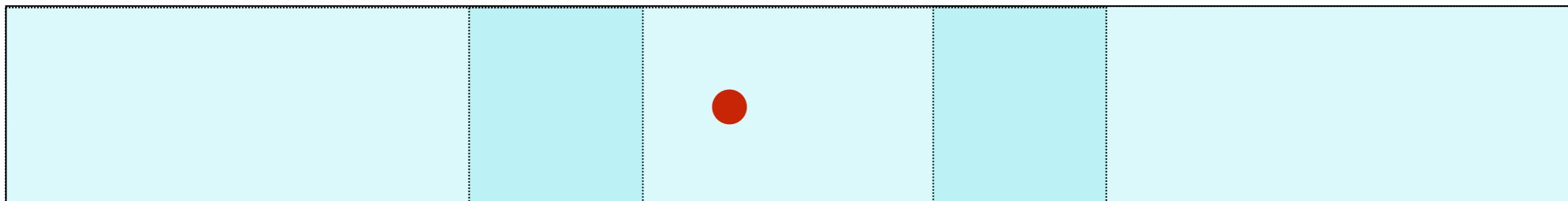
What is the smallest overlap values that *this point* could potentially induce a new complex?



Equivalent: What interval length will this point intersect a new level set?

Consider the problem

Equivalent: What interval length will this point intersect a new level set?



Recall:

$$r = \frac{\bar{z}}{(k - g(k-1))}$$

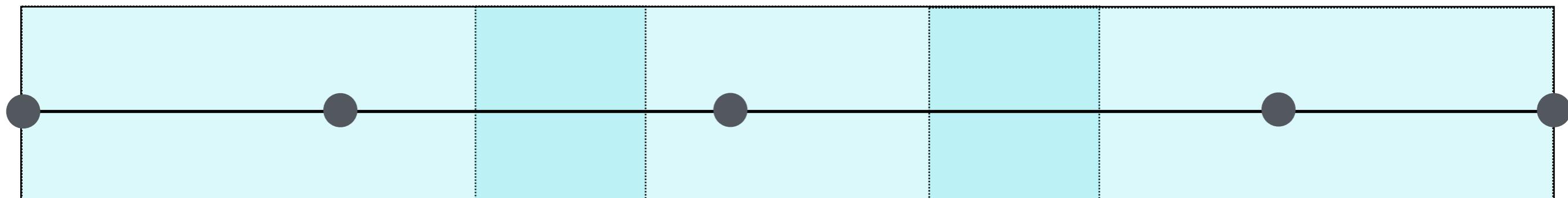
Idea:

$$g = \frac{\bar{z} - rk}{r(-k + 1)}$$

Observation: If k is fixed,
just need r !

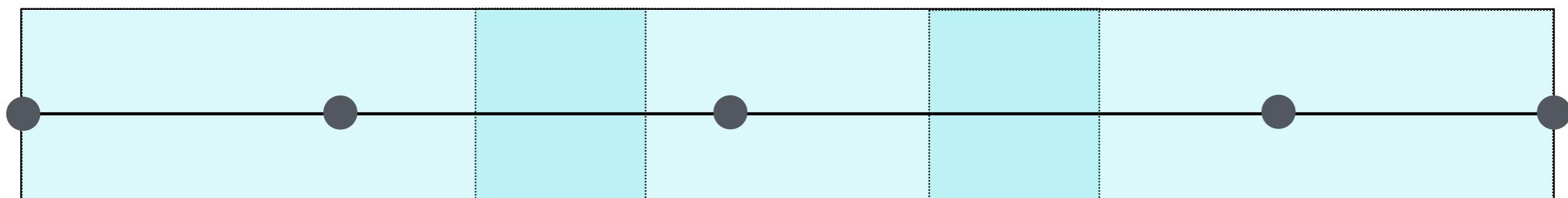
Consider the problem

Problem: Keeping track of each level set may be computationally difficult



Observation: Each point is already associated with an index...

