# Clustering Data

## A Guide for the Perplexed

Leland McInnes
John Healy



leland.mcinnes@gmail.com
jchealy@gmail.com

# Find groups of data that are all similar
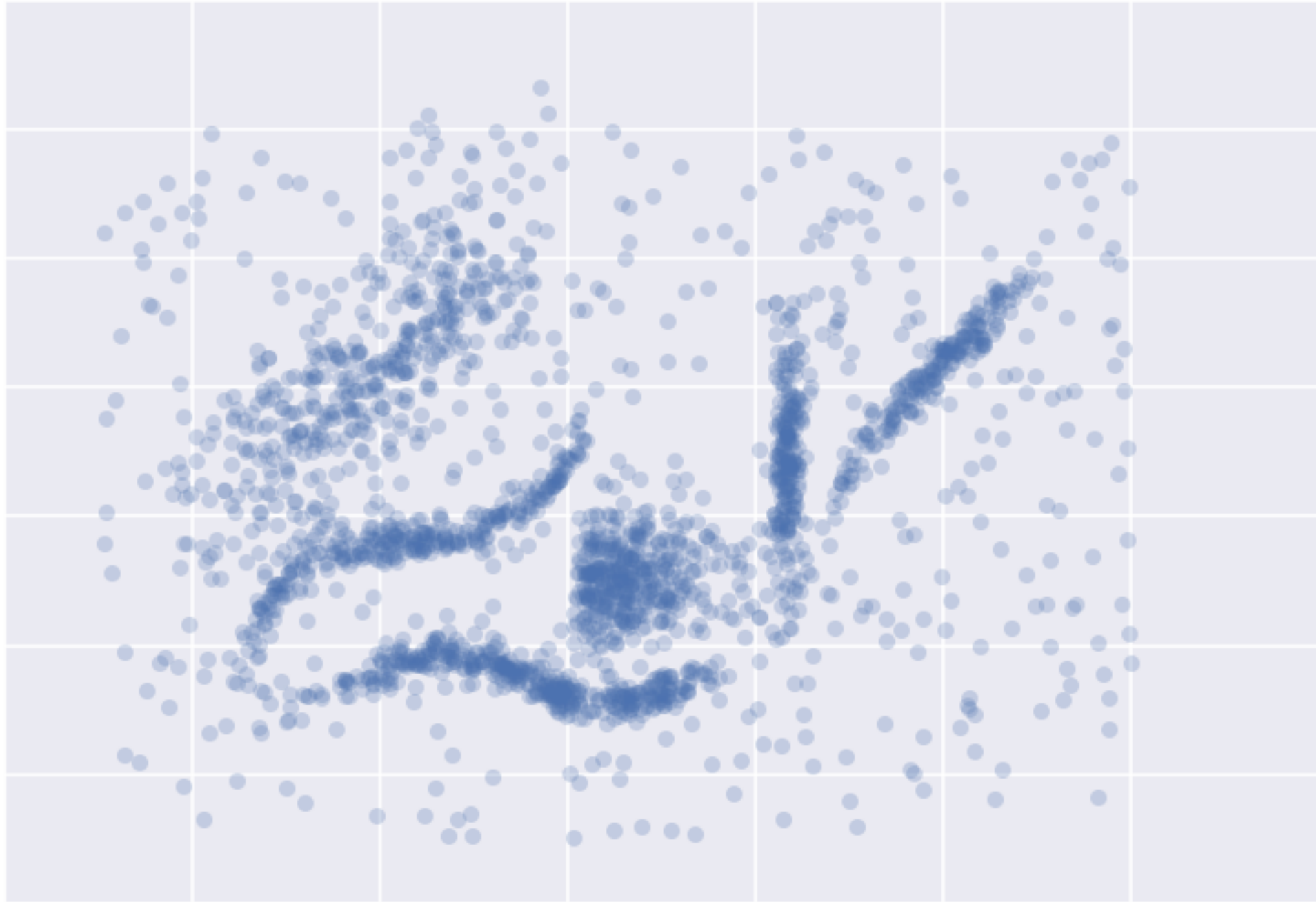
# What that means depends on your application

