

# **GEOG 432/832: Programming, Scripting, and Automation for GIS**

## **Unit 13.01: Clustering**

**Dr. Bitterman**

# Today's schedule

- Open discussion
- Announcement/update
- Slides, discussion, and exercises
- For next class

# Open discussion

# Announcements

- Presentation order is posted to Canvas discussion board
  - We will follow order
  - But exact day may change depending on time
  - PRACTICE, PRACTICE, PRACTICE
  - I WILL CUT YOU OFF AT THE TIME RESTRICTION (UG 10 min, G 13 min)
- Wednesday: we will meet, but mostly work time

# Today's prep

Download *unit13inclass.ipynb* and *unit13data.zip* from Canvas

# Today's prep:

```
%matplotlib inline

import seaborn as sns
import pandas as pd
import geopandas as gpd
import numpy as np
import matplotlib.pyplot as plt
# new ones below
from sklearn import cluster # note the difference
import contextily as cx
```

**Verify all packages are in your environment**

**...and that you're in the correct environment**

**...can you find them all???**

# Alternative sources

Not all packages are in the default channel (like *contextily*)

## Let's add one

1. While selecting the correct environment, click on "Channels", then "Add..."
2. in the new box, type the URL of conda-forge: `https://conda.anaconda.org/conda-forge/`
3. Press "Enter" to add
4. Update index (or now, indices)
5. Add *contextily* to your environment

# Today we'll be tackling clustering

What's a cluster???



# Background

## Real world is complex

- non-linear
- uncertain
- multivariate (as opposed to univariate...)

# Tackling multivariate data

- Can be difficult to "look at" more than a few variables simultaneously
- Often necessary to *reduce dimensionality*
- Many techniques don't require preliminary assumptions about data structure
- Also useful as an exploratory tool

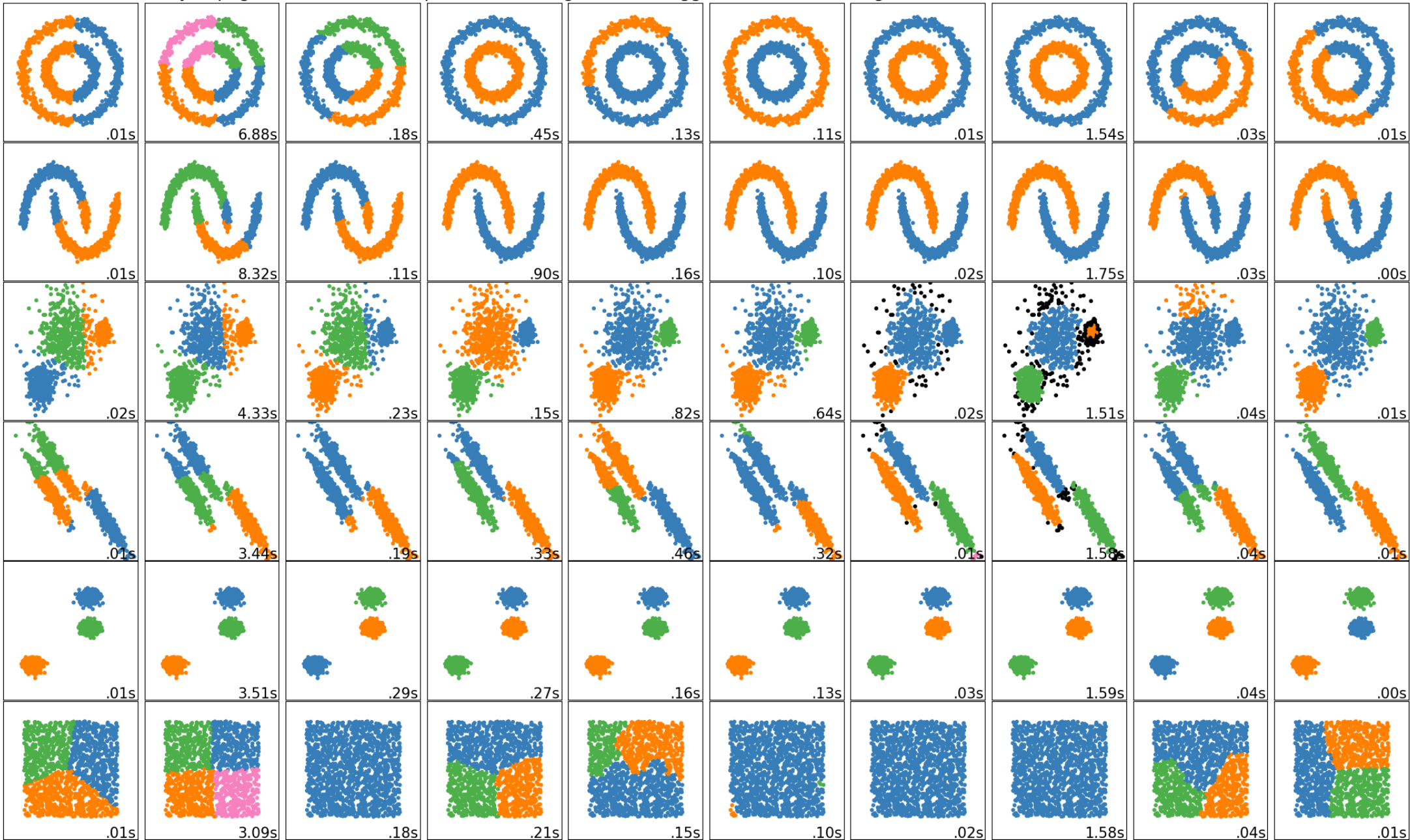
# Statistical clustering

- **Basic idea:** summarize the information multivariate data by creating a (relatively) small number of classes
- Each observation is assigned to one (and only one) class, depending on its position in multidimensional space
- MANY, MANY, MANY techniques to statistically group observations
- Common thread: define classes or categories of observations that are similar within each of them, but differ between groups

## But implementation details differ

- How "similarity" and "dissimilarity" are defined differs among algorithms
- Applications and utility depend on problem set

MiniBatchKMeans AffinityPropagation MeanShift SpectralClustering Ward AgglomerativeClustering DBSCAN OPTICS Birch GaussianMixture



# K-means

- A (very) common statistical clustering technique
- anyone familiar???

## Whiteboard demo

(also: <https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html>)

**What are the big assumptions?**

**(and/or required parameters)**

# Interactive session (Jupyter notebook)



## For next class

- We will meet on Wednesday, but expect work time
- No class meeting on Friday (work on your own, but let me know if questions)
- Lab 8 due Thursday
- Readings are linked/posted on Canvas...
- Completed project signup sheet on Canvas (discussion board)