$$\{U_\alpha\}_{\alpha \in A}$$



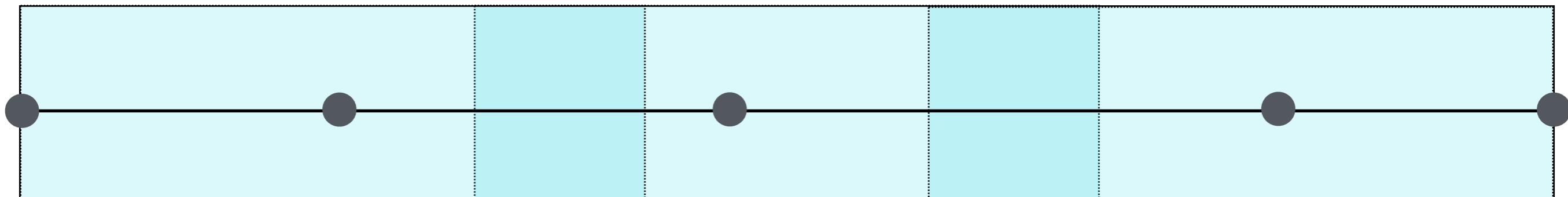
$$\alpha_0 = (0)$$

$$\alpha_1 = (1)$$

$$\alpha_2 = (2)$$

Idea: Each point is already associated with an index...

Observation: Each point is already associated
with a set of indexes... $\{U_\alpha\}_{\alpha \in A}$

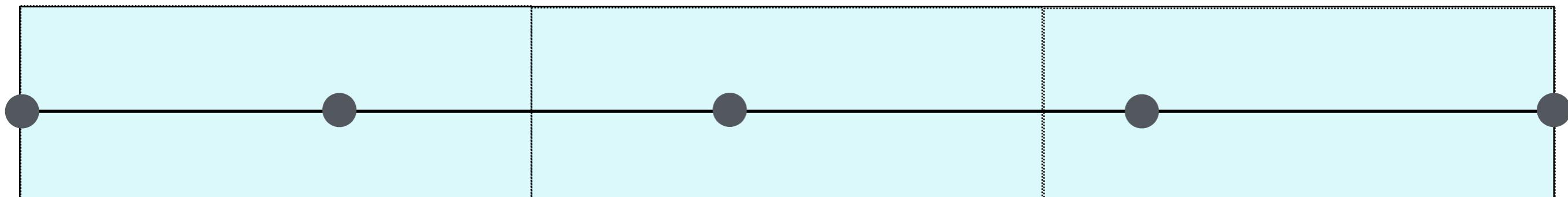


$$\alpha_0 = (0)$$

$$\alpha_1 = (1)$$

$$\alpha_2 = (2)$$

When there is 0 overlap, there is only one index...

