

UMAP

Uniform Manifold Approximation and Projection
for dimension reduction

Who am I?

I am a research mathematician at the Tutte Institute for Mathematics and Computing

My Ph.D. was in Profinite Lie Rings
(no, you don't care)

I now work on applying topological techniques
to unsupervised learning problems

What is Dimension Reduction?

Find the “latent”
features in your data

5	0	4	1	9	2	1	3	1	4
3	5	3	6	1	7	2	8	6	9
4	0	9	1	1	2	4	3	2	7
3	8	6	9	0	5	6	0	7	6
1	8	7	9	3	9	8	5	9	3
3	0	7	4	9	8	0	9	4	1
4	4	6	0	4	5	6	1	0	0
1	7	1	6	3	0	2	1	1	7
9	0	2	6	7	8	3	9	0	4
6	7	4	6	8	0	7	8	3	1



Matrix Factorization

Neighbour Graphs

Matrix Factorization

Principal Component Analysis

Linear Autoencoder

Latent Dirichlet Allocation

Non-negative Matrix Factorization

Generalised Low Rank Models

Word2Vec

GloVe

Neighbour Graphs

Laplacian Eigenmaps

Hessian Eigenmaps

Local Tangent Space Alignment

JSE

Isomap

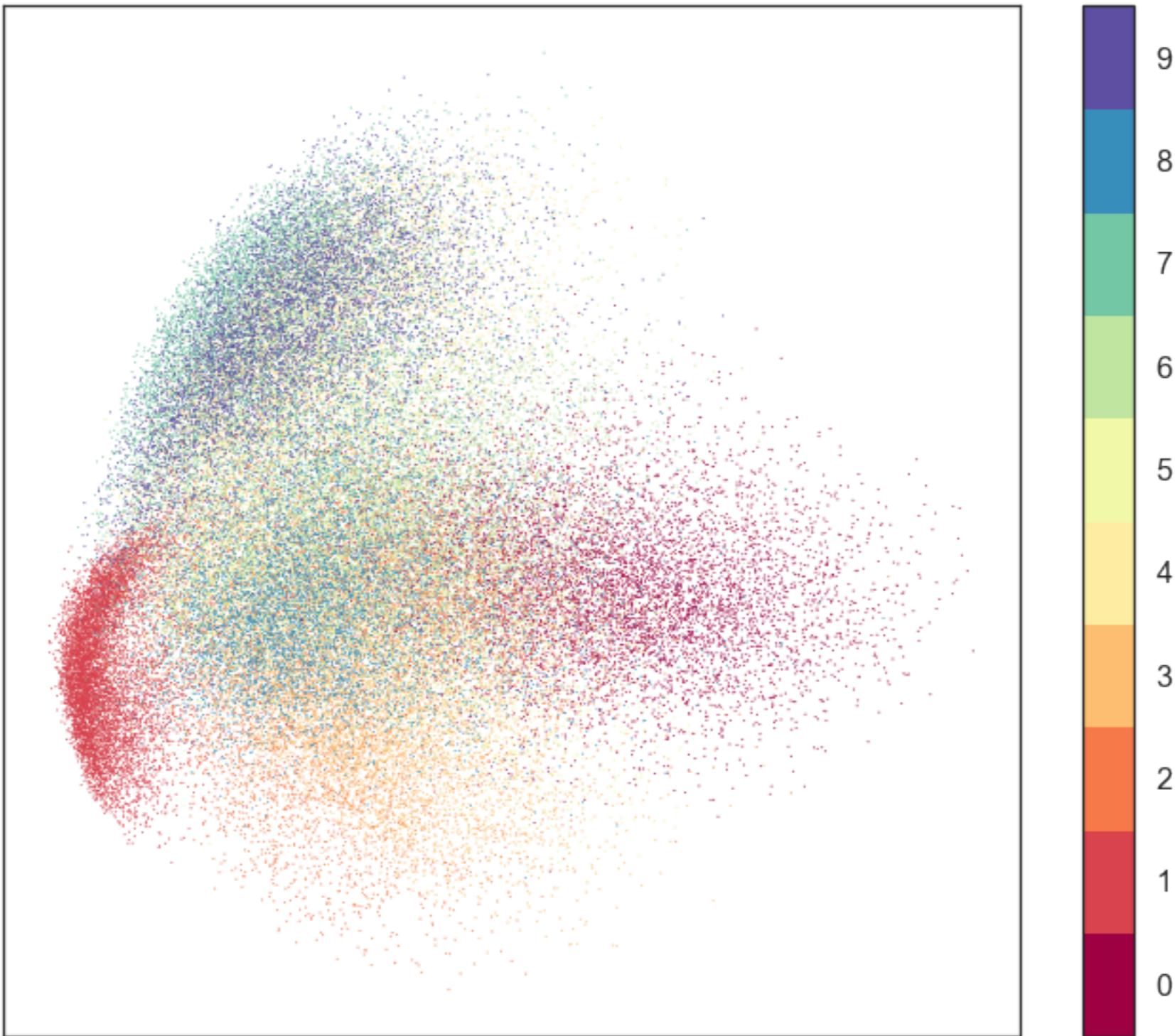
t-SNE

Locally Linear Embedding

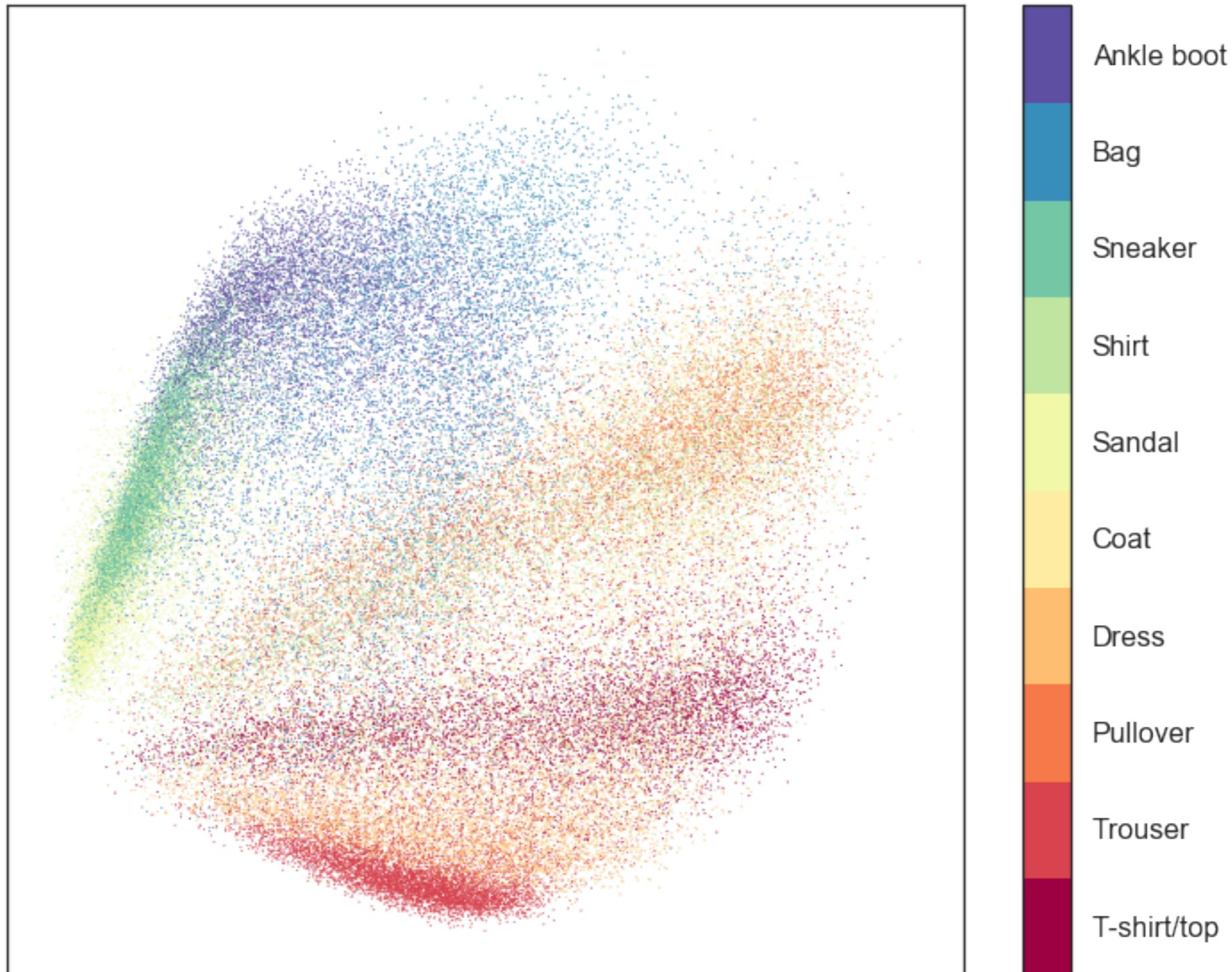
UMAP

PCA is the prototypical
matrix factorization

PCA on MNIST digits

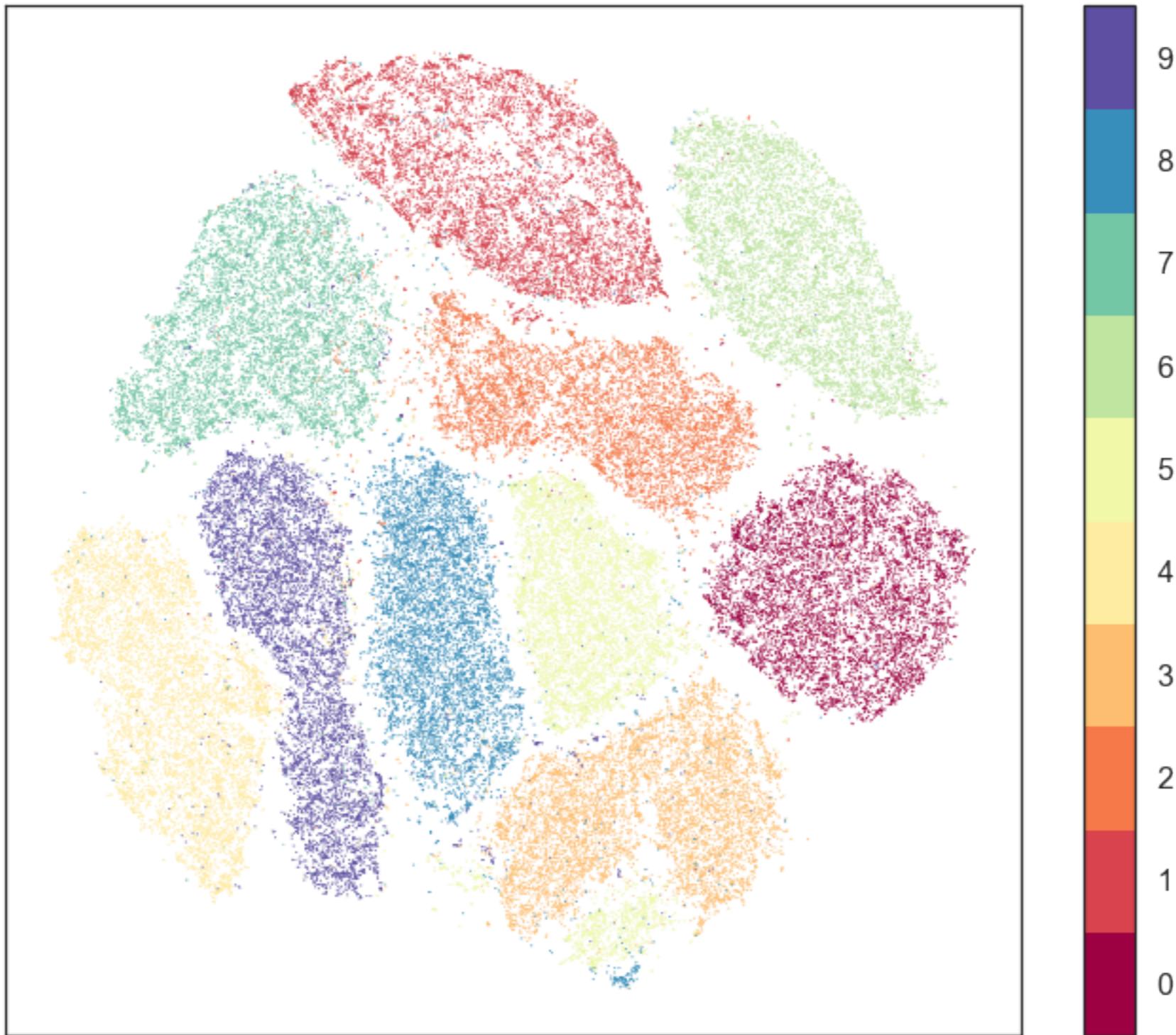


PCA on Fashion MNIST

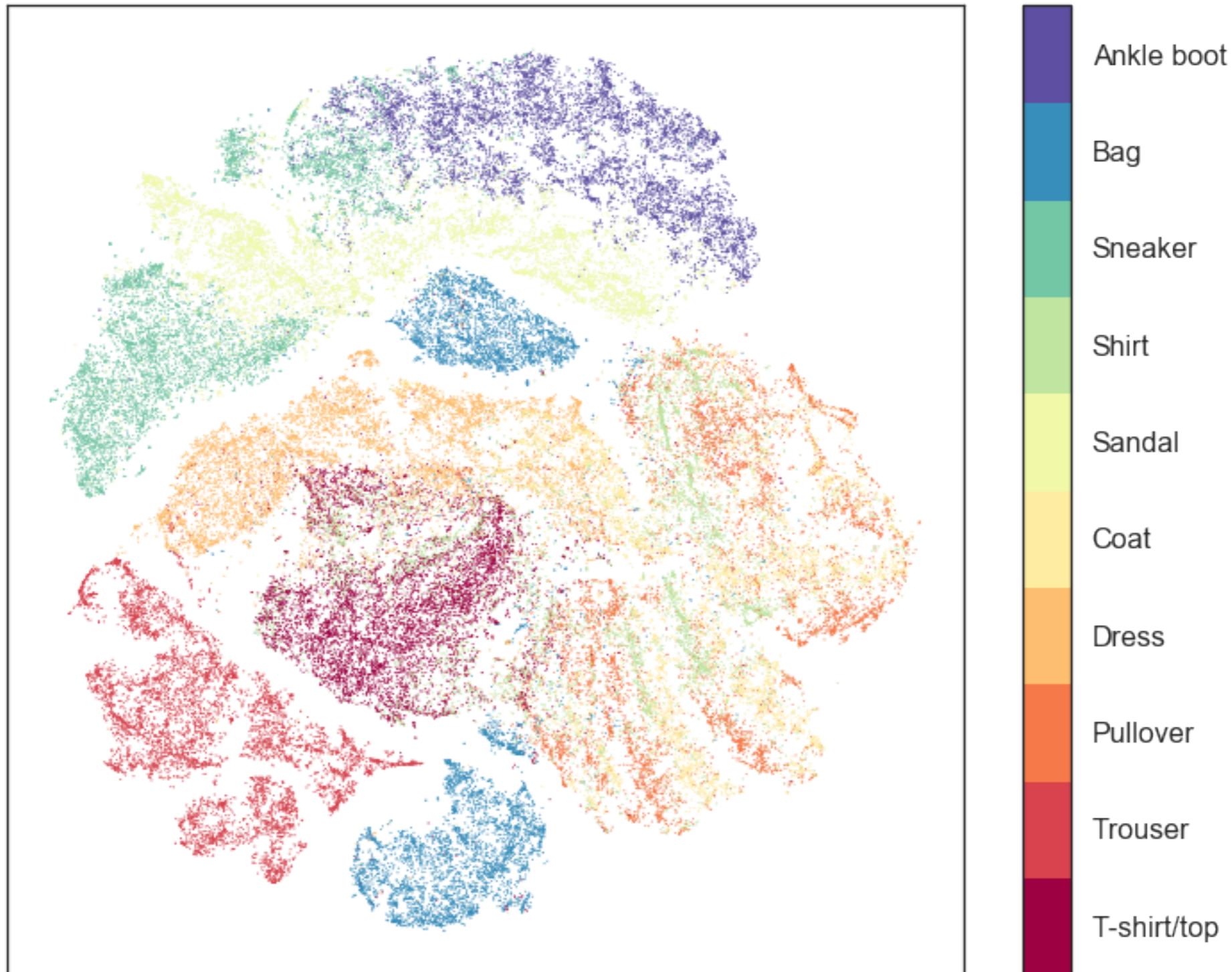


t-SNE is the current
state-of-the art for
neighbour graphs

t-SNE on MNIST digits



t-SNE on Fashion MNIST

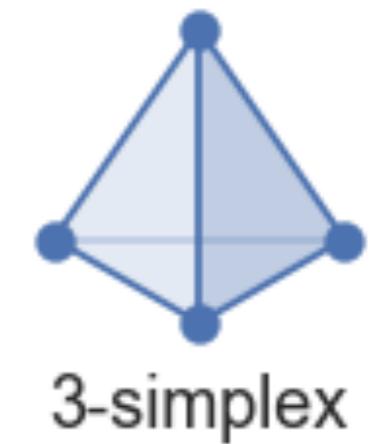


Uniform Manifold Approximation and Projection

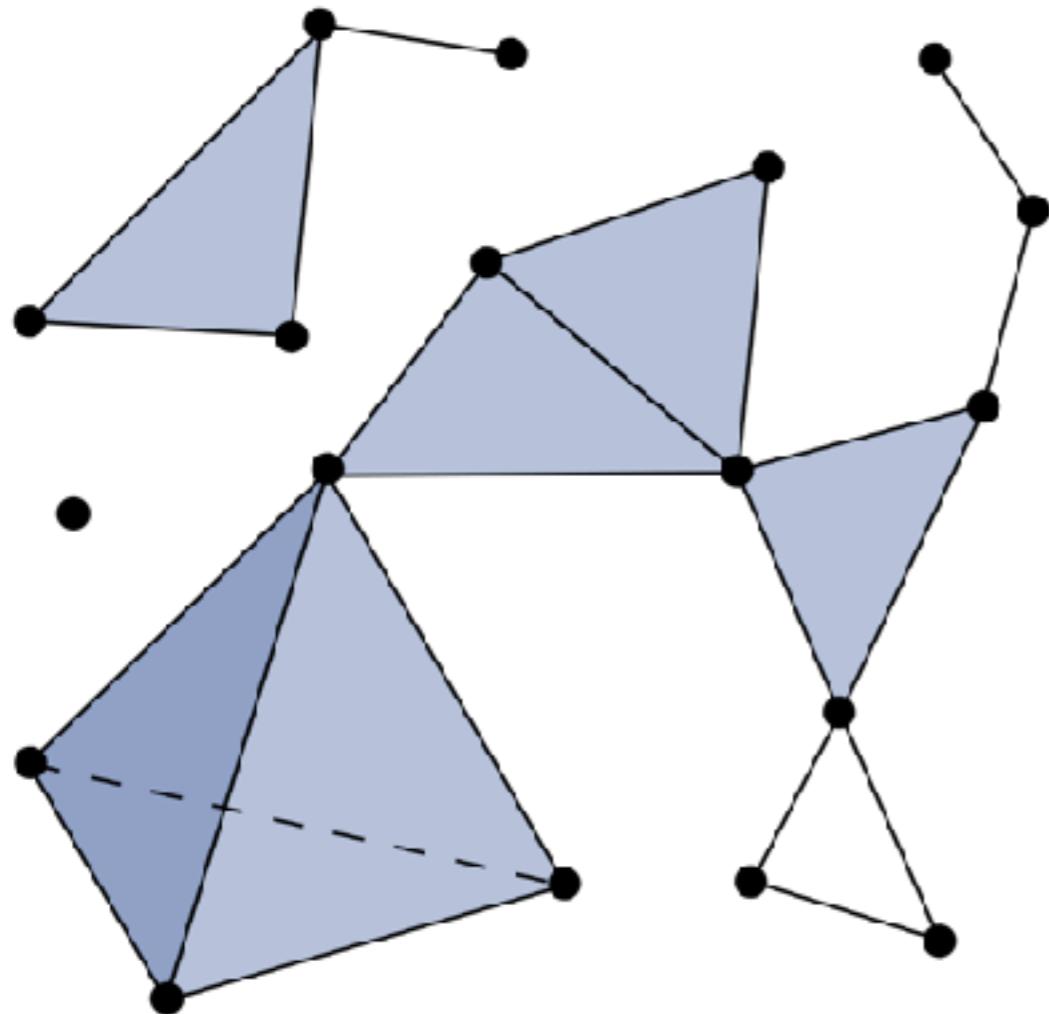
UMAP builds
mathematical theory to
justify the graph based
approach

First, a little bit of
topological data
analysis...