# Goals

- Partition your data
- Summarize your data
- Explore your data
- Embed in a vector space
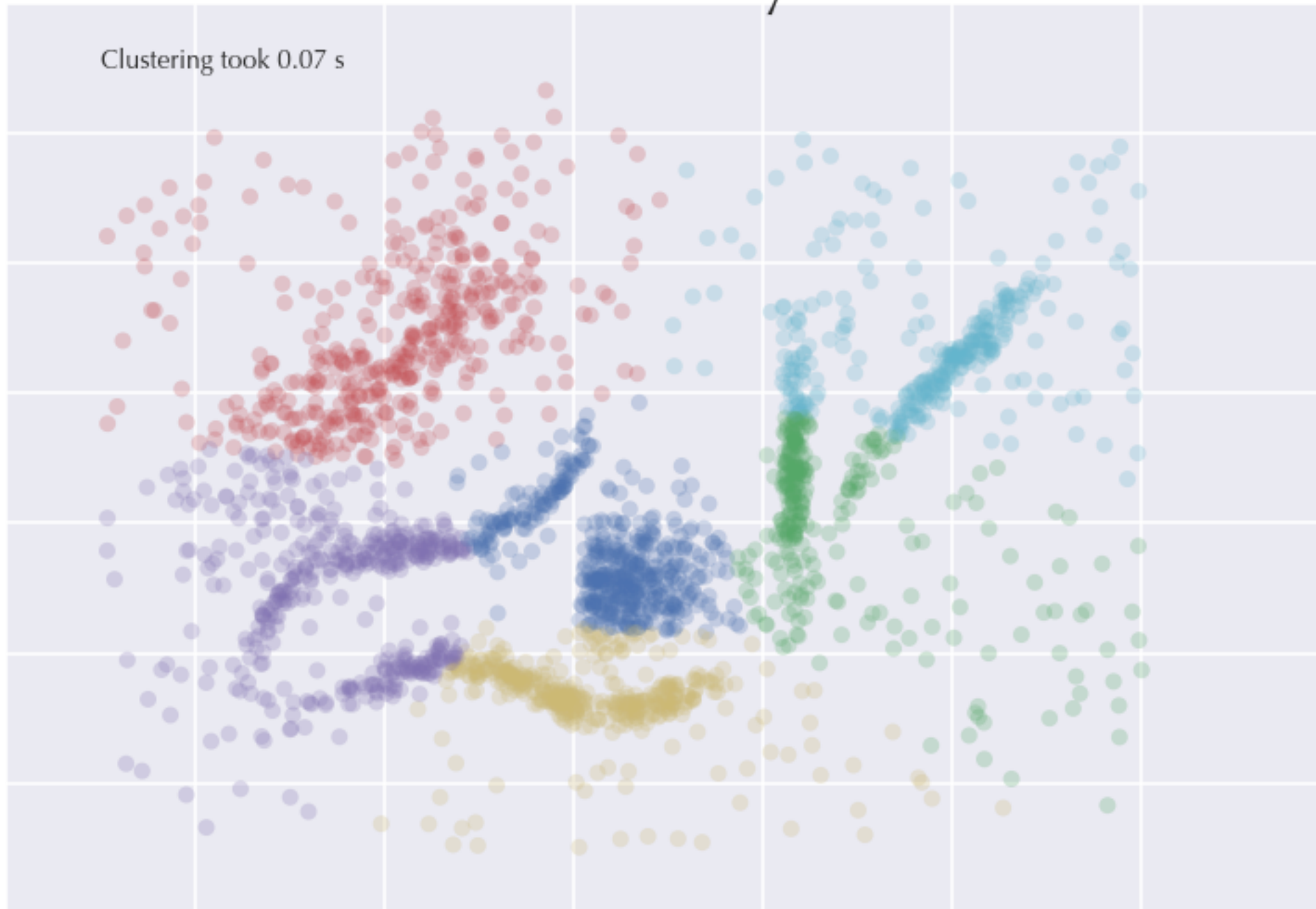- Find patterns in your data
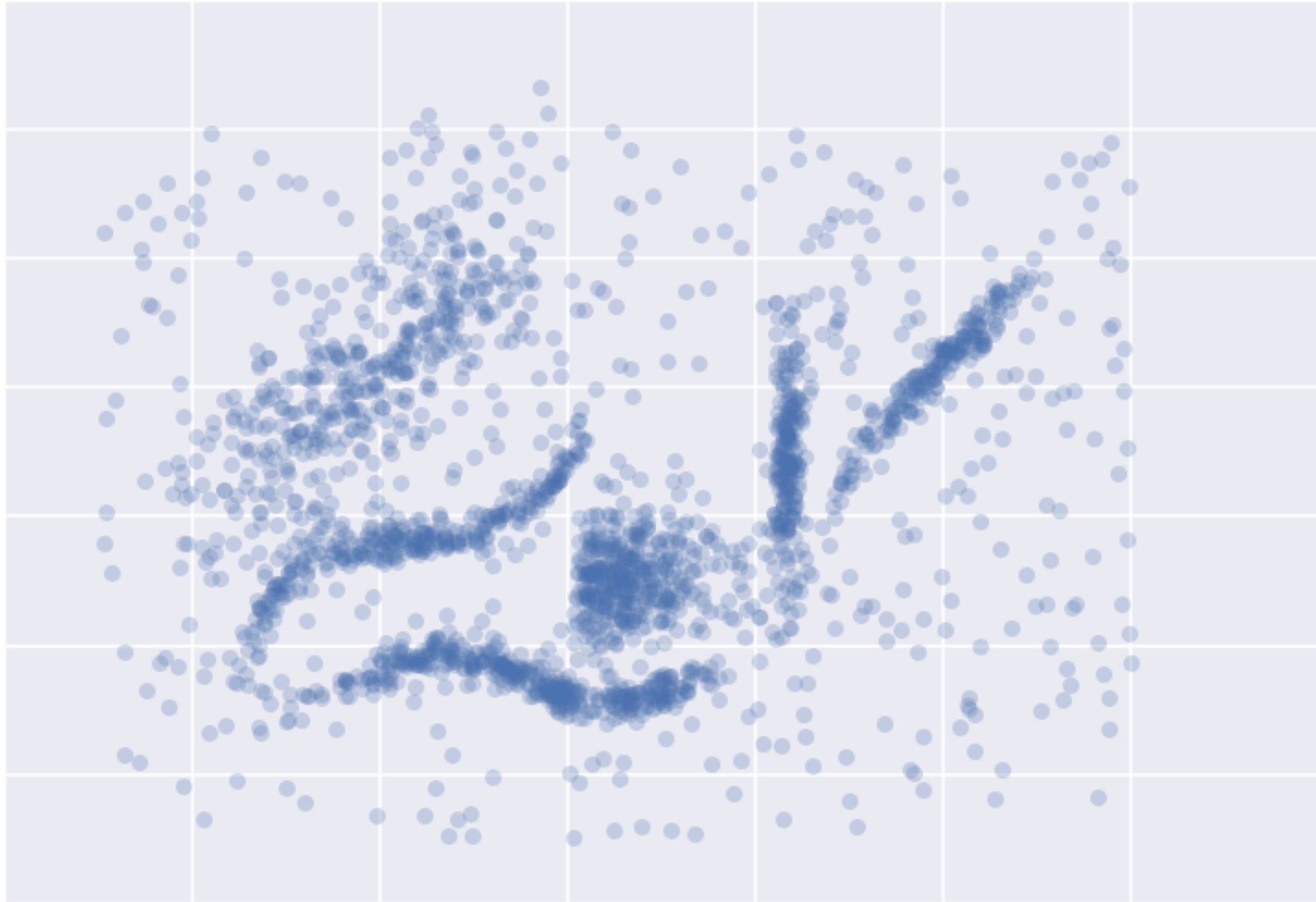- …

