# TOPOLOGICAL DATA ANALYSIS

Joshua Sheldon, Justin Barnwell, Michelle Arubi
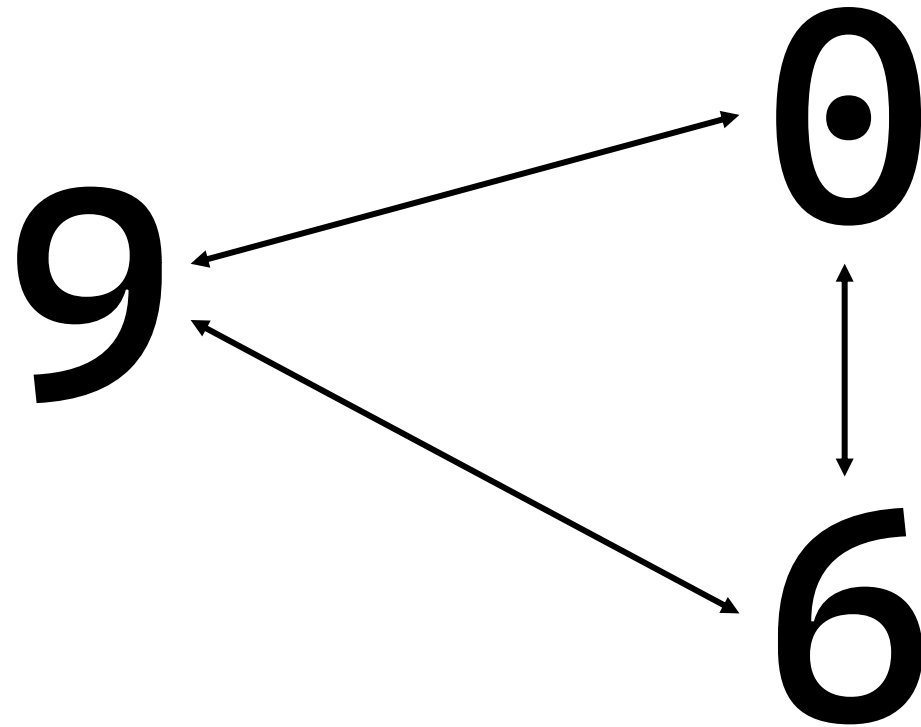
# OUTLINE

- Intro to TDA

- Data & Objective

- Implementation & Results
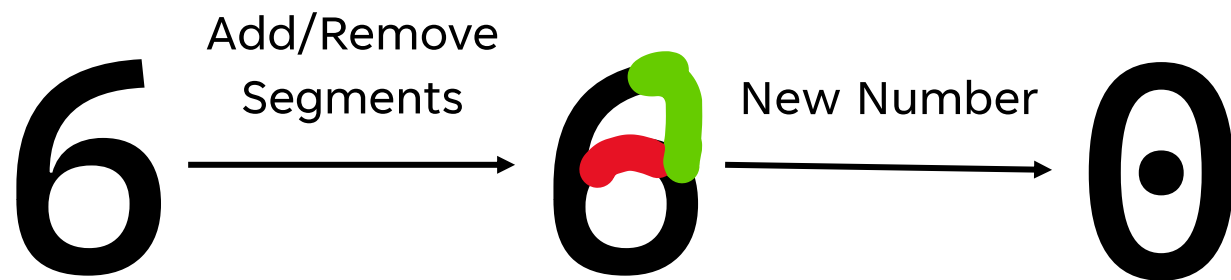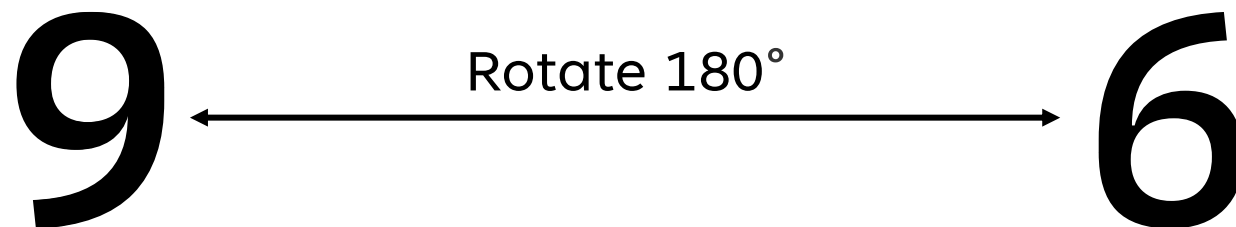
# INTRO TO TDA