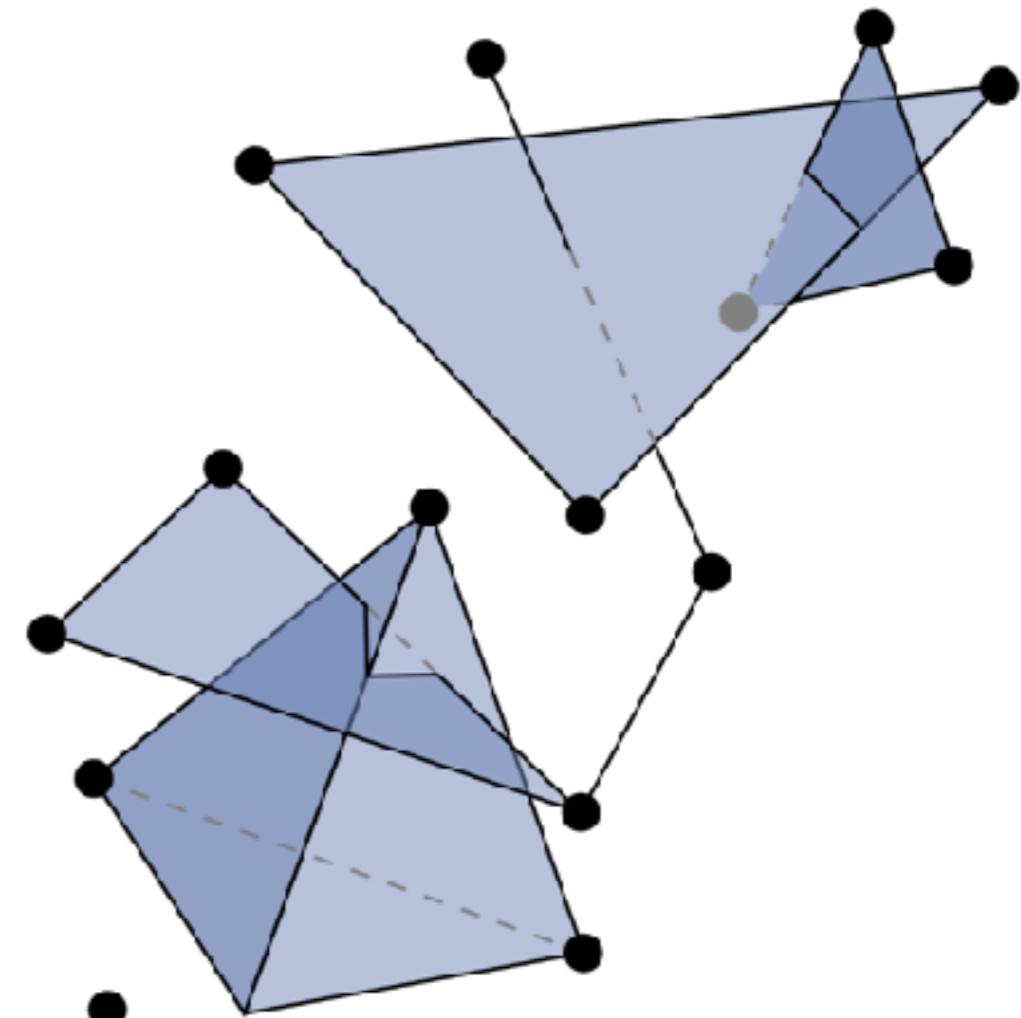
Simplices



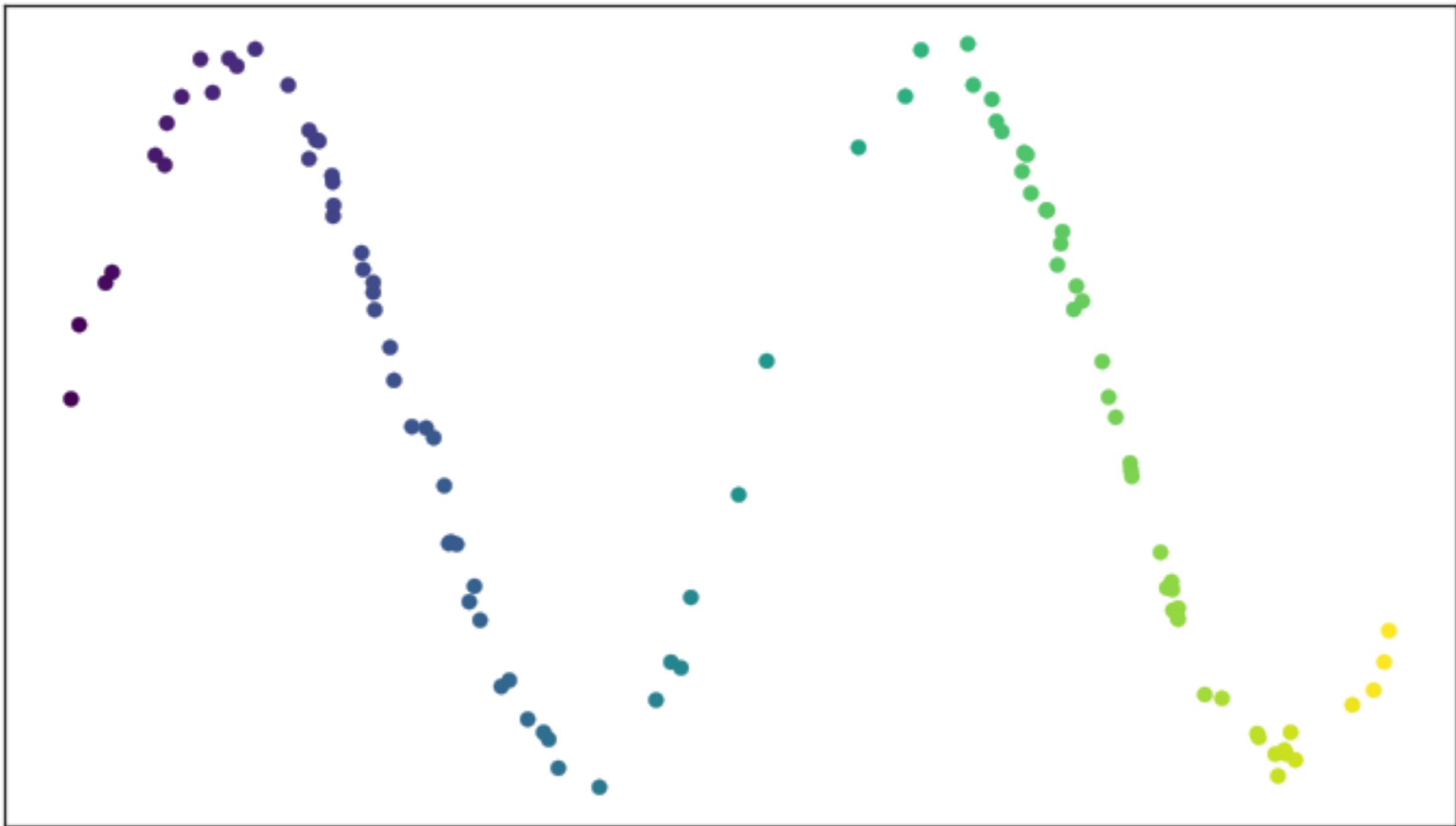
Simplicial complex

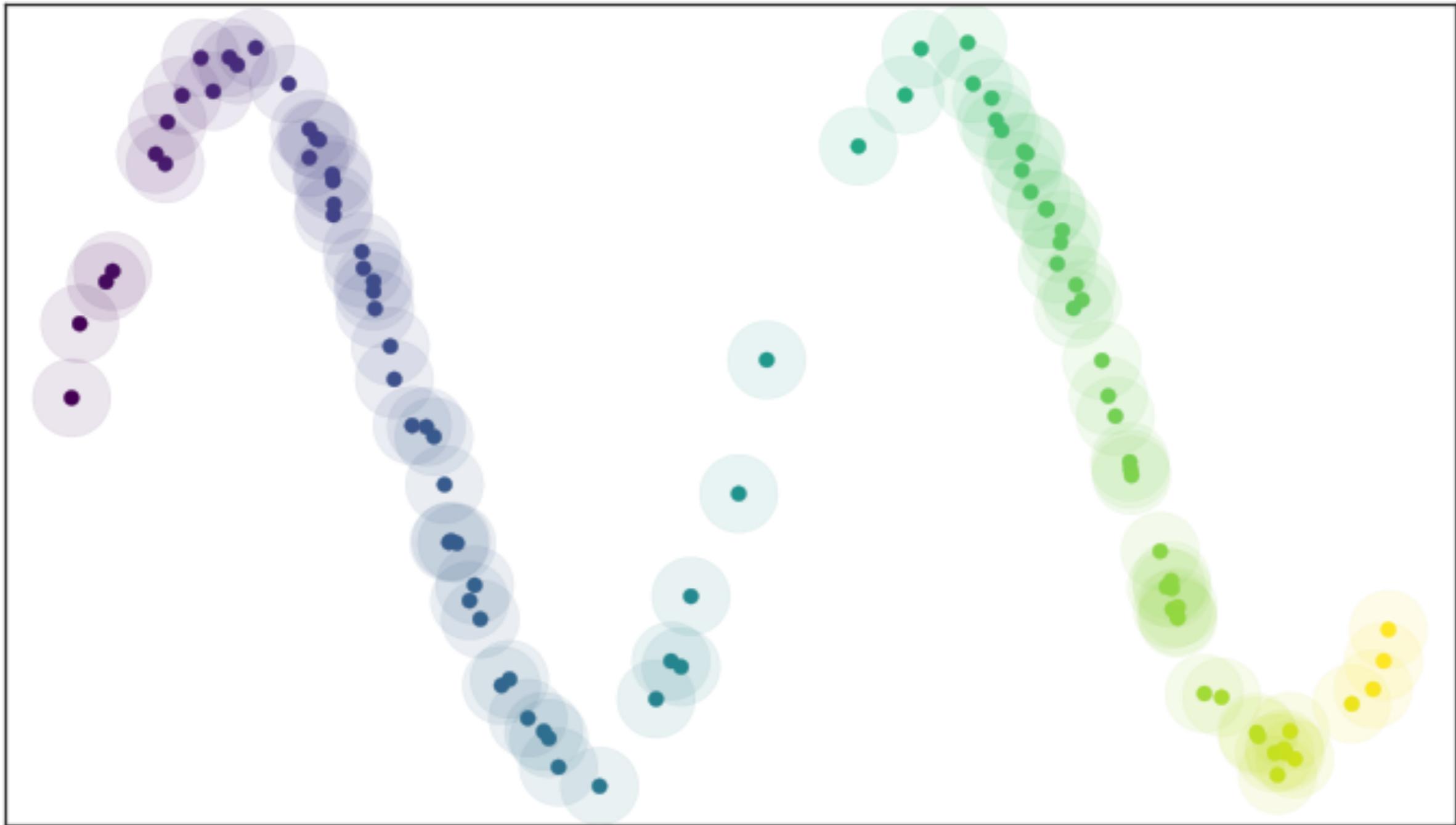


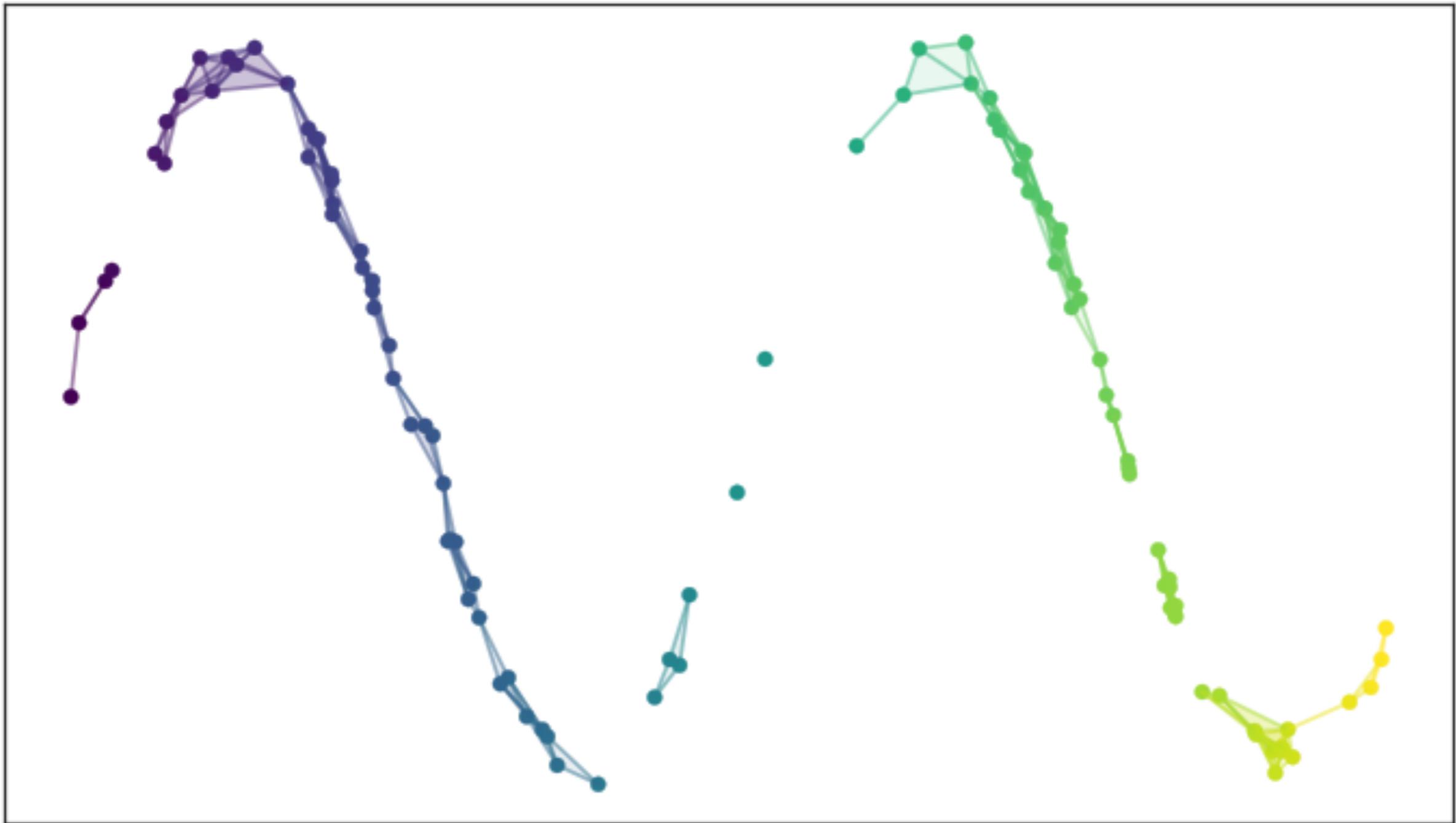
Collection of Simplices



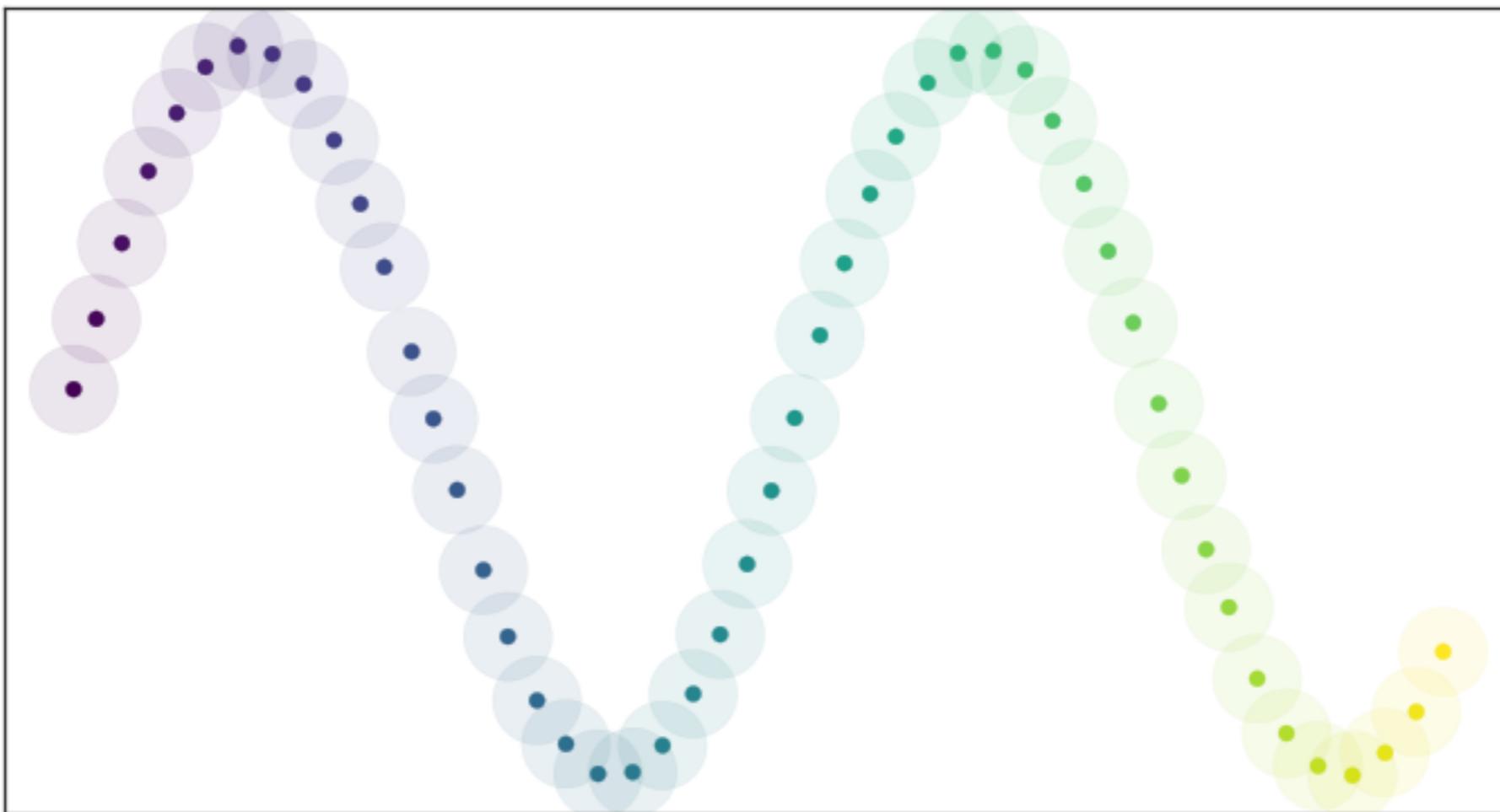
Theorem 1 (Nerve theorem). *Let $\mathcal{U} = \{U_i\}_{i \in I}$ be a cover of a topological space X . If, for all $\sigma \subset I$ $\bigcap_{i \in \sigma} U_i$ is either contractible or empty, then $\mathcal{N}(\mathcal{U})$ is homotopically equivalent to X .*







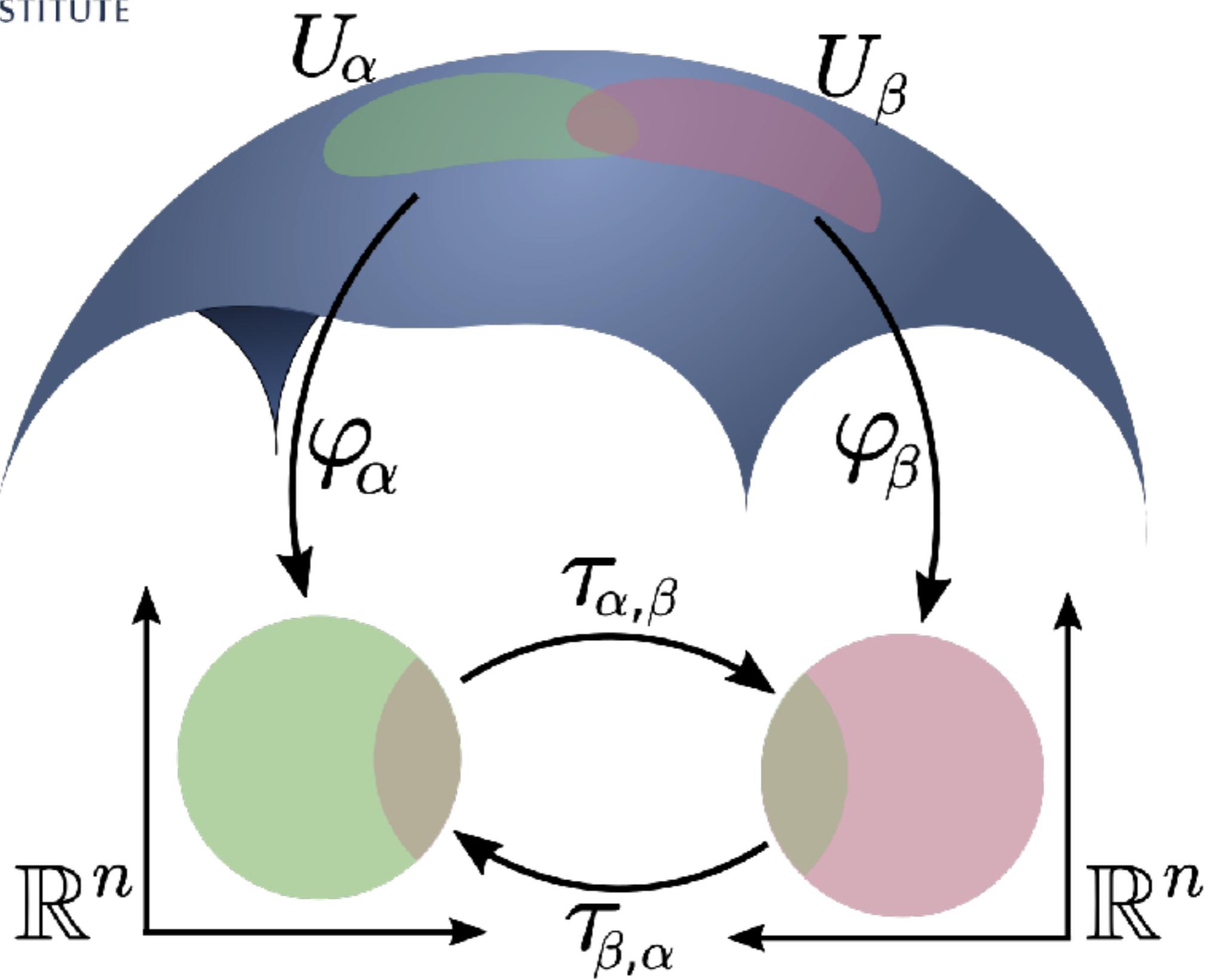
If the data is uniformly distributed on the manifold then the cover will be “good”

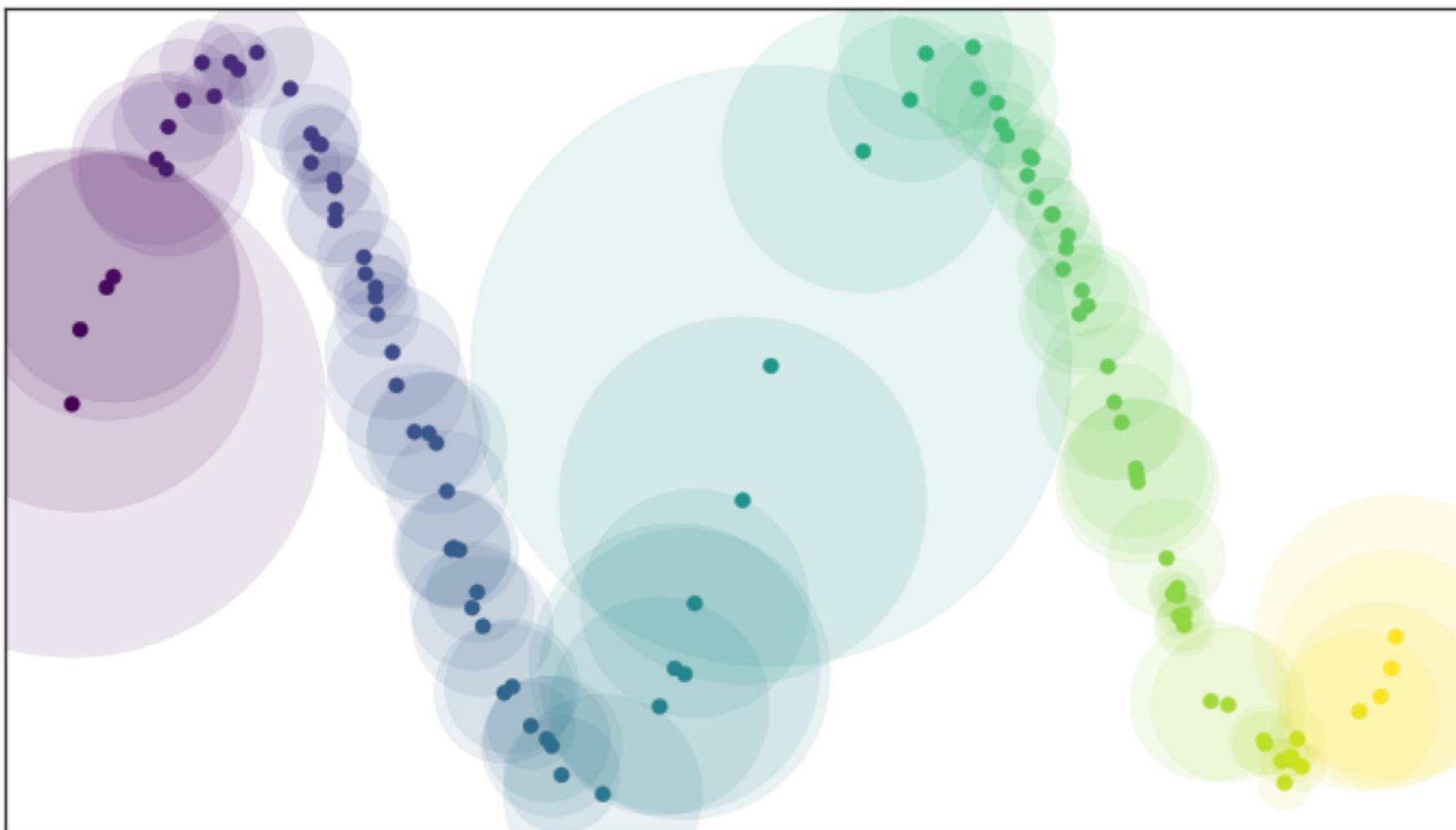


When is data that
nicely behaved?