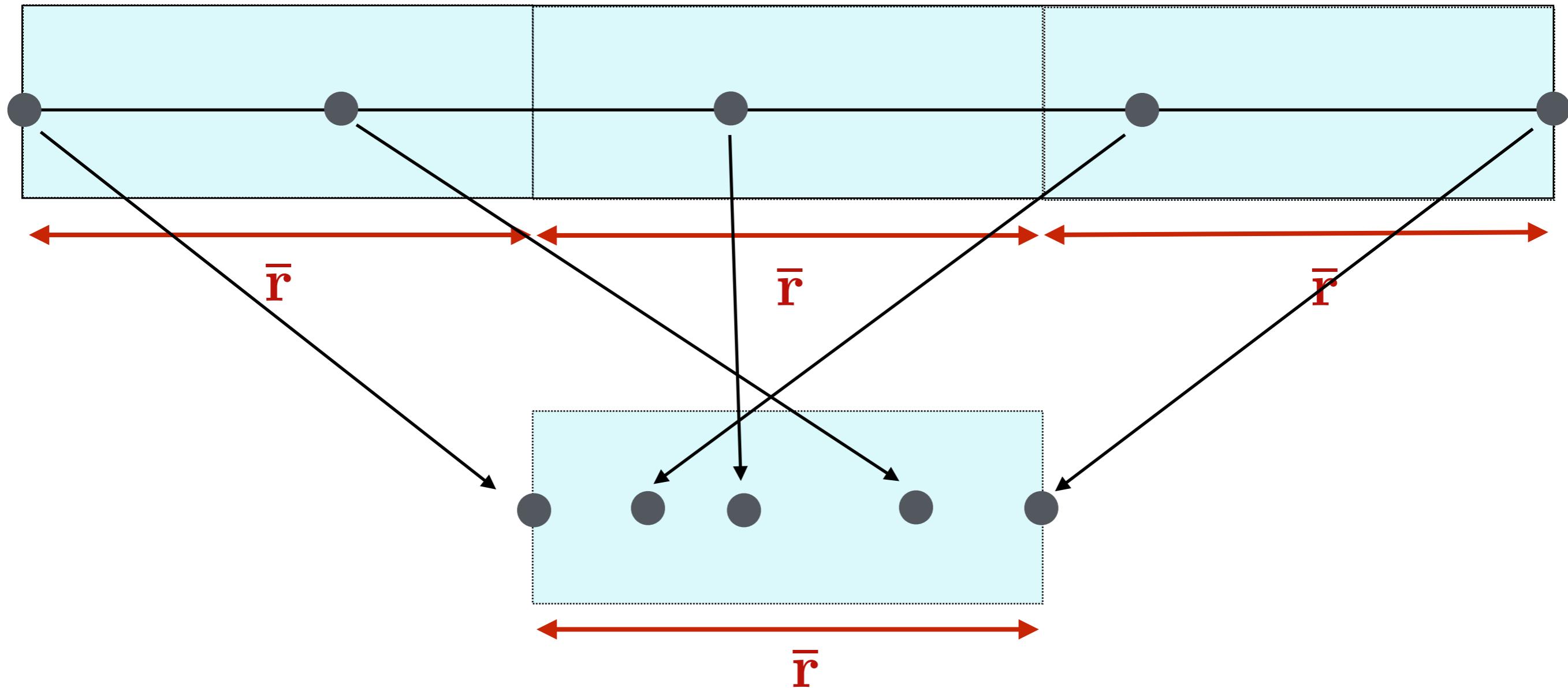


What if we record that index for each point?

$$\mathbf{A} = [\alpha_0, \alpha_0, \alpha_1, \alpha_2, \alpha_2]^\top \longrightarrow \mathbf{A} = [0, 0, 1, 2, 2]^\top$$

And start with this *base* interval length

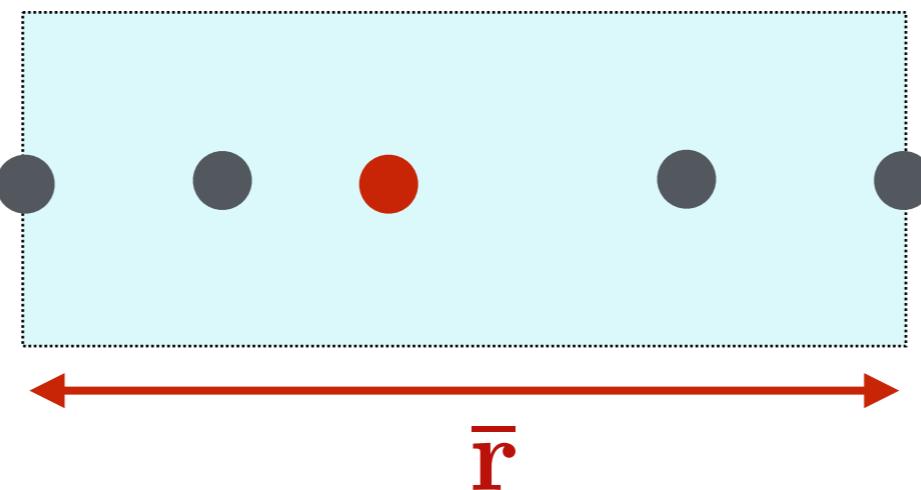
$$k = 3, g = 0 \longrightarrow \bar{r} = \frac{\bar{z}}{k} \quad e = 0$$