The science of shapes

# SHAPING UP

# TRANSFORMERS
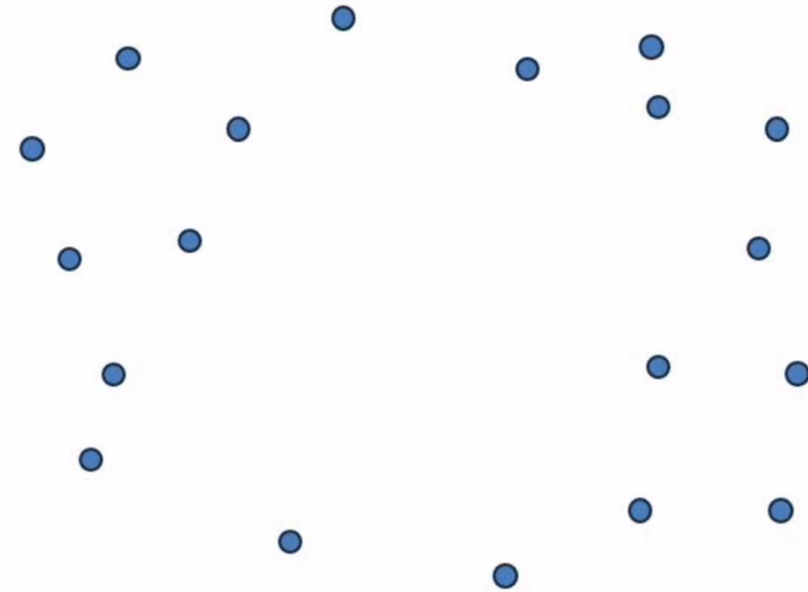
9 ← Rotate 180° → 6



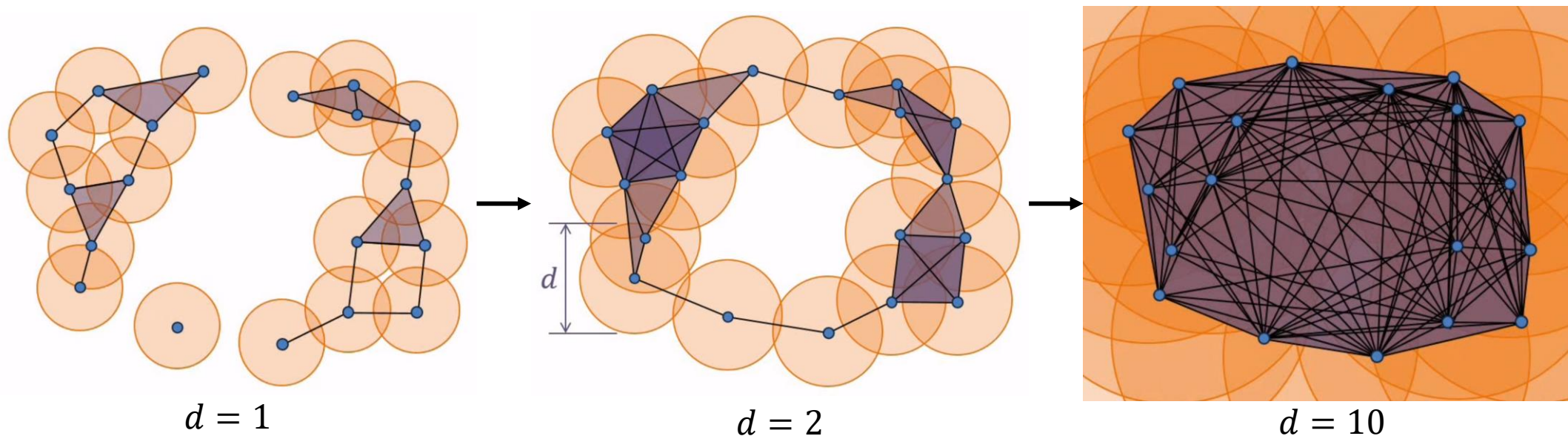6 → Add/Remove Segments →  → New Number → 0

# TOPOLOGICAL DATA ANALYSIS

**to·pol·o·gy**
1. The way in which constituent parts are interrelated or arranged
2. The study of geometric properties and spatial relations unaffected by the continuous change of shape or size of figures

## What shape?

# PERSISTENT HOMOLOGY
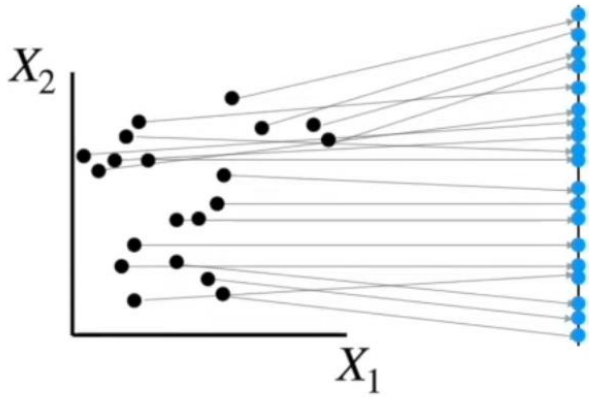


$d = 1$           $d = 2$           $d = 10$

- Middle hole may last from $d = [2, 8]$
- **Persistence** = $d_{end} - d_{start}$
- Higher persistence = feature, lower persistence = noise

# NOW WHAT?

- Point clouds = data sets

- Using homology, we acquire **important features** of our data.

- We can reduce the **dimensionality** of data while maintaining **important features**.