Assumption:
Data is uniformly
distributed on the
manifold

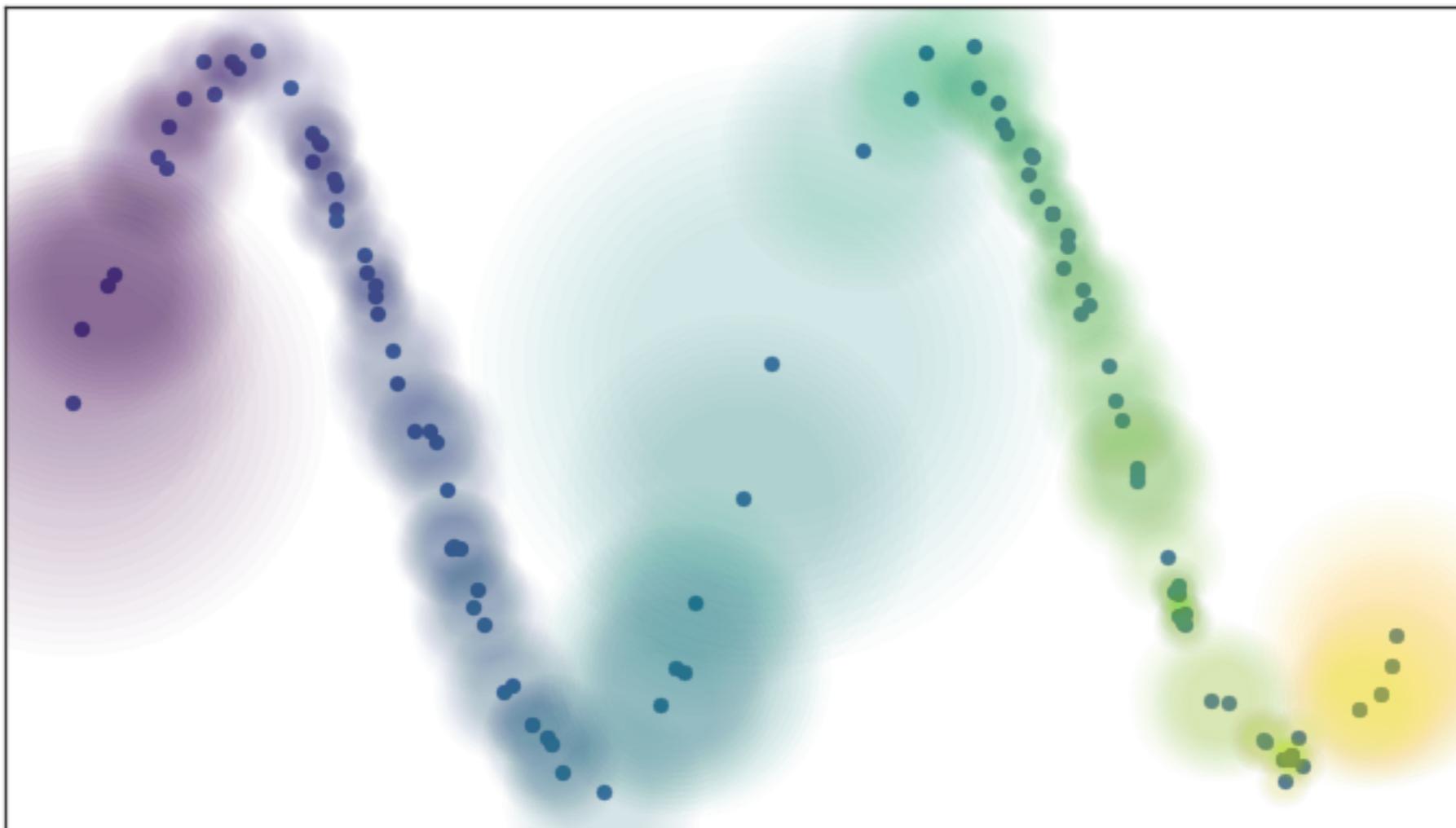
Define a Riemannian
metric on the manifold
to make this
assumption true



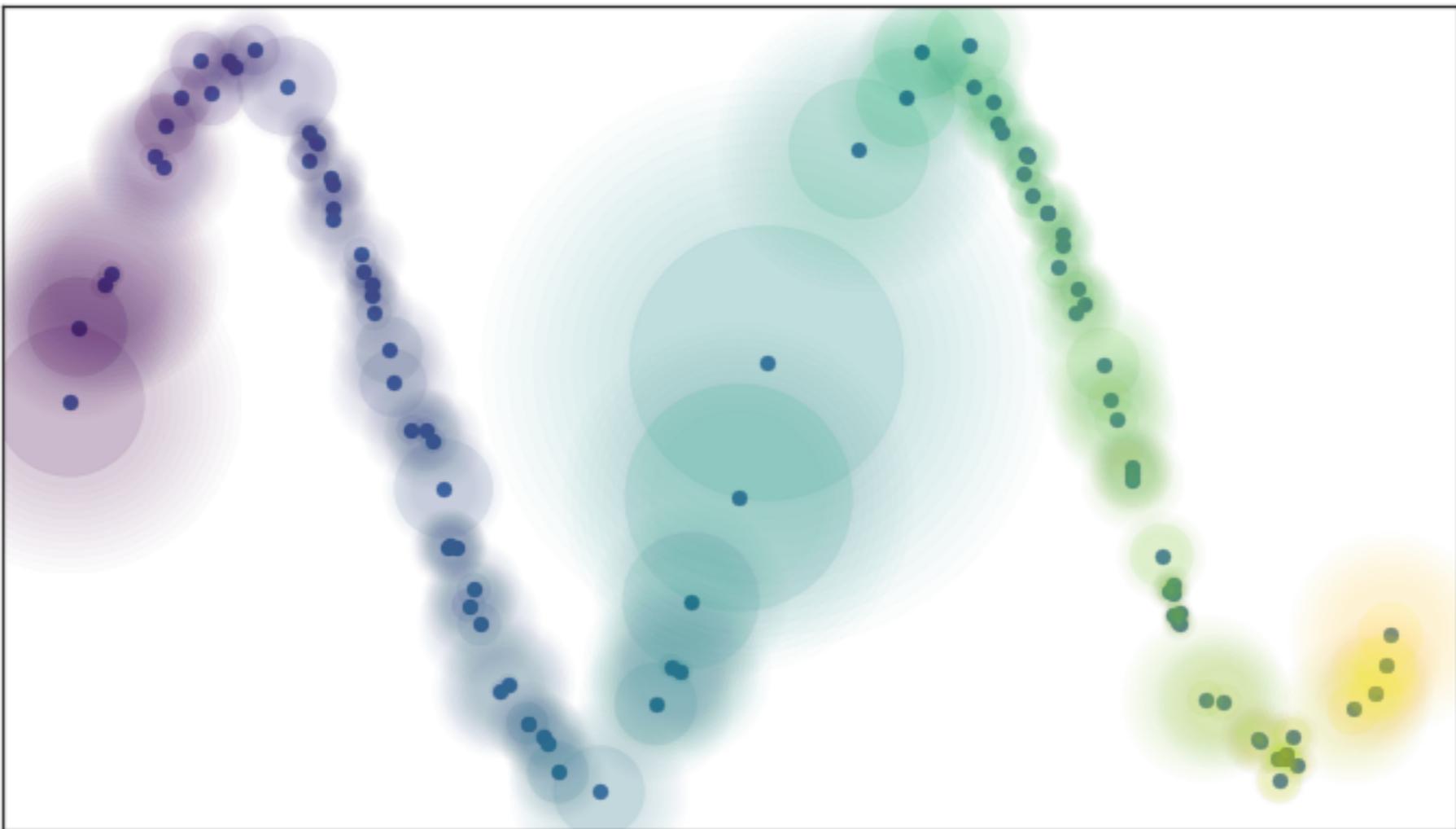


Why choose a fixed
radius? Why not have a
fuzzy cover?

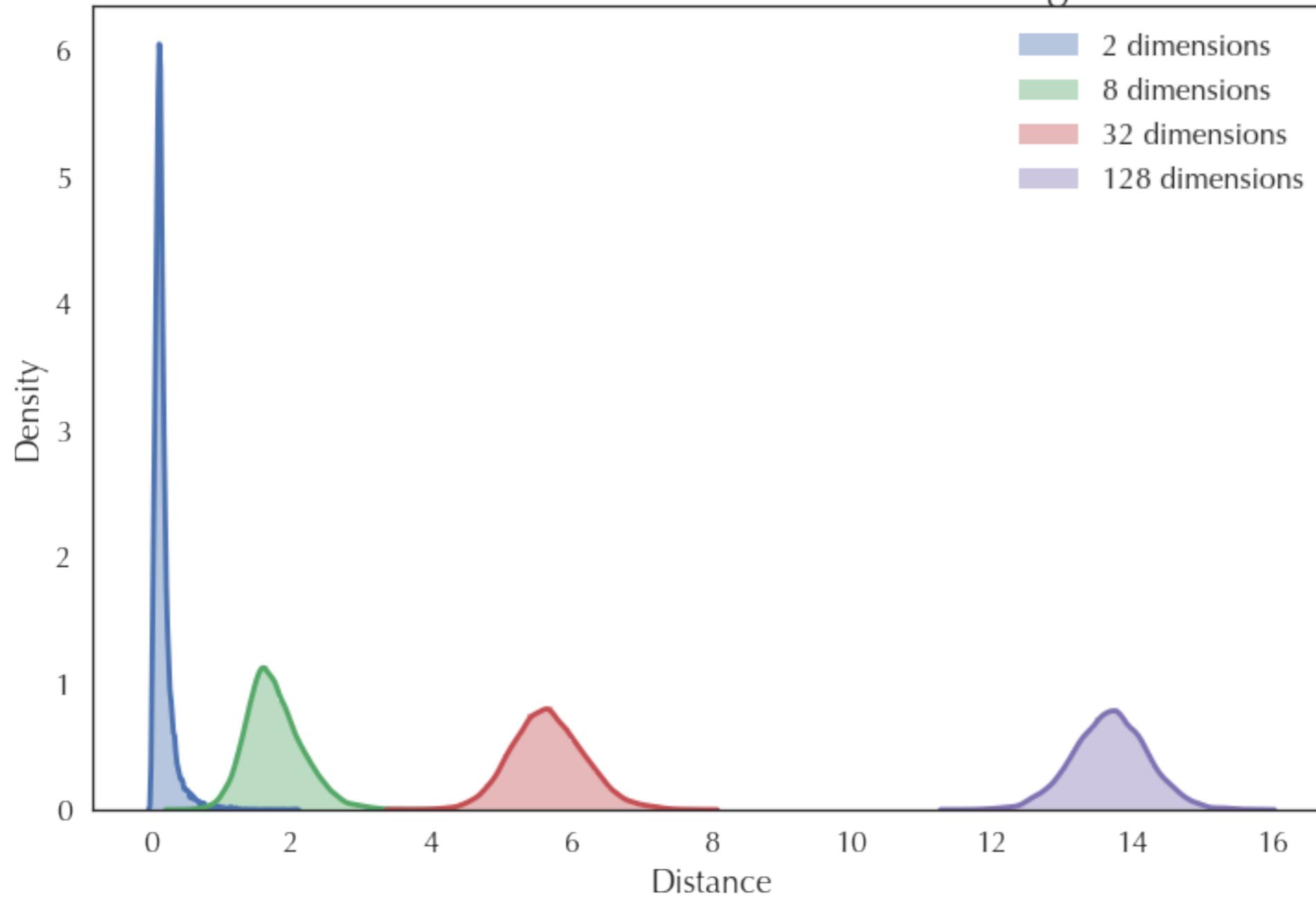
Theorem 2 (UMAP Adjunction). *The functors $\text{FinReal} : \text{sFuzz} \rightarrow \text{FinEPMet}$ and $\text{FinSing} : \text{FinEPMet} \rightarrow \text{sFuzz}$ form an adjunction $\text{FinReal} \dashv \text{FinSing}$.*