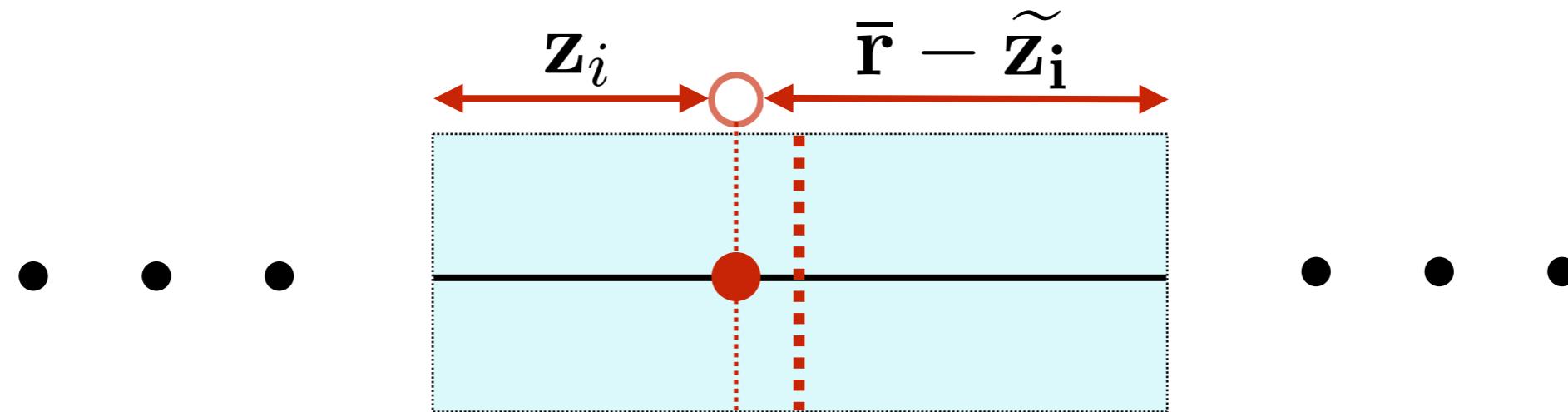
And project to a “unit” box...

$$\tilde{z} = (z - z_{\min}) - A \circ \bar{r}$$

And project to a “unit” box...

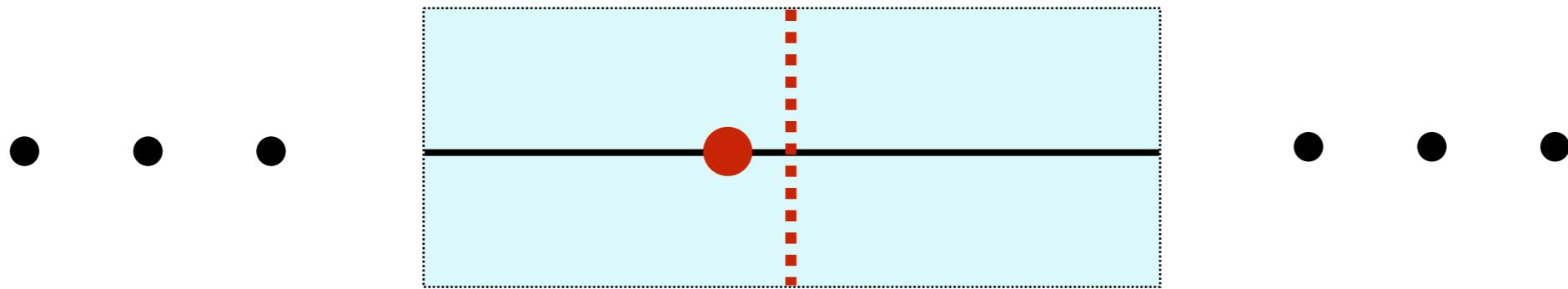


Goal: What's the distance to the closest box?