- Reduced dimensionality makes **data analysis** possible!
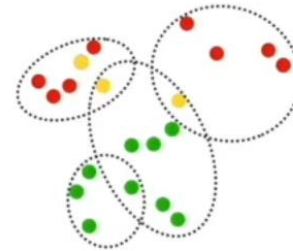
# MAPPER ALGORITHM



1) Data     2) Project data

3) Cover
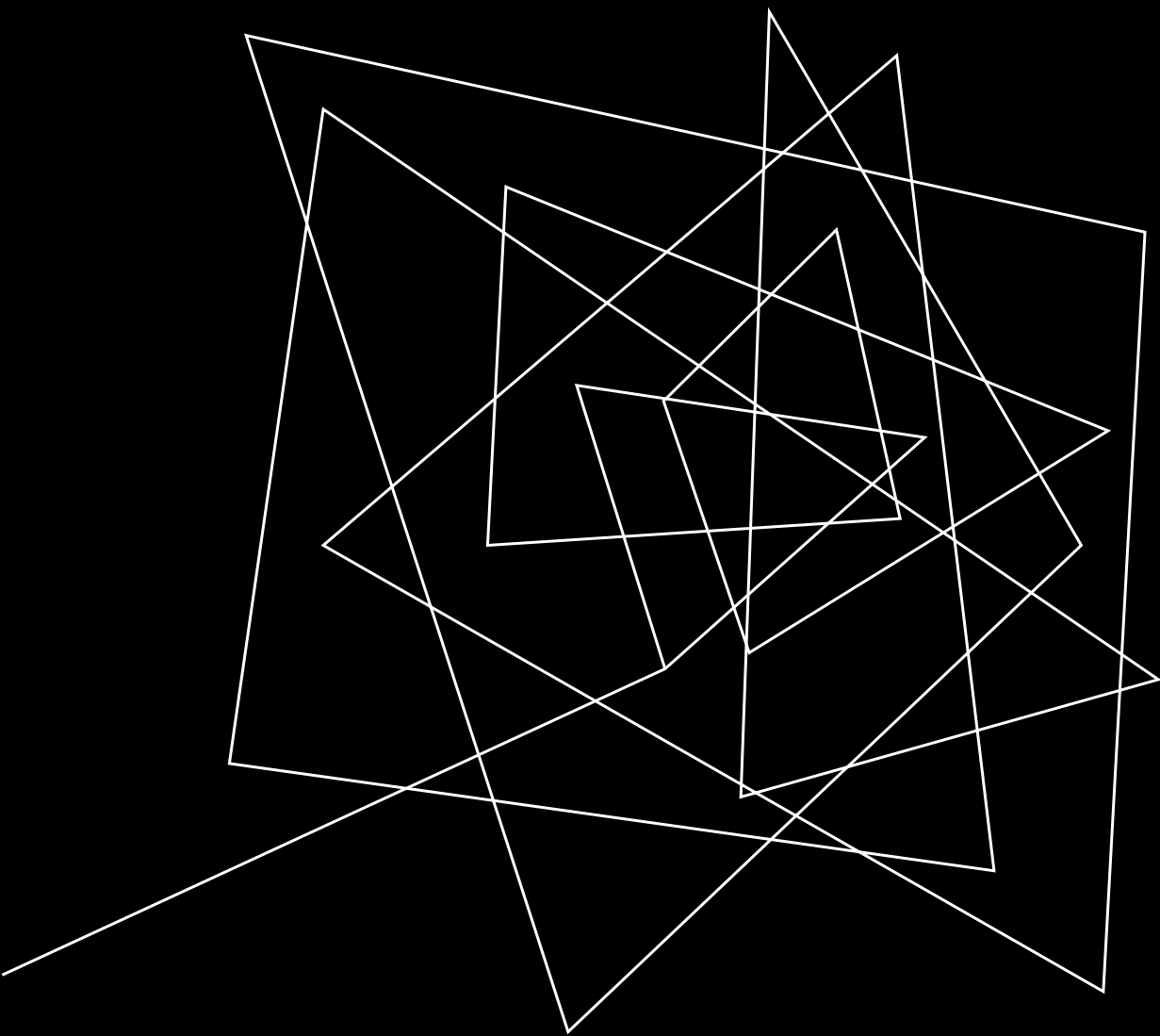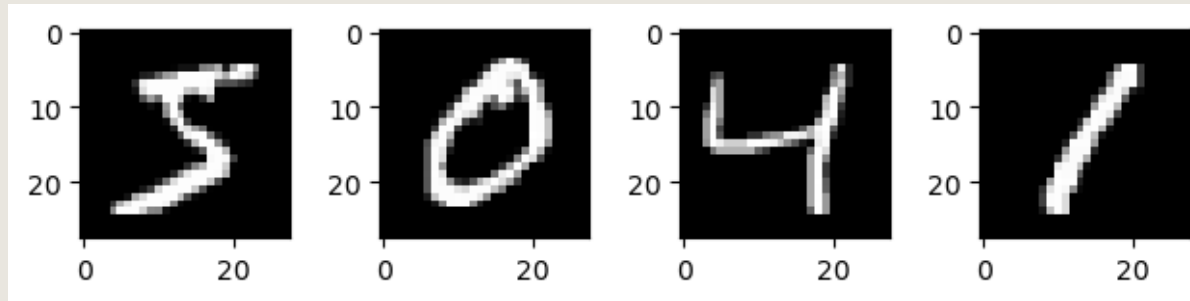
4) Cluster Pre-image

5) Graph output

**Nodes** = clusters
**Edges** = clusters share members
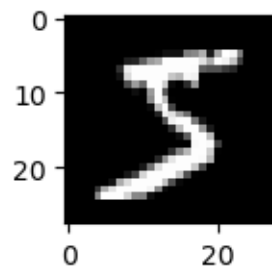
# DATA &
# OBJECTIVE

Digit data with detracted dimensionality

# MNIST



- Library of grasycale handwritten digit images.