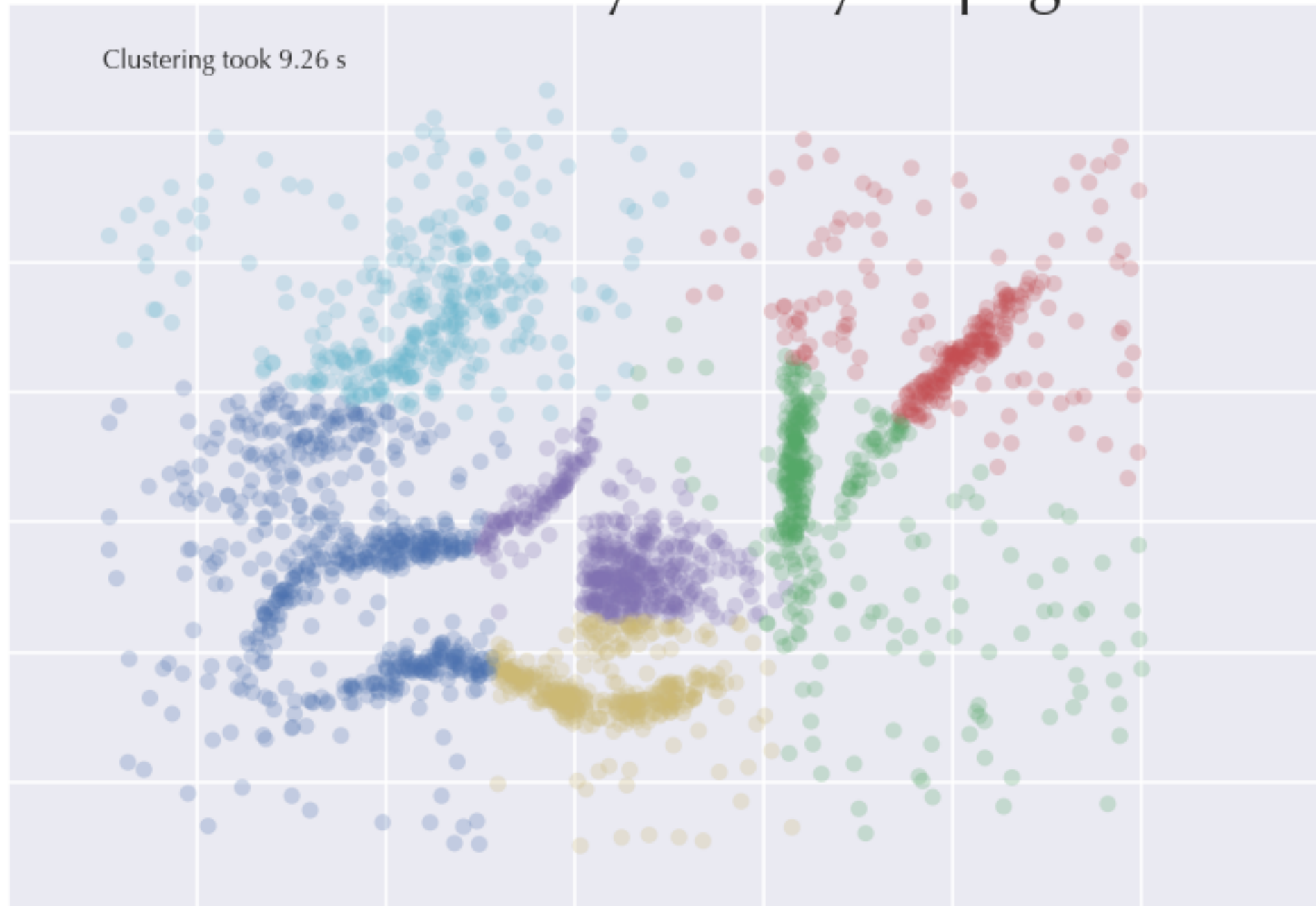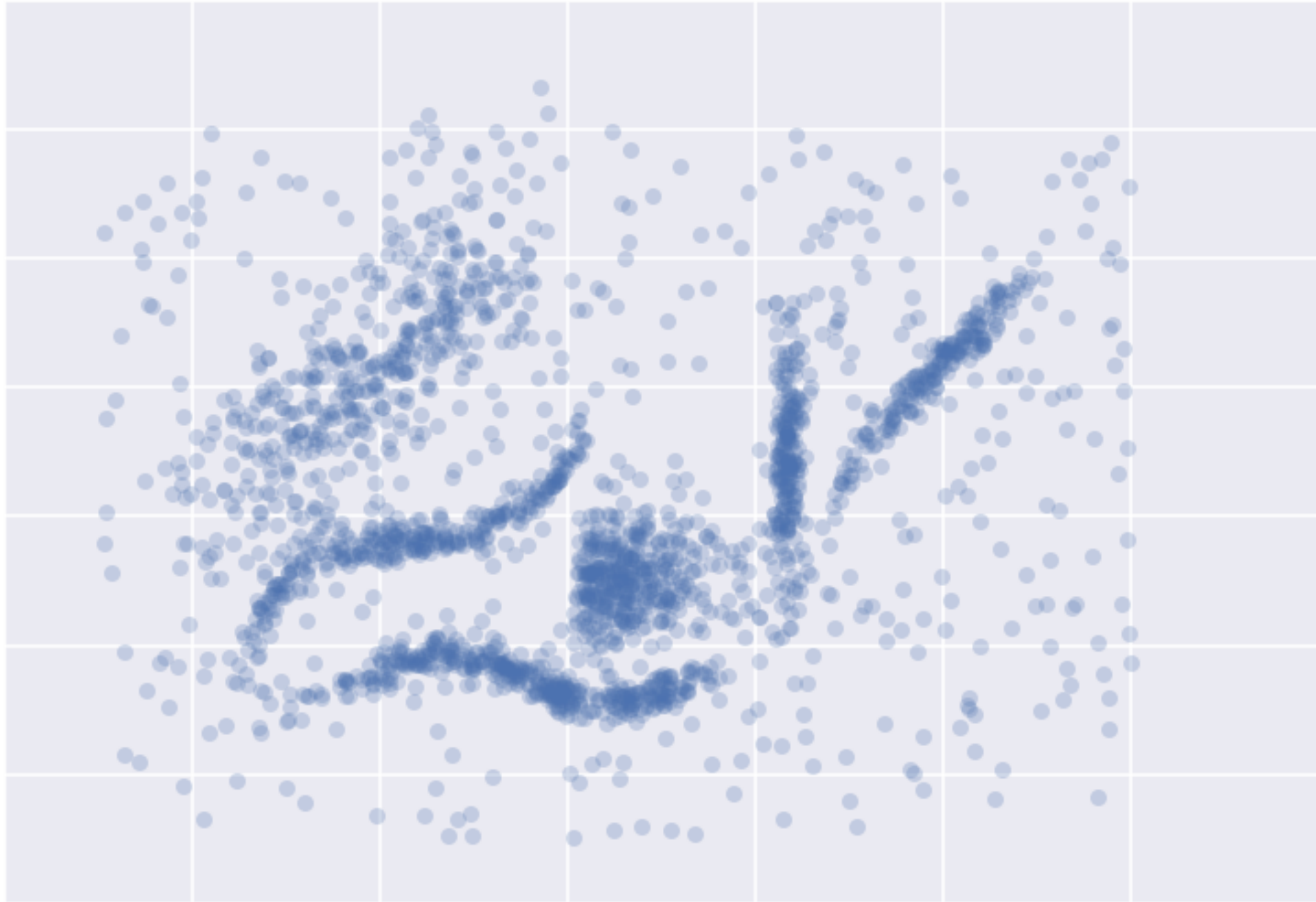Find groups of data that are all similar