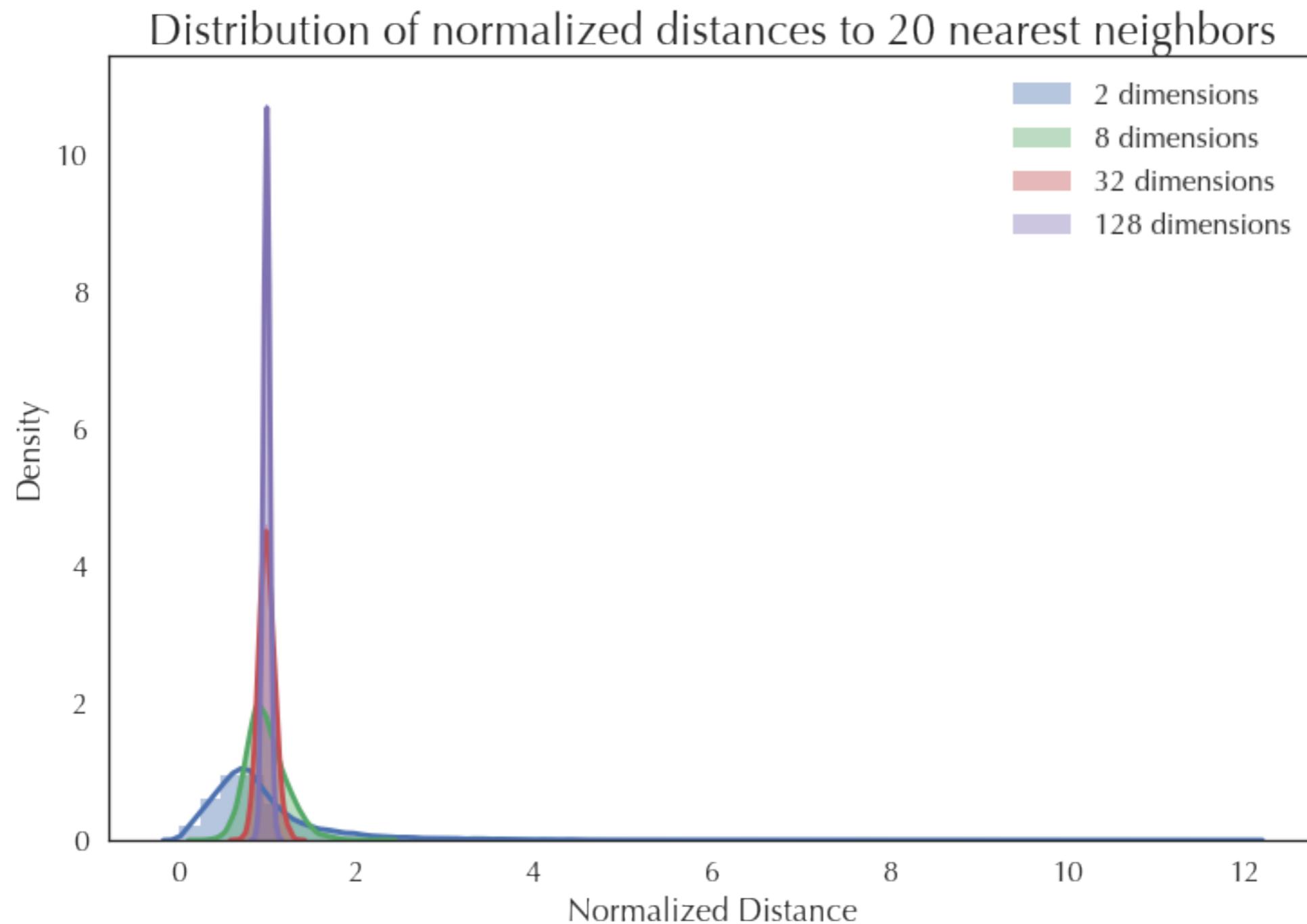


Assumption:
The manifold is locally
connected

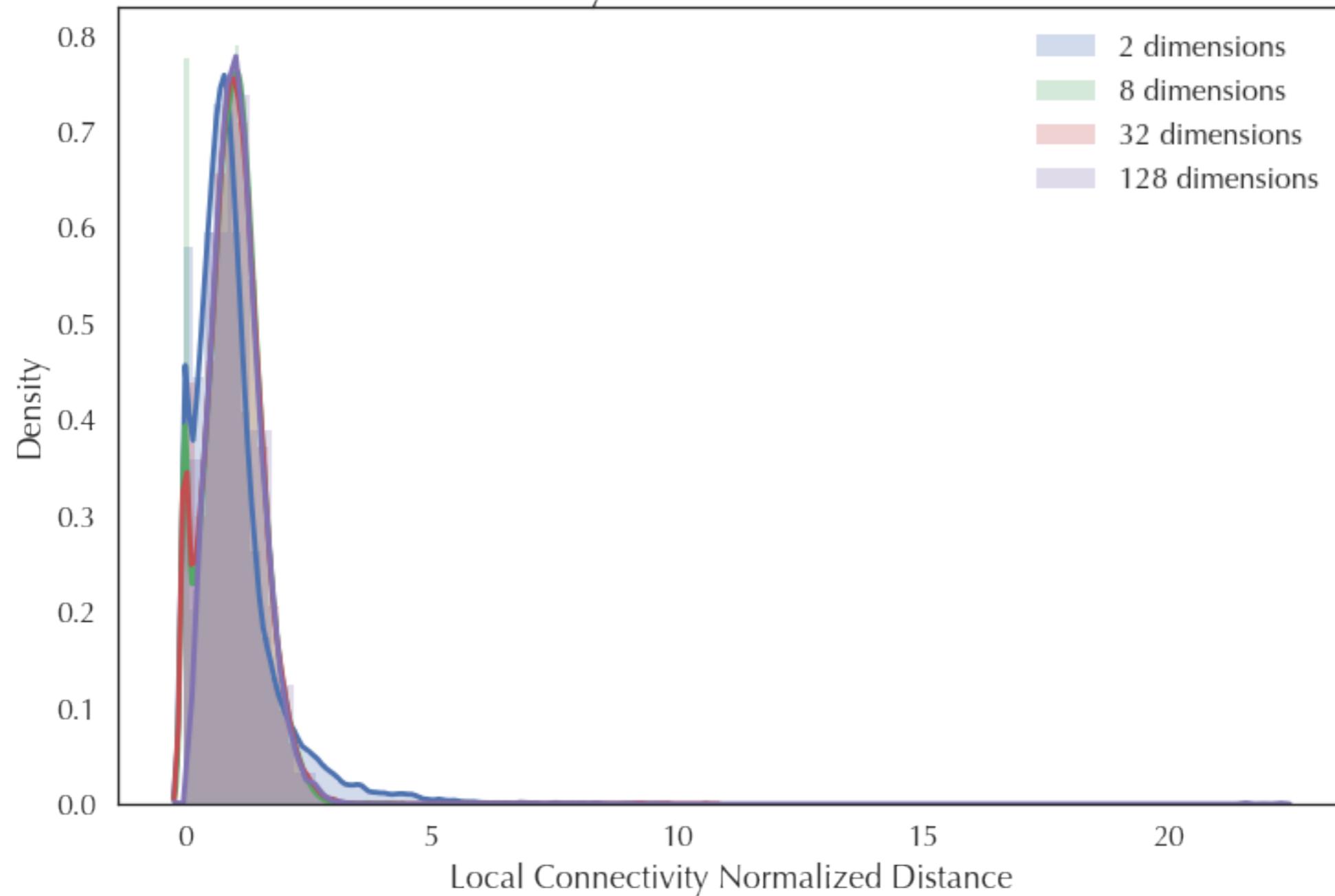


Distribution of distances to 20 nearest neighbors

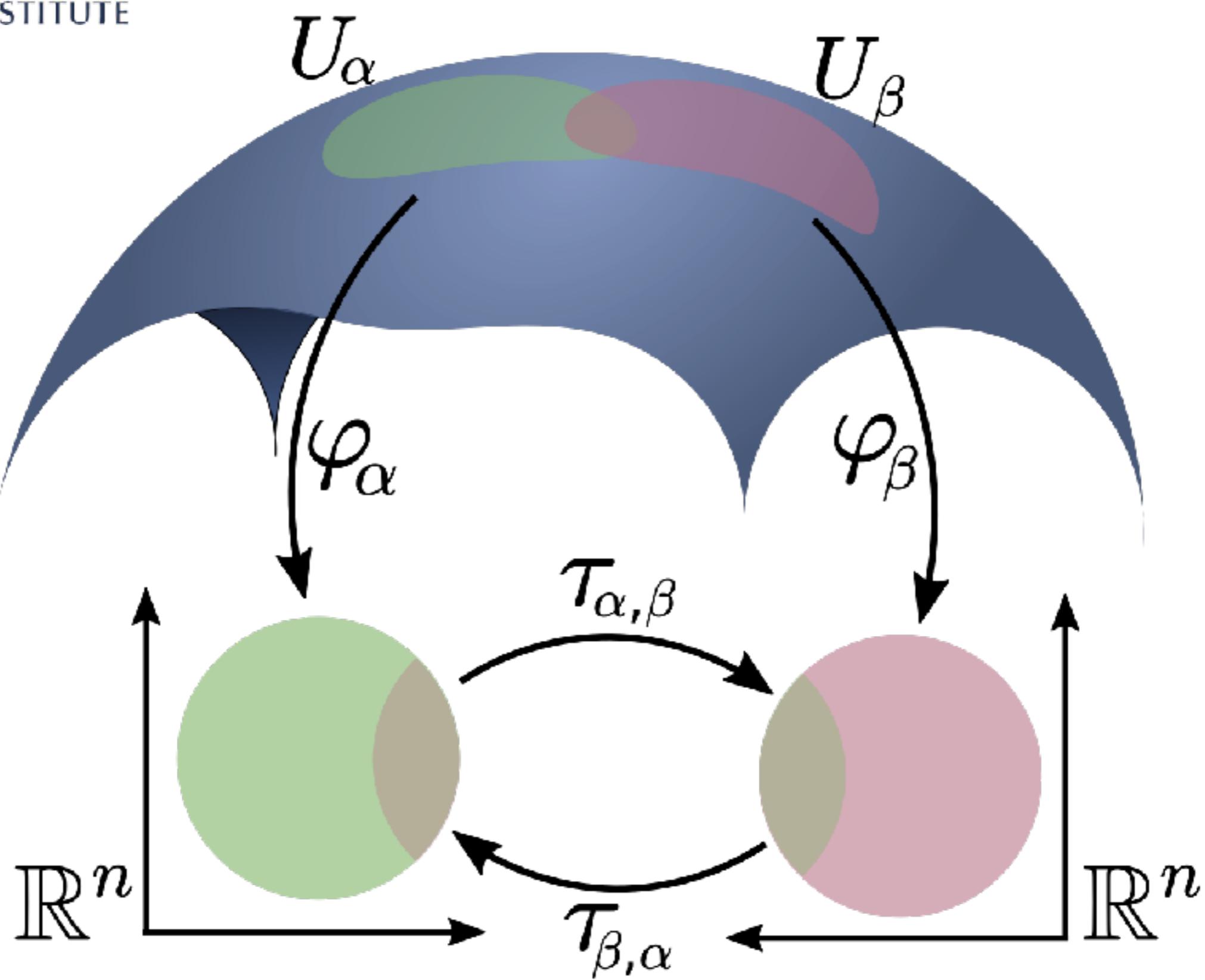


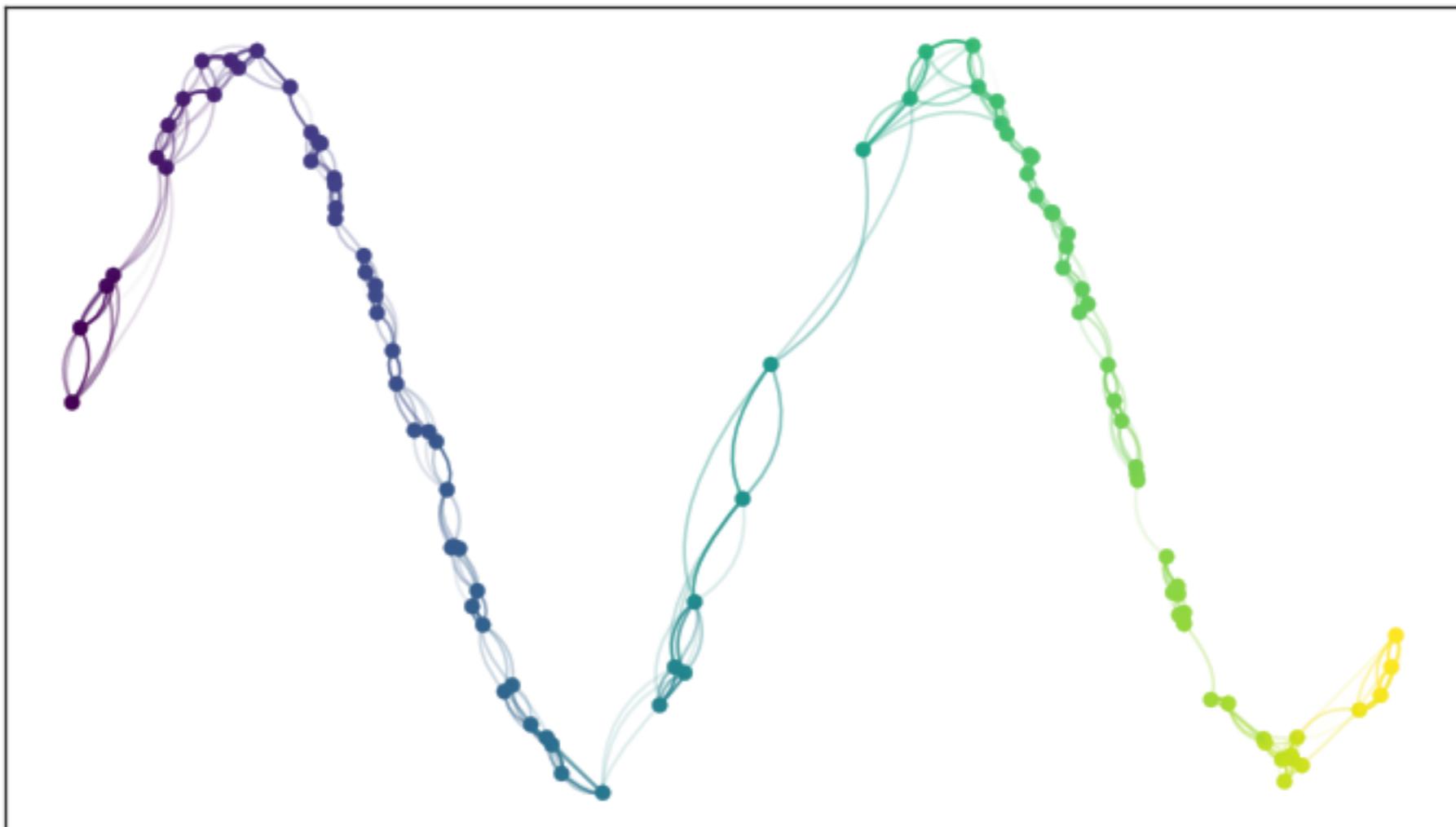


Distribution of local connectivity normalized distances to 20 nearest neighbors



But our local metrics
are all incompatible!