# MNIST DIGIT EXAMPLE



28x28 matrix of
numbers [0, 255]
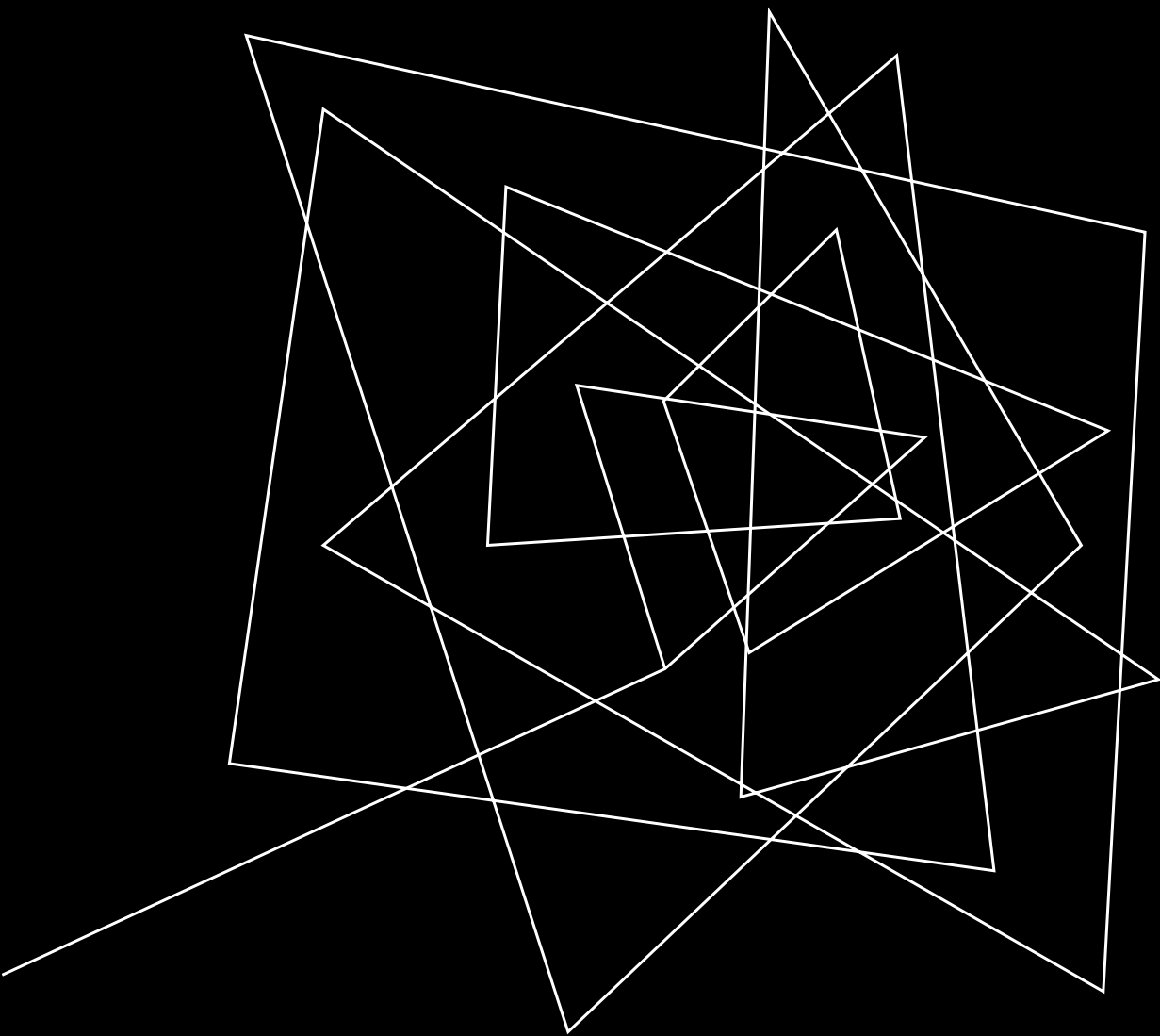(8 bit unsigned ints)

0                                    255

# OUR OBJECTIVE

- Select 10 MNIST images for each digit (100 total)

- Analyze them with the Mapper algorithm

- Visualize the relationships between digits

# IMPLEMENTATION & RESULTS

You're in the weeds: beware Pythons

# TECH STACK

# IMPORTING MNIST



```python
# Demonstrating how the x array contains the numbers,
# and the y array contains the Labels
for i in range(4):
    print(train_y[i])
    plt.subplot(330 + 1 + i)
    plt.imshow(train_x[i], cmap=plt.get_cmap('gray'))
    plt.show()
```

- Keras provides the MNIST dataset through 4 arrays:
  - x_train : (60000,28,28)
  - y_train : (60000,)
  - x_test  : (60000,28,28)
  - y_test  : (60000,)

- x = image data

- y = labels

- Converting 28x28 matrices to 784-dimension vectors
- Three new data structures
  - `data    : array(100,784)`
  - `labels : array(100,)`