Easy in theory
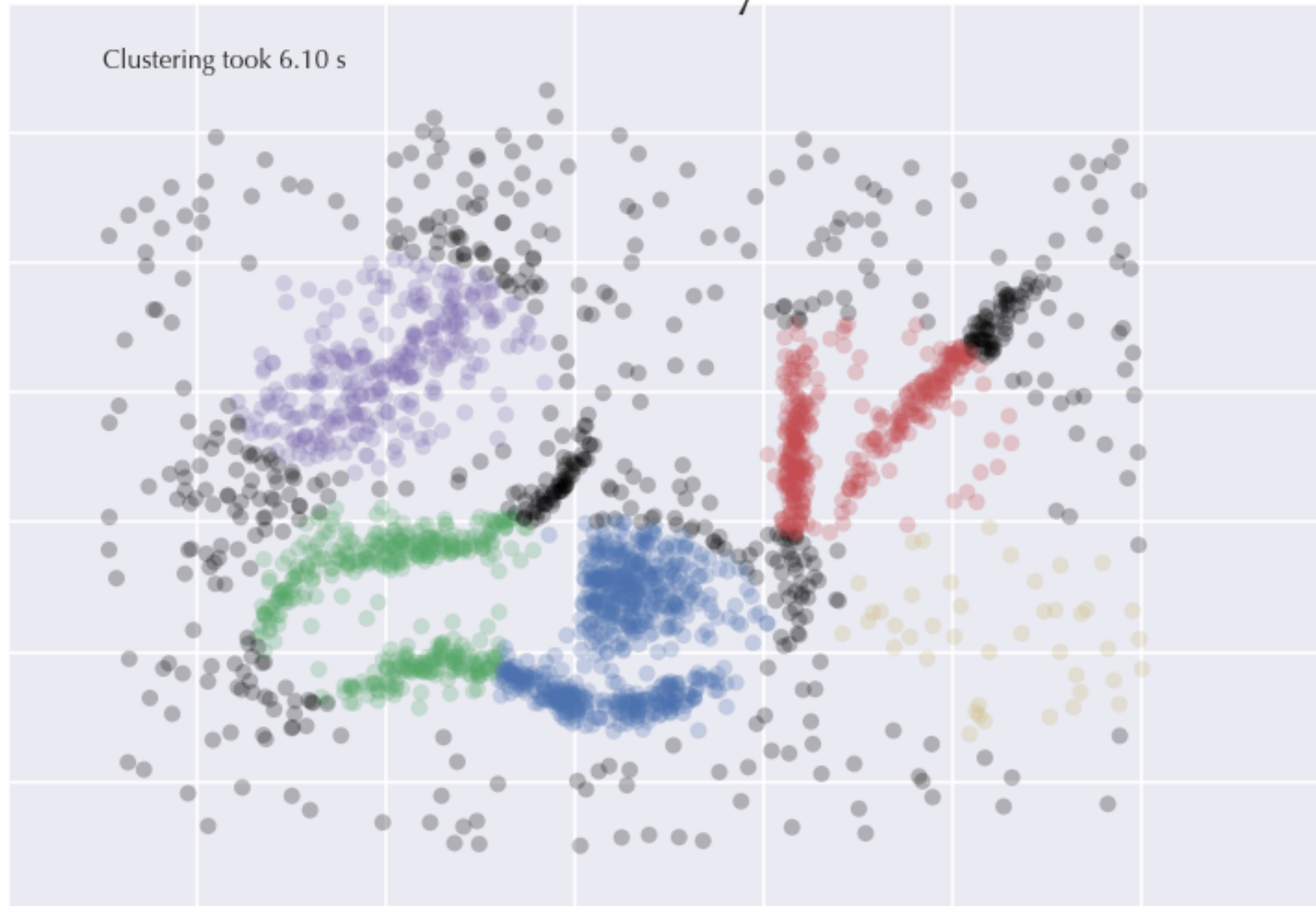