Observation: Boxes expand ‘uniformly’ in both directions
Just need the distance to the outside of the interval

$$\mathbf{Z}_\Delta = [\dots, \min(\tilde{\mathbf{z}}_i, \bar{r} - \tilde{\mathbf{z}}_i), \dots]^T$$



The minimum interval length for each point to intersect its closest level set is thus:

$$\hat{\mathbf{R}} = \mathbf{Z}_\Delta + \bar{\mathbf{r}}$$

The corresponding overlap values?

$$G = \frac{(\bar{\mathbf{z}} - \hat{\mathbf{R}}k)}{\hat{\mathbf{R}}(-k + 1)}$$

Why is this projection useful?

- Consider the following strategy to compute Mapper
 1. Construct “base” cover $\mathcal{U}[\bar{r}, 0]$ ($\bar{r} = \bar{z}/k$)
 2. Compute all overlap/interval values, sort by increasing value
 - * Note: We know which level sets each point will intersect (and at what finite set of overlap values this will occur!)
 3. Idea: Start with base Mapper, then *incrementally* increase the interval size, i.e. $\bar{r} + \epsilon$
 - * Note: For each computed ϵ , only need update the clusters in level sets whose point memberships have changed
 4. (3) implies that we only need to update that simplexes that are *immediately* path-connected with the level sets that just changed
- Could this be used to efficiently compute an approximation to Multi-scale Mapper???

Why is having all overlap values useful?

Consider Mappers complexity...

Filter function	Form cover	Distance Matrix	Cluster Hierarchically	Cut Tree
$f : X \rightarrow Z$ + (e.g. $O(n^2)$)	$\{U_\alpha\}_{\alpha \in A}$ + $O(nd)$	$D(X_\alpha)$	$C_H(X_\alpha)$	$C(X_\alpha)$
(needed every instance of Mapper)	(needed every instance of Mapper)		(needed every instance of Mapper)	
Form connected components	Form 1-skeleton		Form n -skeleton	
+ $f^*(\mathcal{U})$	+ $M(\mathcal{U}, f)^{(1)}$		$M(\mathcal{U}, f)^{(k)}$	
$O(n\alpha)$	$O(n^3)$		$O(3^{n/3} \times n^2)$	
(needed every instance of Mapper)	(needed every instance of Mapper)		(needed every instance of Mapper)	

Why is having all overlap values useful?

Incremental Mappers complexity...

$$\kappa = |f^*(U_i)| \cdot |f^*(U_j)| \quad \forall \quad i \xrightarrow{P} j$$

Filter function	Form cover	Distance Matrix	Cluster Hierarchically	Cut Tree
$f : X \rightarrow Z$	$\{U_\alpha\}_{\alpha \in A}$	$D(X_\alpha)$	$C_H(X_\alpha)$	$C(X_\alpha)$
(e.g. $O(n^2)$)	$O(nd)$	$O(n^2) + O(n^2) + O(n\alpha)$		
(needed once)	(needed once)		(needed per updated level set)	
Form connected components	Form 1-skeleton		Form n -skeleton	
$+ f^*(\mathcal{U})$	$+ M(\mathcal{U}, f)^{(1)}$		$M(\mathcal{U}, f)^{(k)}$	
$O(n\alpha)$ (needed once)	$O(n^3)$ (needed once)	or	$O(3^{n/2} \times n^2)$	
$O(\kappa\alpha)$ (per update)	$O(\kappa)$ (per update)			(needed every instance of Mapper)

Conclusion

- TDA is an emerging field
- Some see it as a **more formal generalization of clustering**
 - *How did I get here? :)*
- Mapper is becoming **remarkably popular!**
 - Ayasdi was awarded **\$106 million** in funding towards their solution that uses Mapper for TDA [2]
 - Has proven *useful* in several data-analysis domains already!
- Mapper provides **useful summary information**, e.g. high-dimensional data sets
- ***Can Topology provide a solution to all unsupervised things in ML?***
 - “The reader should not conclude that the subject is quickly or painlessly learned. [The field of topology] is the impetus of future work: hard, slow, and ***fruitful***. “ - Robert Ghrist, Applied Elementary Topology
 - Check out **my extension** to the **open-source Mapper** if you want to actually use for experiments: <https://github.com/peekxc/TDAmapper>

Questions?