# SELECTION ALGORITHM

```
# Keep track of how many of each digit we've collected
added = new size 10 array, initialized to 0s

# Select digits
old_index = 0 # For traversing through Keras MNIST arrays
new_index = 0 # For traversing through our arrays
while not all(value == 10 for value in added):
    digit = train_y[old_index]

    if added[digit] < 10:
        # Add label in format: <digit> (#<occurrence>)
        labels.append(label)

        # Reduce dimensionality and add to array
        data[new_index] = train_x[old_index].reshape(-1)

        # Increment occurrences of digit and position in data array
        added[digit] += 1
        new_index += 1

    old_index += 1
```
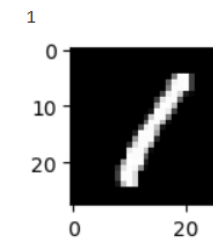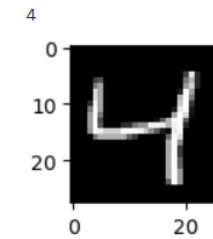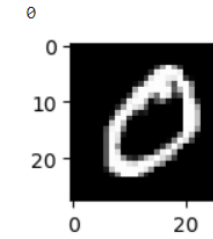
# REDUCING DIMENSIONALITY
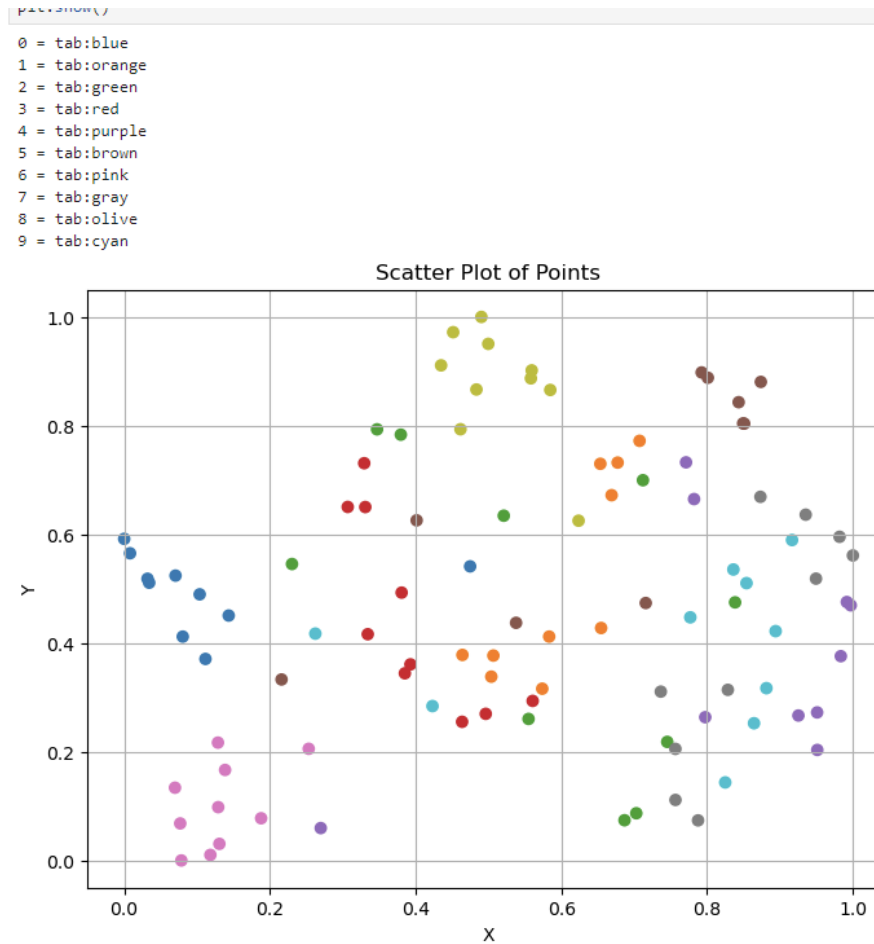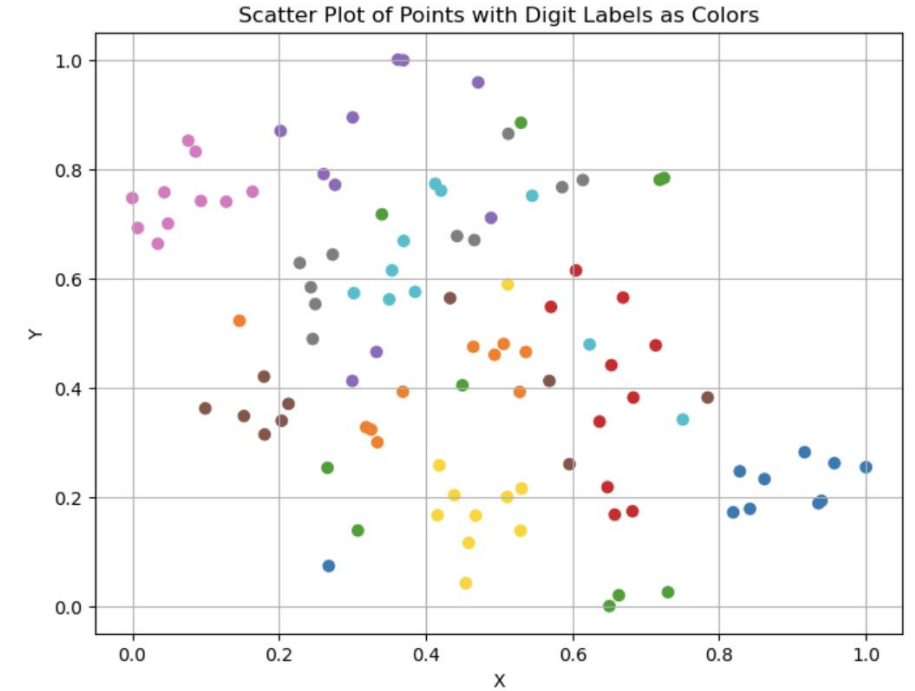
- Now we apply a series of **projections** to reduce our data from 784 dimensions to 2.
- Example:
  - Isometric Mapping – 784 -> 100
  - UMAP* – 100 -> 2
- Other popular projections: PCA, t-SNE, Feature Scaling
- Each projection may have different strengths/objectives.