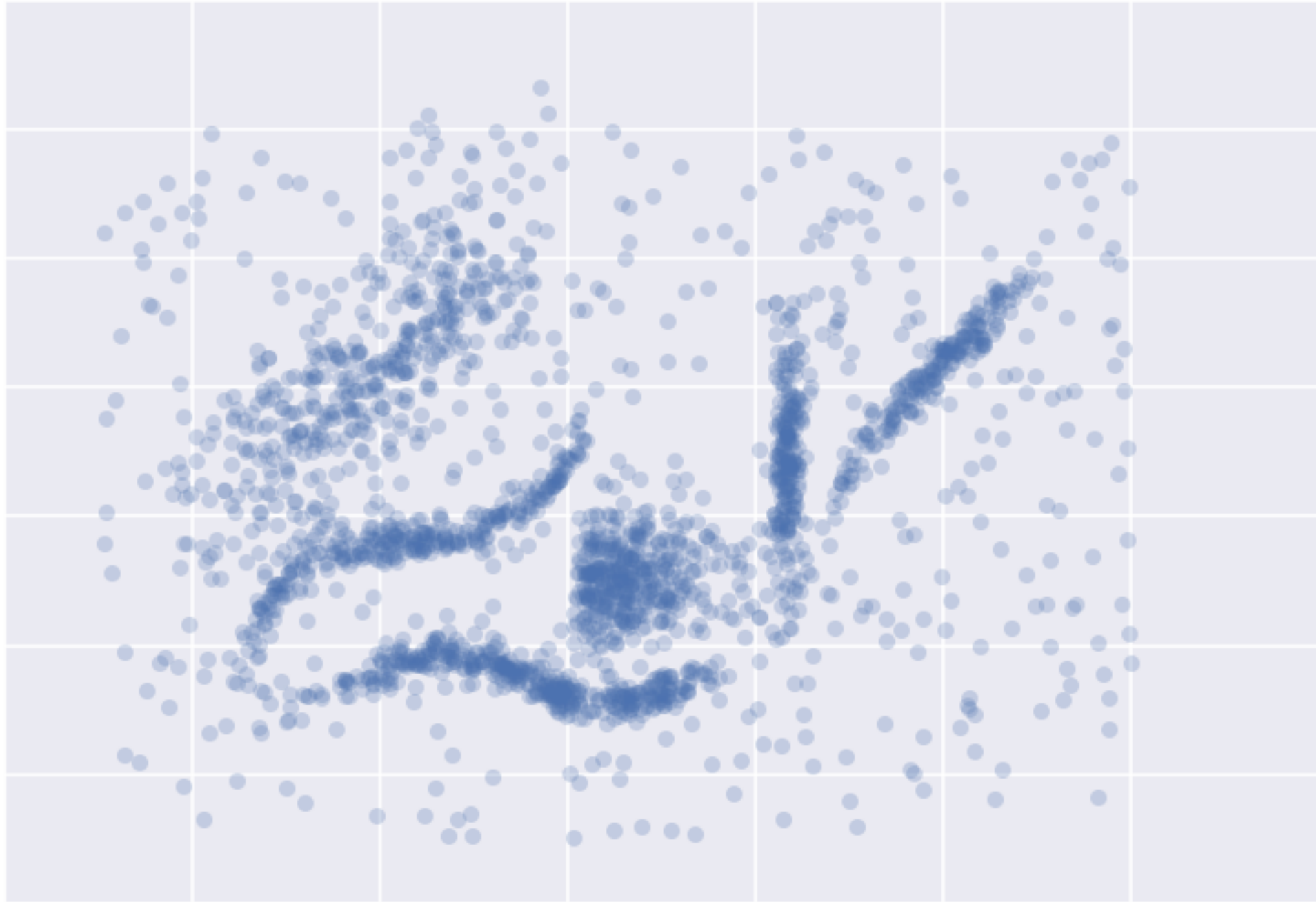
Hard in practice