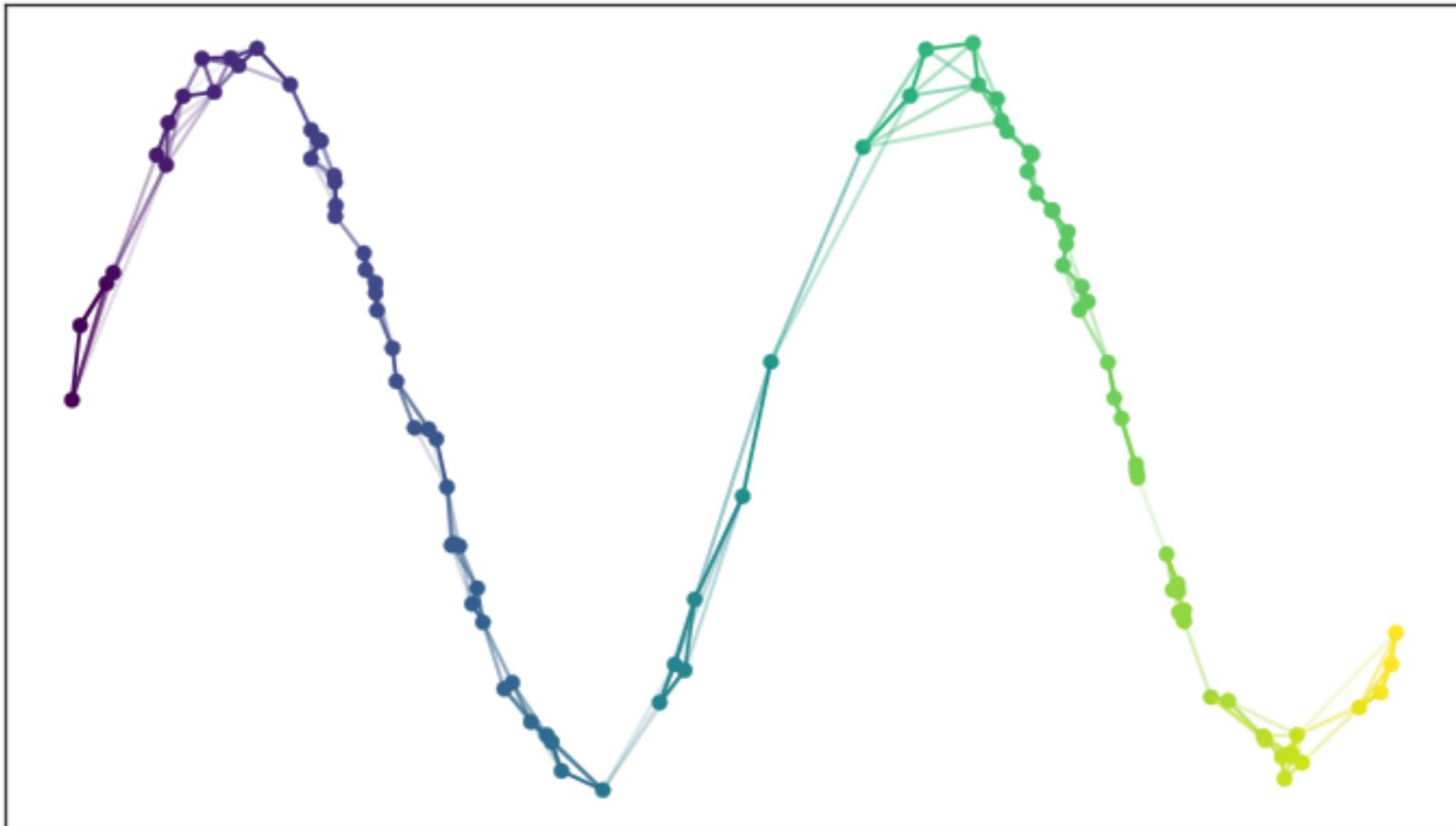




Theorem 2 (UMAP Adjunction). *The functors $\text{FinReal} : \text{sFuzz} \rightarrow \text{FinEPMet}$ and $\text{FinSing} : \text{FinEPMet} \rightarrow \text{sFuzz}$ form an adjunction $\text{FinReal} \dashv \text{FinSing}$.*

Under a probabilistic fuzzy union the combination of weights on edges is given by

$$f(\alpha, \beta) = \alpha + \beta - \alpha \cdot \beta$$



Suppose we were given
a low dimensional
representation

We can apply the same process to get a fuzzy graph!

Except we *know* the manifold, and *don't* know the “correct” nearest neighbour distance

Now measure the
distance between the
graphs using cross-
entropy and optimize

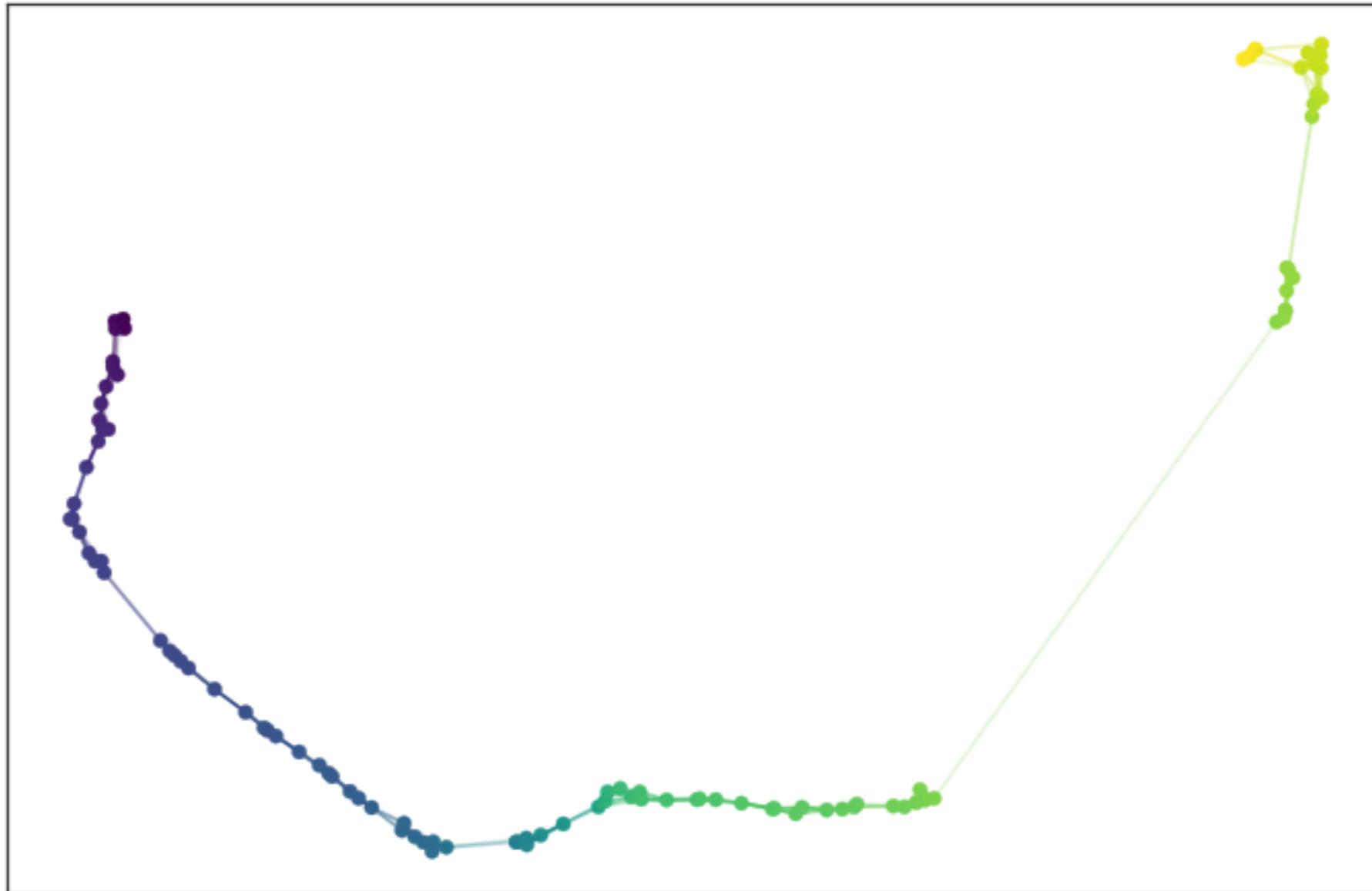
$$\sum_{a \in A} \mu(a) \log \left(\frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left(\frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

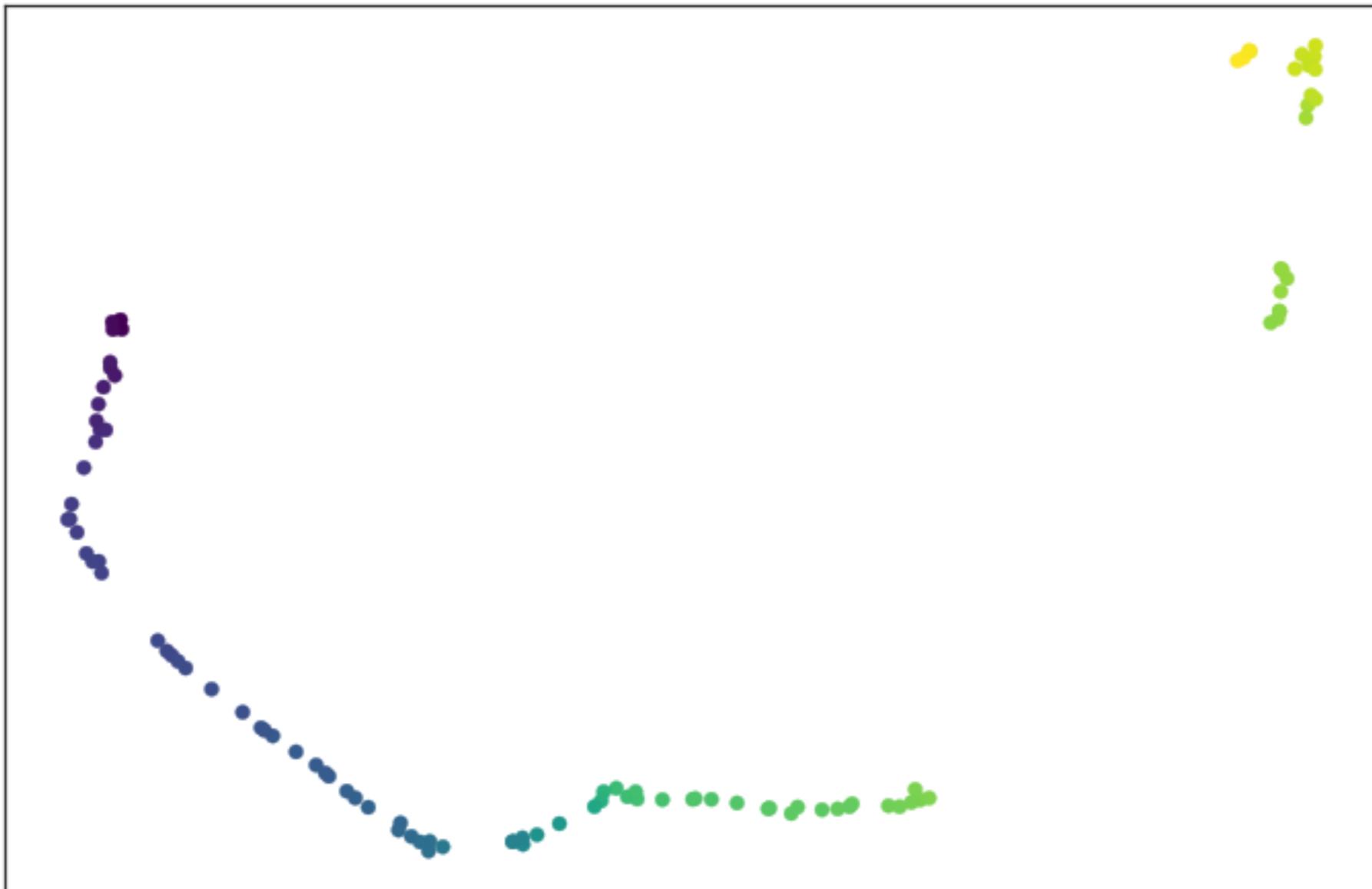
We are just embedding
the graph

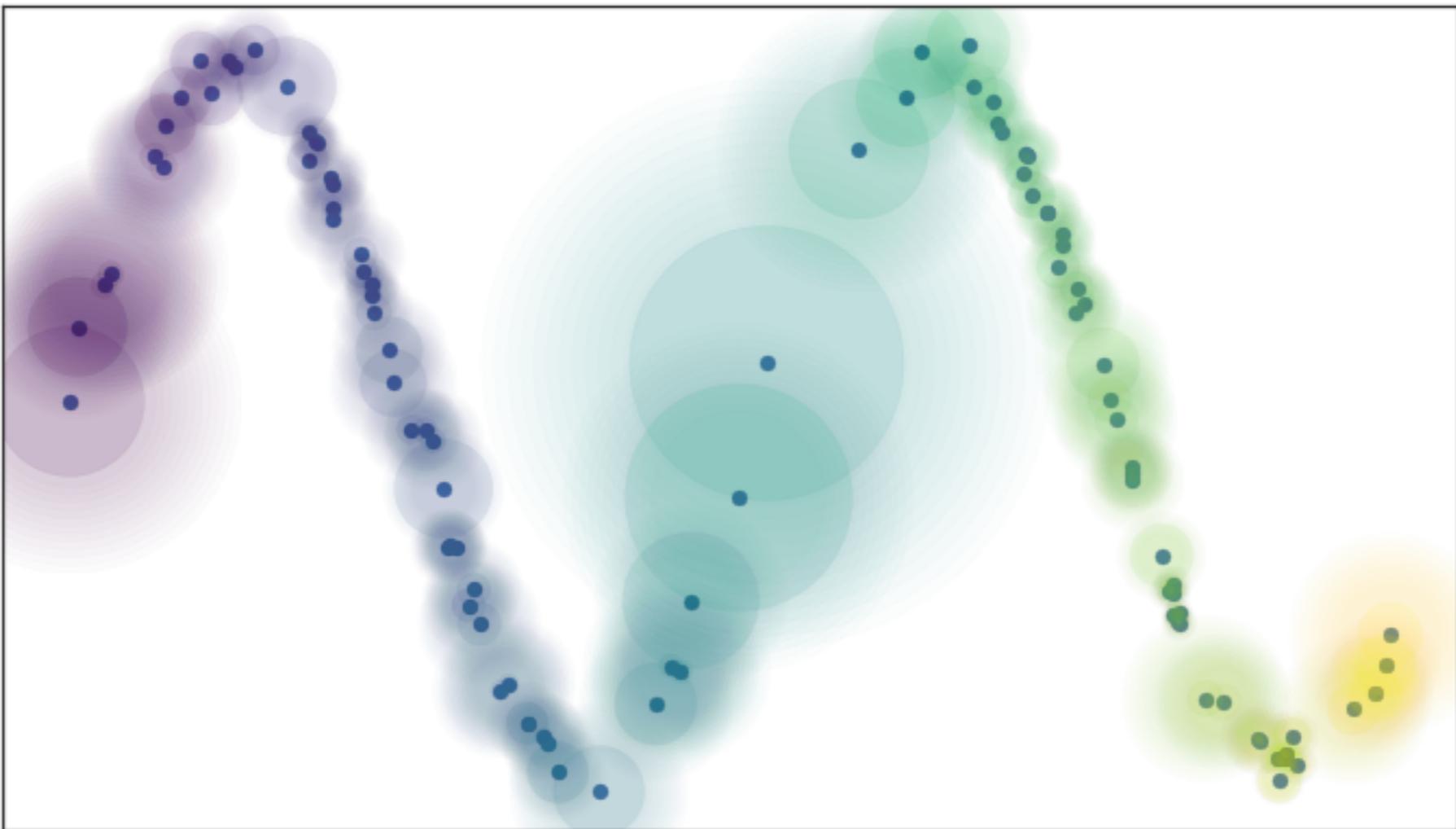
Get the clumps right

$$\sum_{a \in A} \mu(a) \log \left(\frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left(\frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