# PROJECTION AS A LENS



- Pipeline 1
  - Isometric Mapping
  - UMAP
- Pipeline 2
  - MinMax Scaler
  - T-SNE

# COVER, CLUSTER, AND GRAPH

- Now that our data is in 2 dimensions, we can create a cover, cluster the data, and construct a graph.
- Will refer to these operations as **mapping**.
- We tried two types of mapping: **informed** and **uninformed**.
- The difference between these two is knowledge of the digit that each data point represents.

# MANUAL MAPPING

## Step 1: Create concave hulls around all points of a digit.



Scatter Plot of Points with Digit Labels as Colors

For 0



Scatter Plot of Points with Digit Labels as Colors

## For 1



Scatter Plot of Points with Digit Labels as Colors

For 2



Scatter Plot of Points with Digit Labels as Colors

For 3

For 4



Scatter Plot of Points with Digit Labels as Colors

For 5



Scatter Plot of Points with Digit Labels as Colors

For 6



Scatter Plot of Points with Digit Labels as Colors

## For 7



Scatter Plot of Points with Digit Labels as Colors

For 8

# MANUAL MAPPING

For 9



Scatter Plot of Points with Digit Labels as Colors

# MANUAL MAPPING

- Cluster for Digit $x$ – concave hull created from all the data points of digit $x$
- Node for Digit $x$ – contains all data points within the cluster for digit $x$, even if some data points aren't of digit $x$
- Edge – formed between two nodes when a data point is in both nodes
- May be improved by calculating edges by cluster overlap instead of data point overlap.
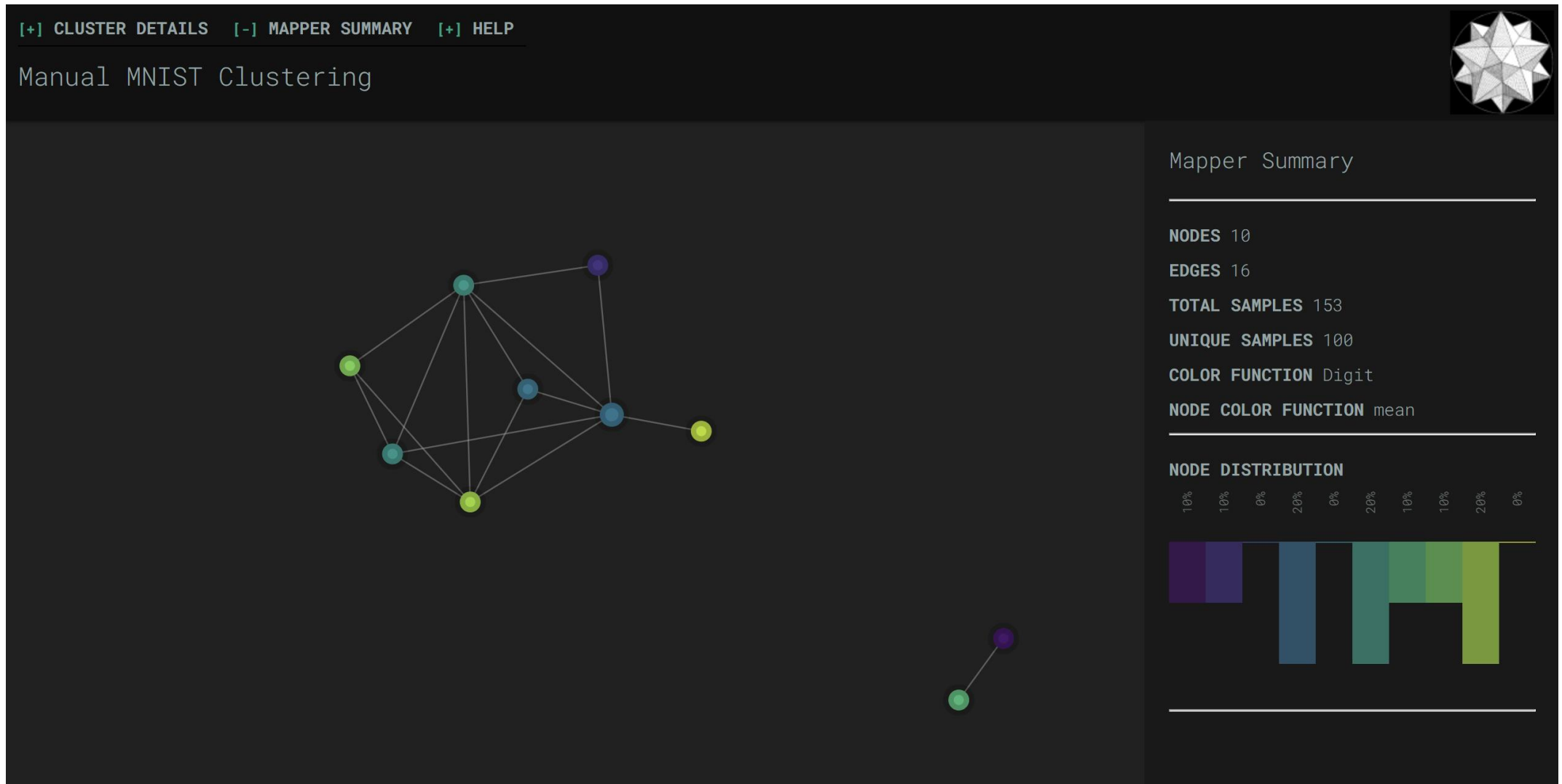
# MANUAL MAPPING