# Clusters found by KMeans

Clustering took 0.07 s

# Clusters found by AffinityPropagation



Clustering took 9.26 s

# Clusters found by MeanShift

Clustering took 6.10 s

# Clusters found by SpectralClustering

Clustering took 1.09 s

# Clusters found by Birch

Clustering took 0.04 s

# Clusters found by AgglomerativeClustering

Clustering took 4.14 s

# Clusters found by DBSCAN

Clustering took 0.01 s

# What makes it so hard?

# Which clustering is better?

A Possible Clustering

A Possible Clustering

What do I mean by a cluster?

# Let's start by looking at what not to do

# Clusters found by KMeans

Clustering took 0.07 s

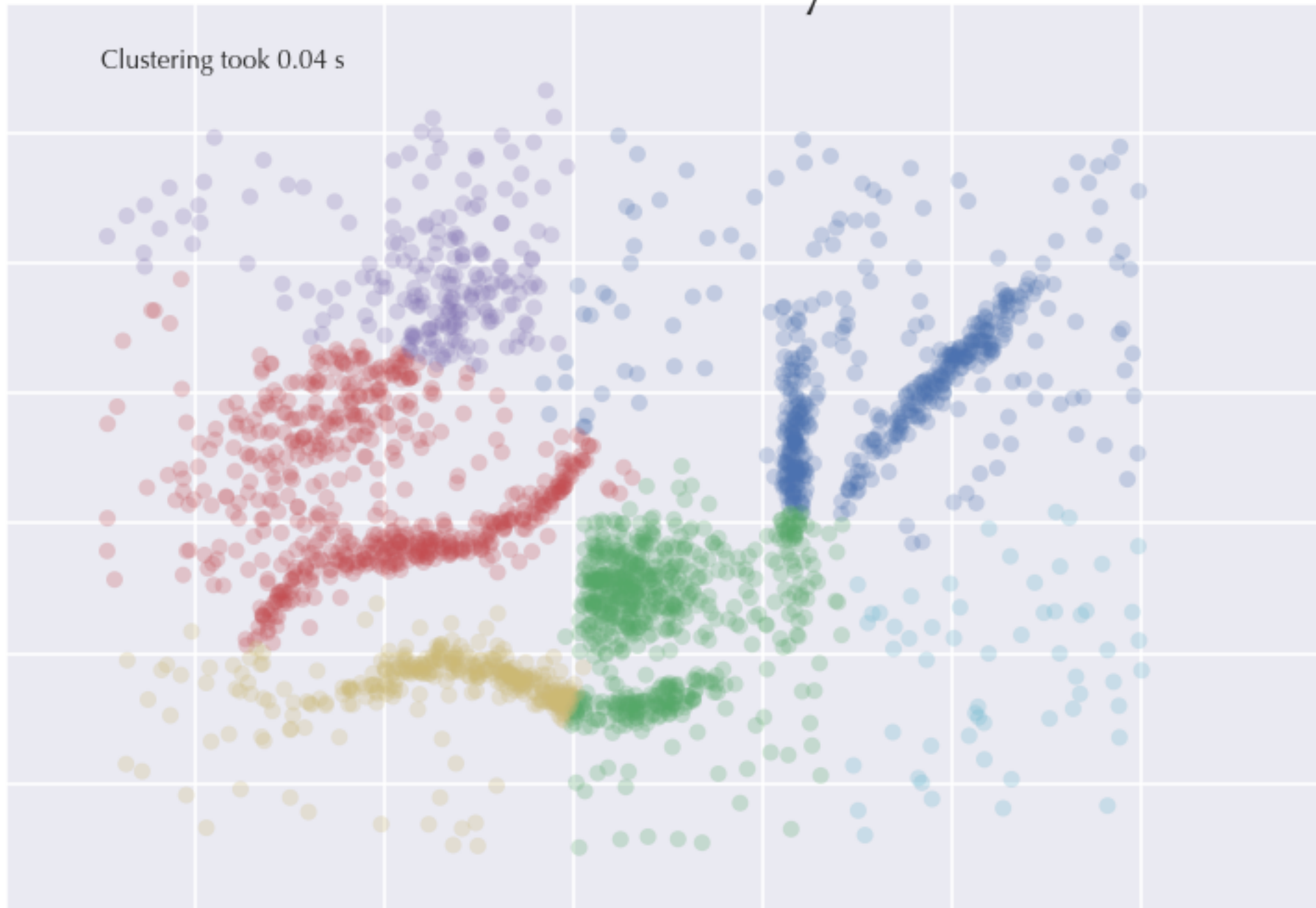Clusters need not be balls!