Get the gaps right

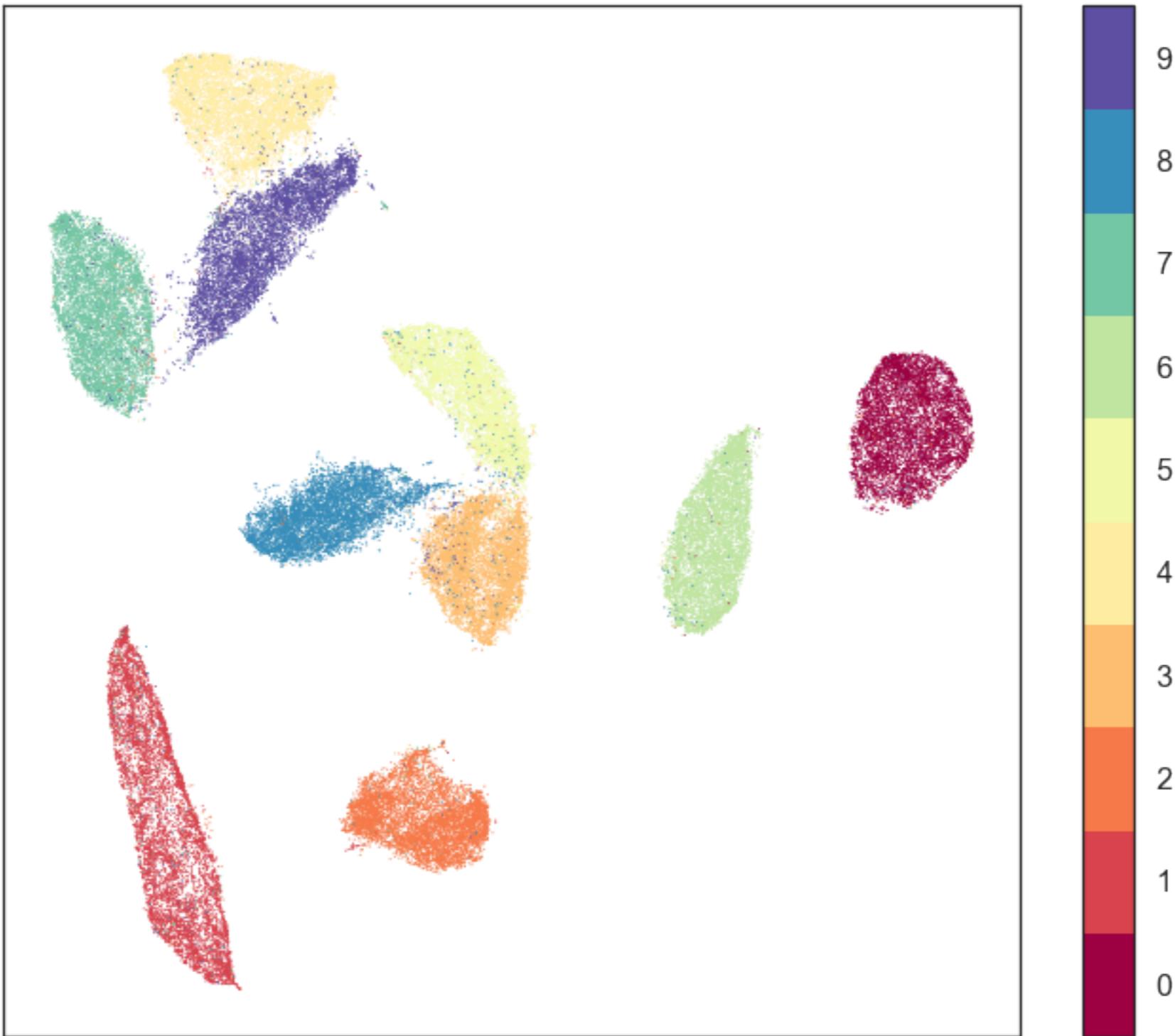




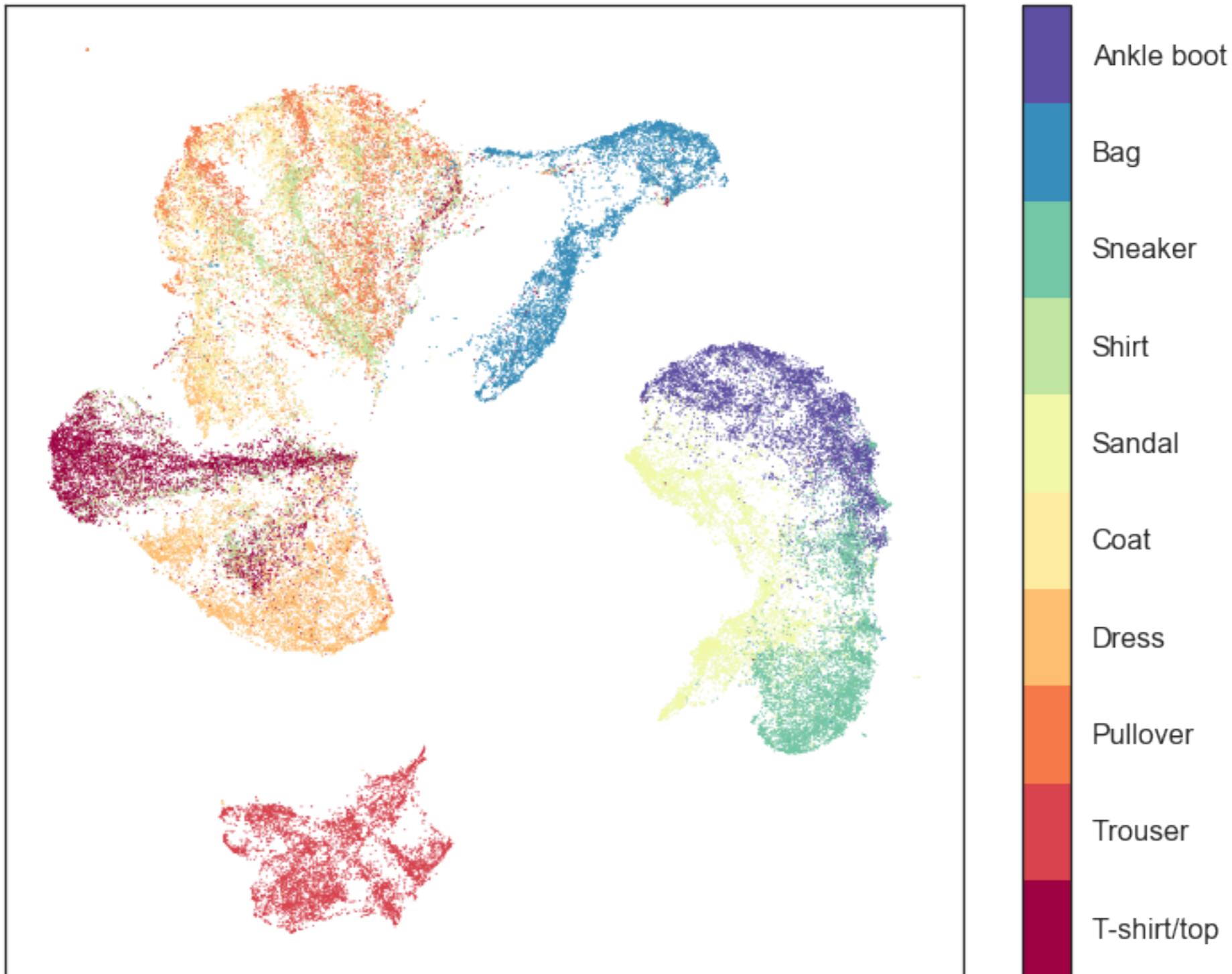


On real data?

UMAP on MNIST digits



UMAP on Fashion MNIST



Implementation

Need to find (approximate)
nearest neighbours very
efficiently

Even in high dimensional space

RP-trees + NN-descent

Dasgupta + Freund 2008

Dong, Charikar + Li 2011

Need to optimize the
layout subquadratically

SGD + negative sampling

Mikolov et al 2013
Tang et al 2016

Need to be high level
but still fast



 python™ +  Numba

Numba is awesome!

- High performance
- Clean code
- Custom distance metrics

Performance Comparison

UMAP speed up over t-SNE

COIL20

3x

MNIST

13x

Fashion MNIST

11x

GoogleNews

19x

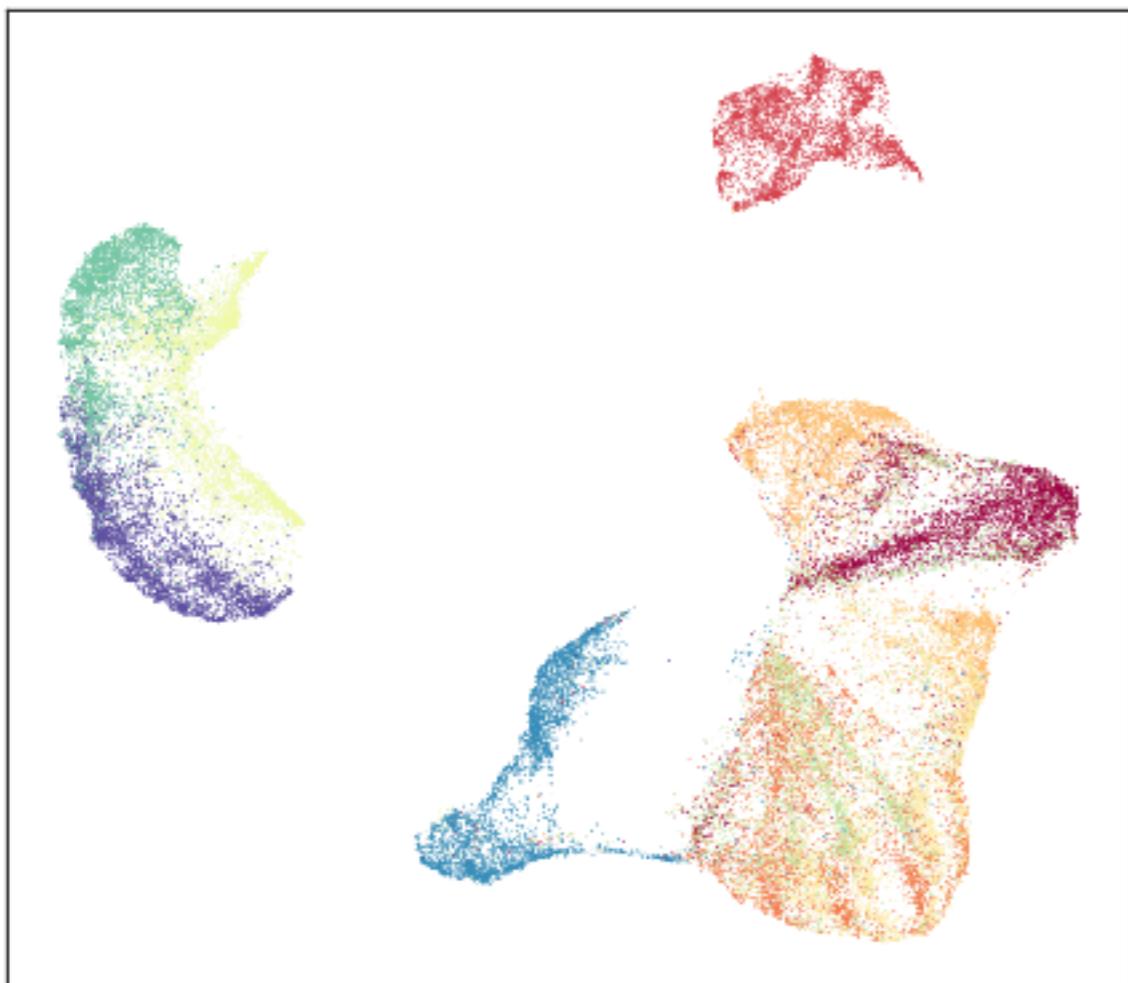


Where Next?

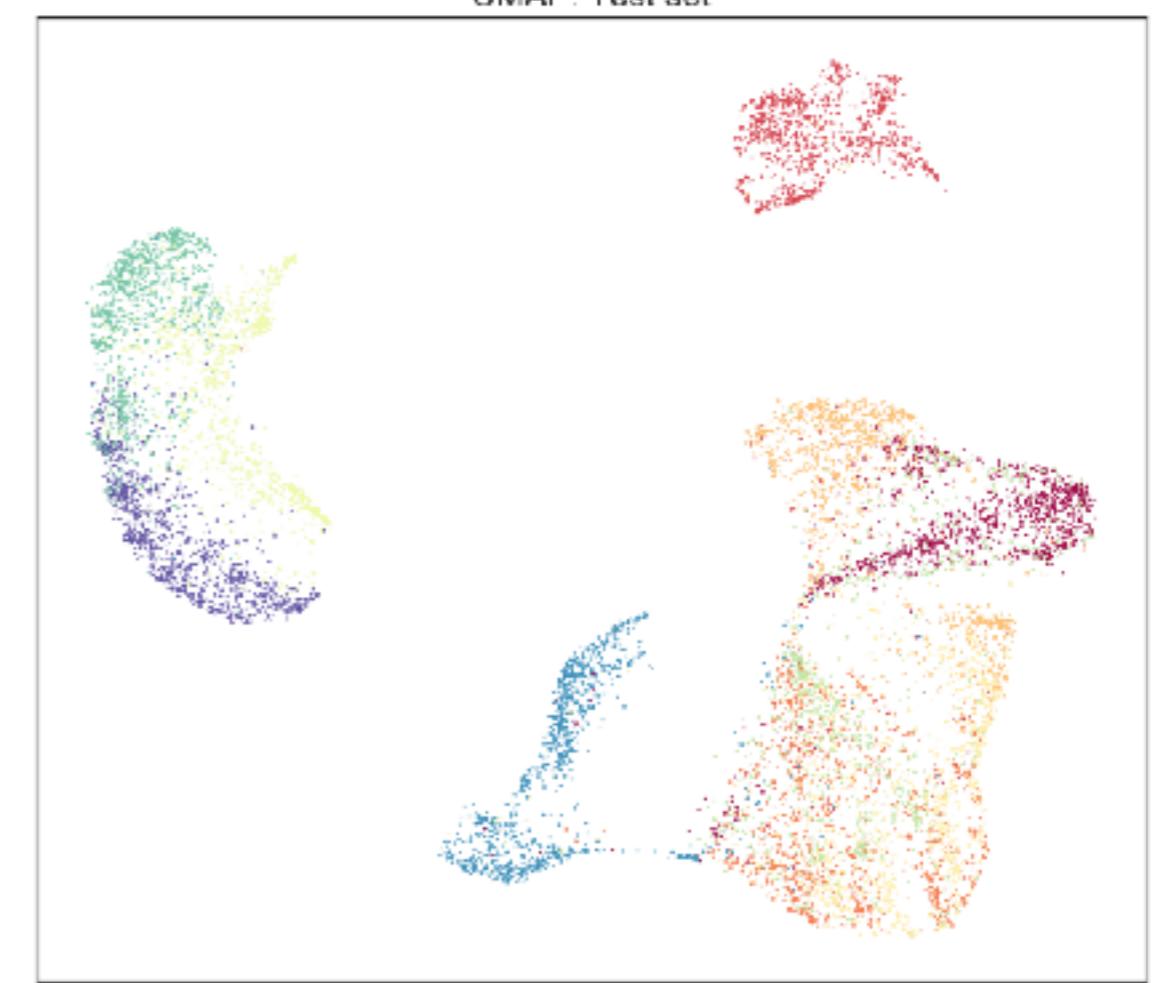
Given the mathematical foundation, a number of options are available