Manual MNIST Clustering

## Mapper Summary

**NODES** 10

**EDGES** 16

**TOTAL SAMPLES** 153

**UNIQUE SAMPLES** 100

**COLOR FUNCTION** Digit

**NODE COLOR FUNCTION** mean

**NODE DISTRIBUTION**

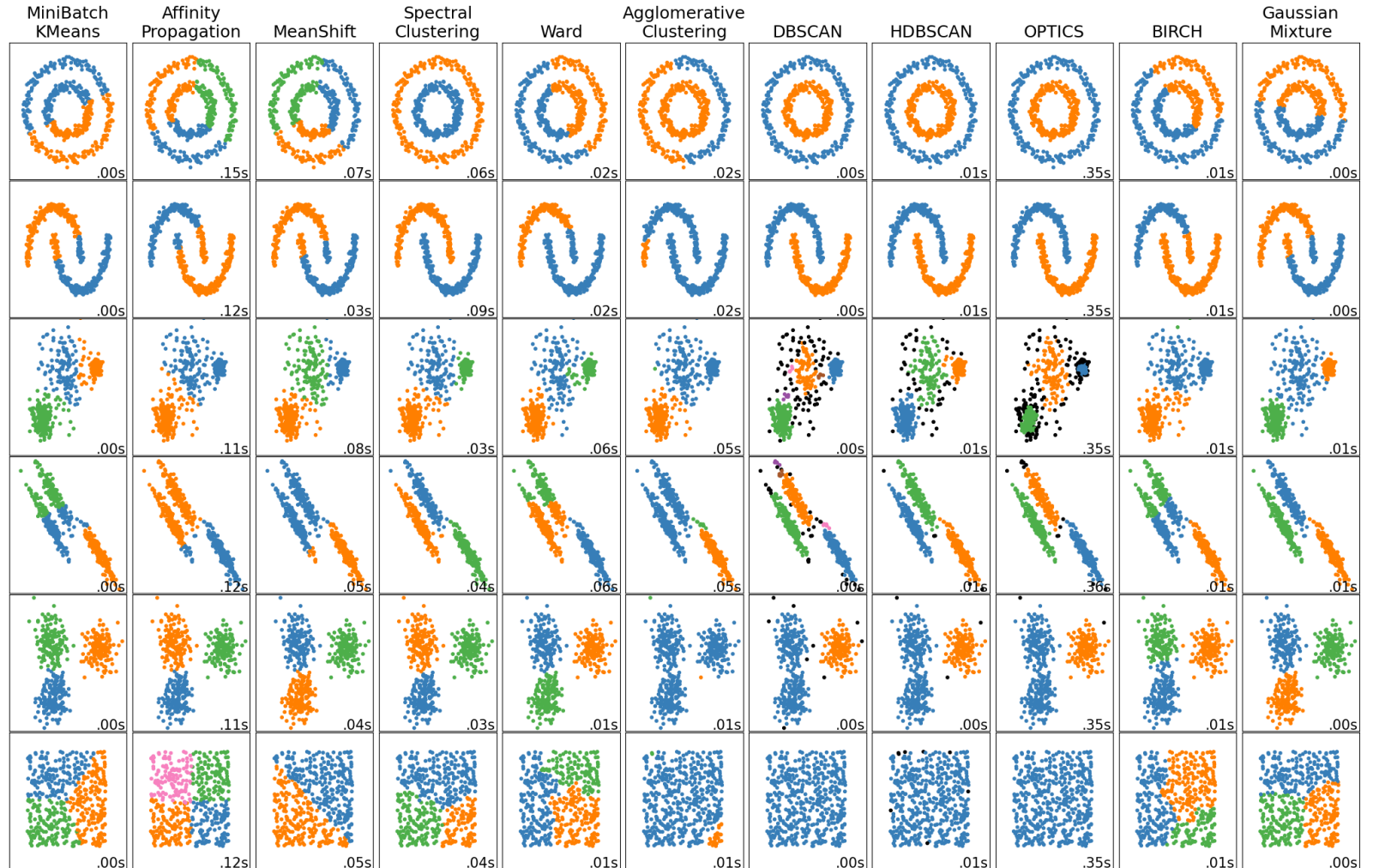10%   10%   0%   20%   0%   20%   10%   10%   20%   0%
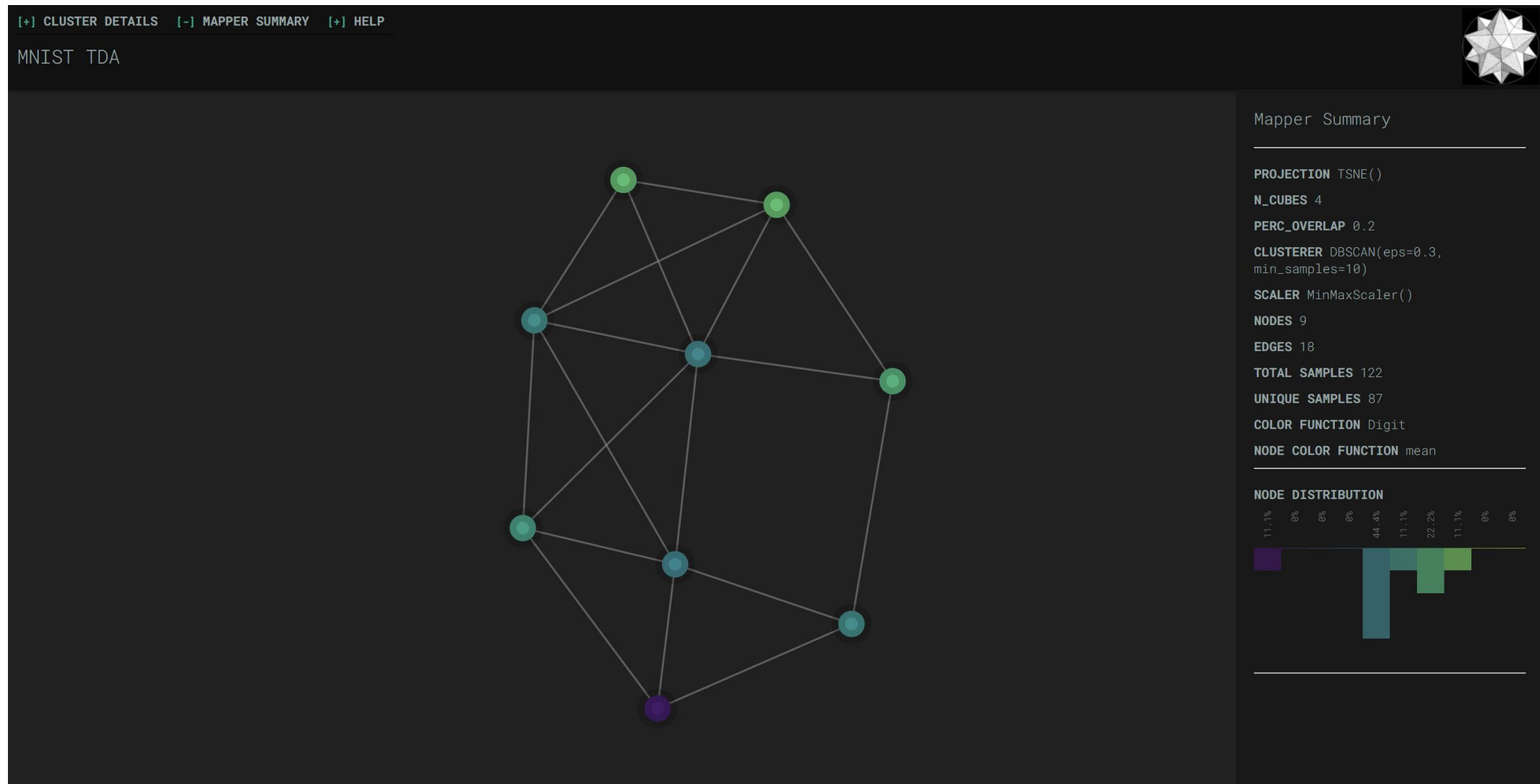
# AUTOMATIC MAPPING