# Clusters found by SpectralClustering

Clustering took 1.09 s

# Not every point is in a cluster

# Real data has noise!

# Clusters found by DBSCAN



Clustering took 0.01 s

# Clusters are dense areas?

# Separated by less dense areas?

That's what I mean when I talk about clustering

You should think hard about what you mean before doing any clustering

Can we be more specific about density-based clusters?

A connected component of a level set of the probability density function of the underlying (and unknown) distribution from which our data samples are drawn.

# What?

Is that process a clustering algorithm?

1. We don't know the PDF
2. Which level sets to choose?
3. Computational complexity

What can we do?

# Locally approximate the density

The connected components of level sets form a tree

Approximate the level set tree, and use 'Excess of Mass' to select clusters

What does this look like in practice?

Density Cluster Tree

Simplified Cluster Tree

Simplified Cluster Tree

And the resulting clusters…

# Clusters found by HDBSCAN



Clustering took 0.04 s

# This is the HDBSCAN* clustering algorithm

# But what about performance?

We don't want to run connected components for every possible epsilon!

# Minimum spanning trees!

The weighted graph is complete!

# We can use spatial indexing to compute fewer distances

# Spatial indexing is great for neighbour queries

If we use Boruvka's algorithm we can reduce minimum spanning tree computations to repeated neighbour queries

# A modified Dual Tree Boruvka algorithm provides O(n log n) performance

March, Ram, Gray 2010

Performance Comparison of Clustering Implementations

Legend:
- Sklearn K-Means
- Sklearn DBSCAN
- Scipy K-Means
- HDBSCAN
- Fastcluster Single Linkage
- Scipy Single Linkage
- DeBaCl Geom Tree
- Sklearn Spectral
- Sklearn Agglomerative
- Sklearn Affinity Propagation

Y-axis: Time taken to cluster (s)
X-axis: Number of data points

Performance Comparison of Fastest Clustering Implementations

Performance Comparison of K-Means and DBSCAN

| | Number of data points |
|---|---|
| Interactive | 100,000 |
| Over coffee | 500,000 |
| Over lunch | 1,000,000 |
| Over night | 5,000,000 |

We still have to choose a "number of points" value for the density estimate

Topology to the rescue!

INSTITUT TUTTE INSTITUTE

Building the initial HDBSCAN* tree can be described in terms of persistent homology

Multi-dimensional persistent homology allows for persistence over multiple variables

Unfortunately the result is not a tree

We can re-interpret HDBSCAN* using sheaves instead of trees

The resulting algorithm can be generalized to the multi-dimensional case!

More robust
Fewer parameters
Similar performance

# K-Means shouldn't be your first choice

# K-Means probably shouldn't be your second choice either

# You may as well run HDBSCAN* while you're thinking!

https://github.com/scikit-learn-contrib/hdbscan

conda install -c conda-forge hdbscan

pip install hdbscan