Embed new unseen
points into an existing
embedding

UMAP: Train set

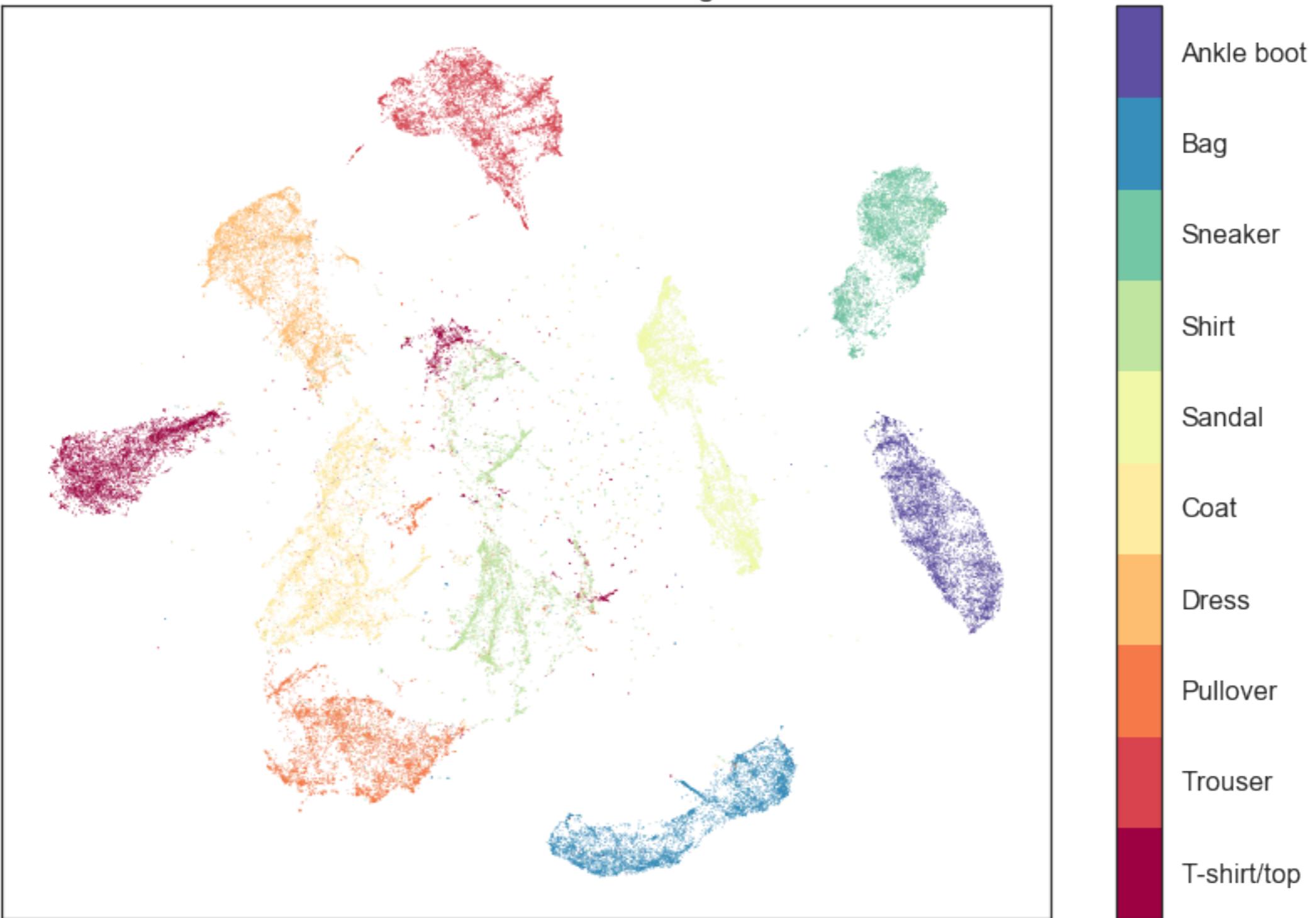


UMAP: Test set



Make use of labels for
supervised dimension
reduction

Fashion MNIST Embedded via UMAP using label information

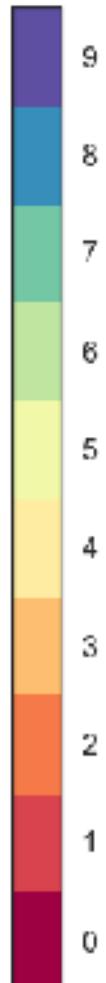
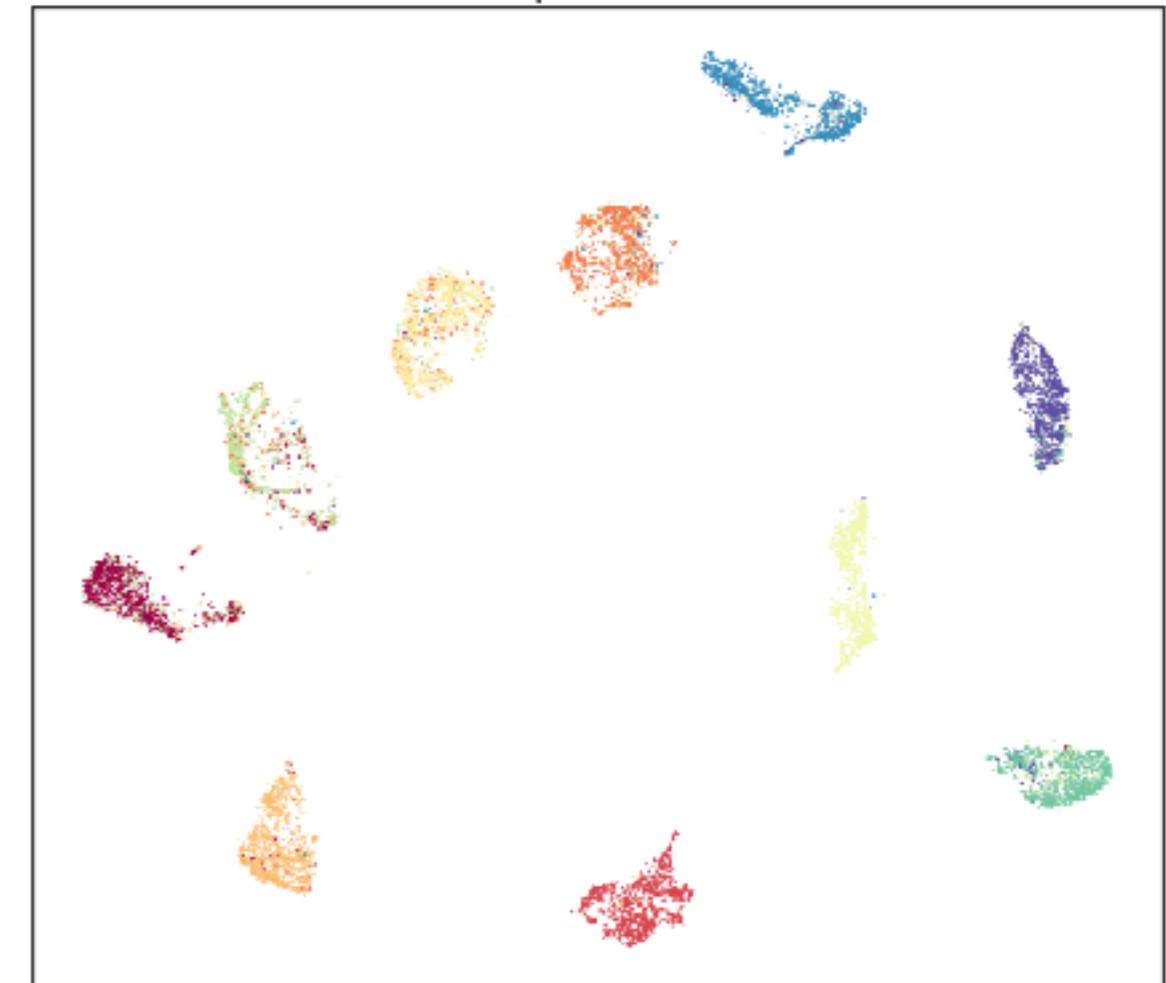


And combine those for
metric learning

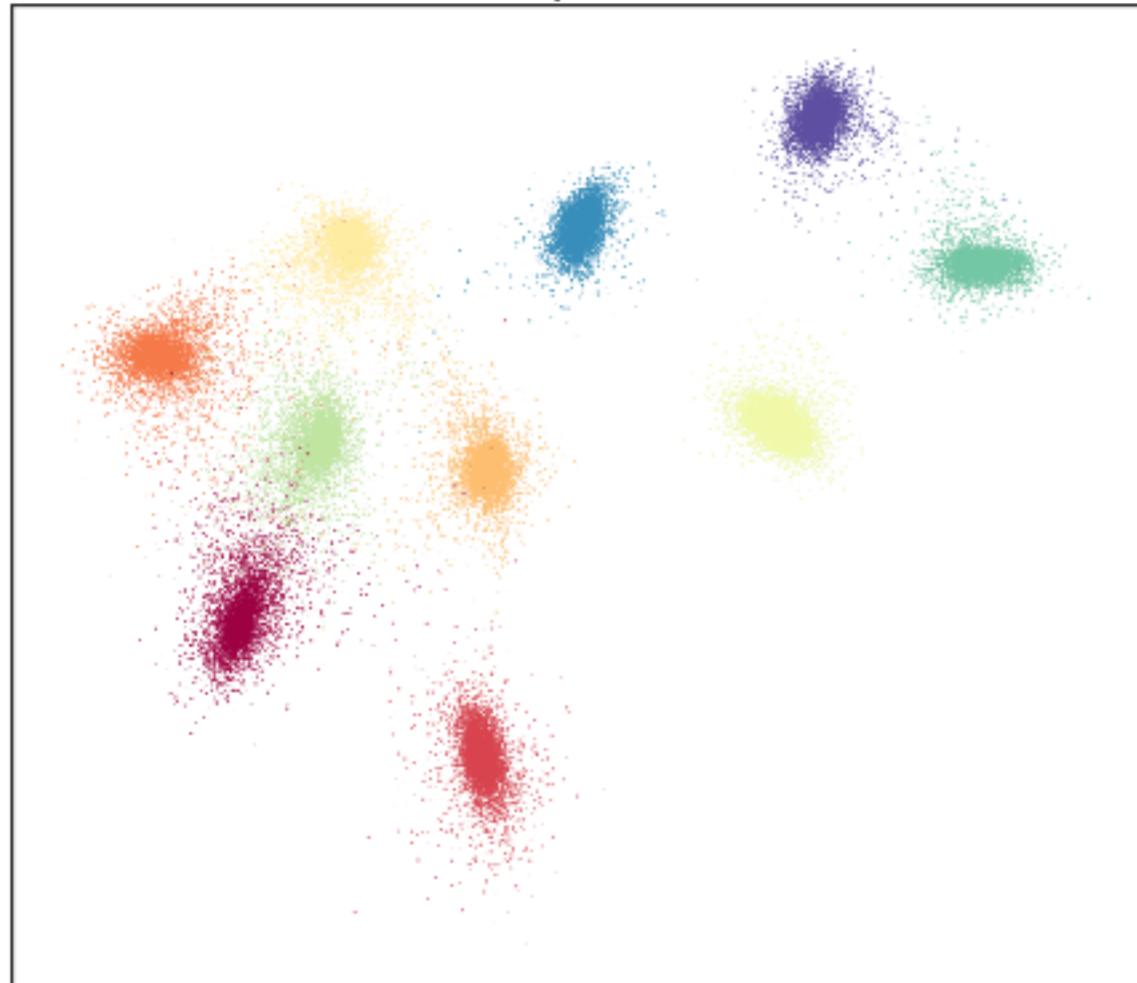
UMAP Supervised: Train set



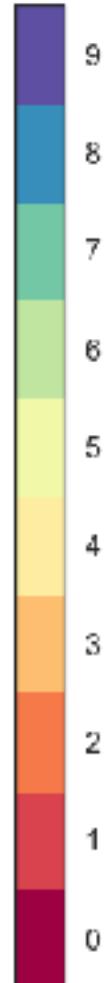
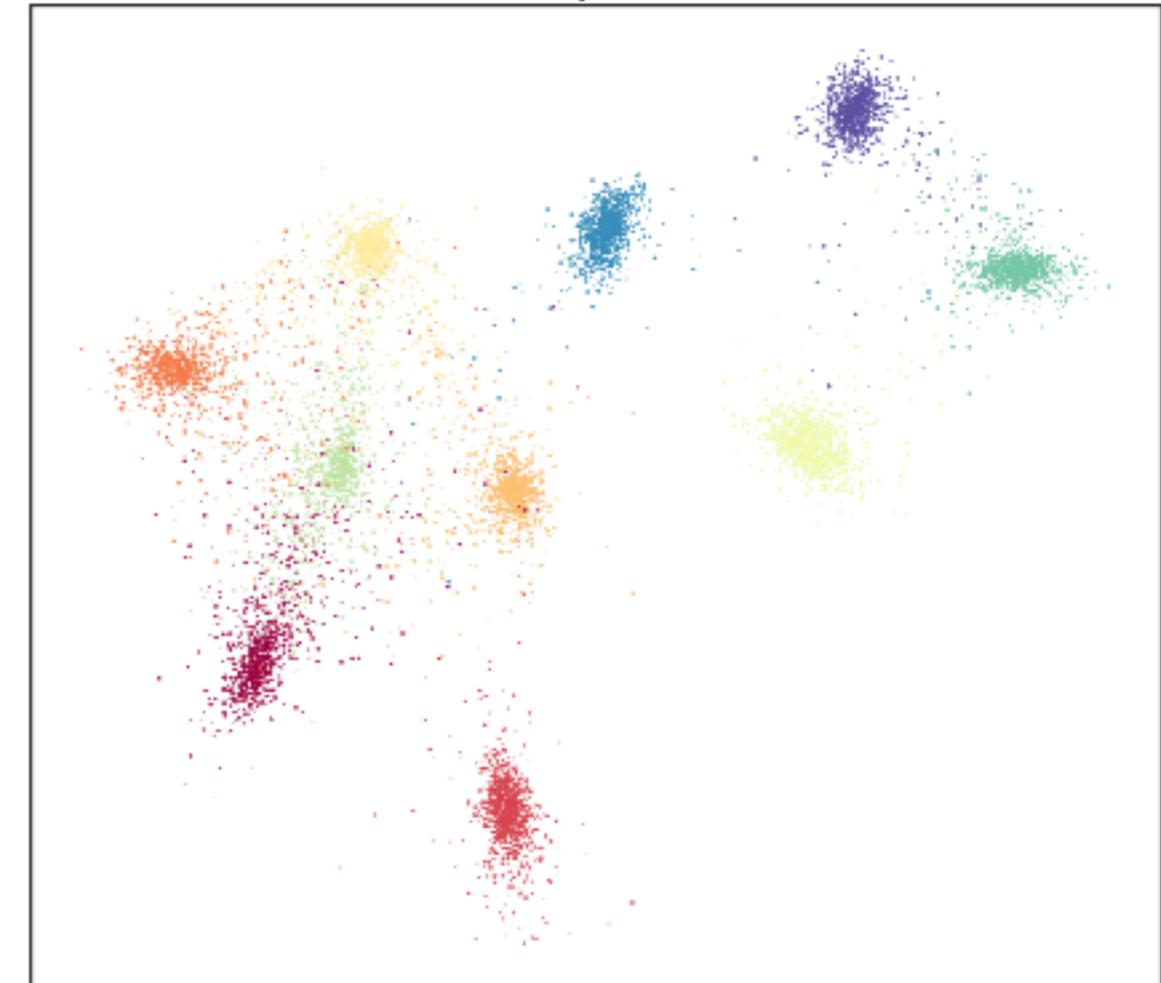
UMAP Supervised: Test set



Online Selection Triplet Network: Train set



Online Selection Triplet Network: Test set



Derived from Adam Bielski's Siamese/Triplet repository:
<https://github.com/adambielski/siamese-triplet>

Adding one categorical
variable is no harder
theoretically than
adding many

Combine spaces with different metrics

Continuous, categorical,
ordinal, Haversine, Levenshtein,
and more ...

As long as you can
provide a metric for the
datatype UMAP can
combine it with other
datatypes!

UMAP for pandas dataframes!

<https://github.com/lmcinnes/umap>

conda install -c conda-forge umap-learn

pip install umap-learn



leland.mcinnes@gmail.com



[@leland_mcinnes](https://twitter.com/leland_mcinnes)