- KeplerMapper provides covering scheme
  - Each dimension spanned by **hypercubes**
  - Adjacent cubes have **% of overlap**
  - Hypercubes per dimension and % of overlap are parameters

- Many different clustering algorithms, each with different strengths and use cases.
- **DBSCAN** worked best for us.
- **Spectral Clustering** looks promising.
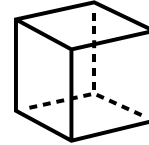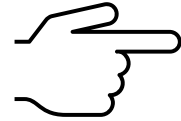
# AUTOMATIC MAPPING

# DATA PIPELINE

**Selection    Reshaping**

Raw Keras MNIST Data

```
x_train = array(60000,28,28)
y_train = array(60000,)
```
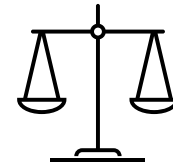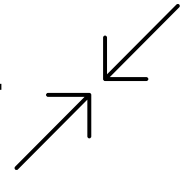
Selected & Reshaped Data

```
data = array(100,784)
labels = array(100,)
```
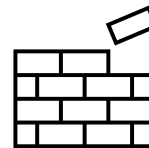
**Dimension
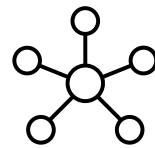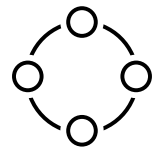Reduction    Scaling**

Projected Data

```
data = array(100,784)
labels = array(100,)
projected_data = array(100,2)
```

**Graph
Creation**

**Cover    Clustering**

Graph

```
data = array(100,784)
labels = array(100,)
projected_data = array(100,2)
graph = dict(nodes,links,…)
```

# INTERESTING FINDINGS

- 0 and 6 typically close on the scatter plot but no overlap

- 9 has large cluster (similar topology to many digits) and almost never overlap with 0 or 6

- 5 and 4 typically close to 9 on the scatter plot

- Both 1 and 3 typically within 2's cluster

- 9 and 7 typically close on the scatter plot

- **All of this could change with different/more data, but consistent across lenses/projection pipelines!**

# THANK YOU



*Interactive HTML files,
Jupyter Notebook w/ code,
Project proposal,
Presentation file, etc.*