# Bradley Voytek, Ph.D.
UC San Diego
Cognitive and Neural Dynamics Laboratory

Department of Cognitive Science
Neurosciences Graduate Program
Halıcıoğlu Data Science Institute

bvoytek@ucsd.edu
@bradleyvoytek

UC San Diego

# COGS 108
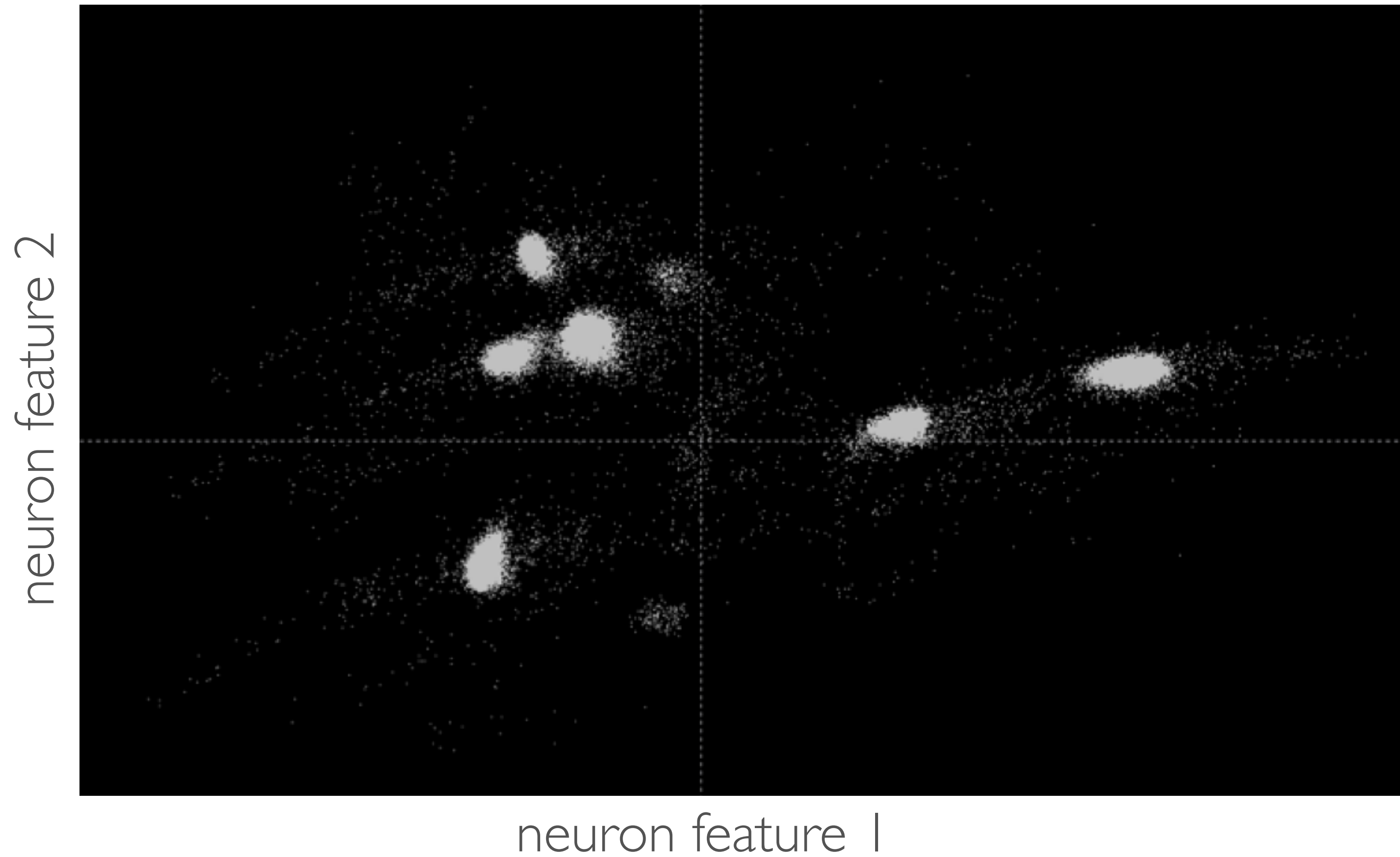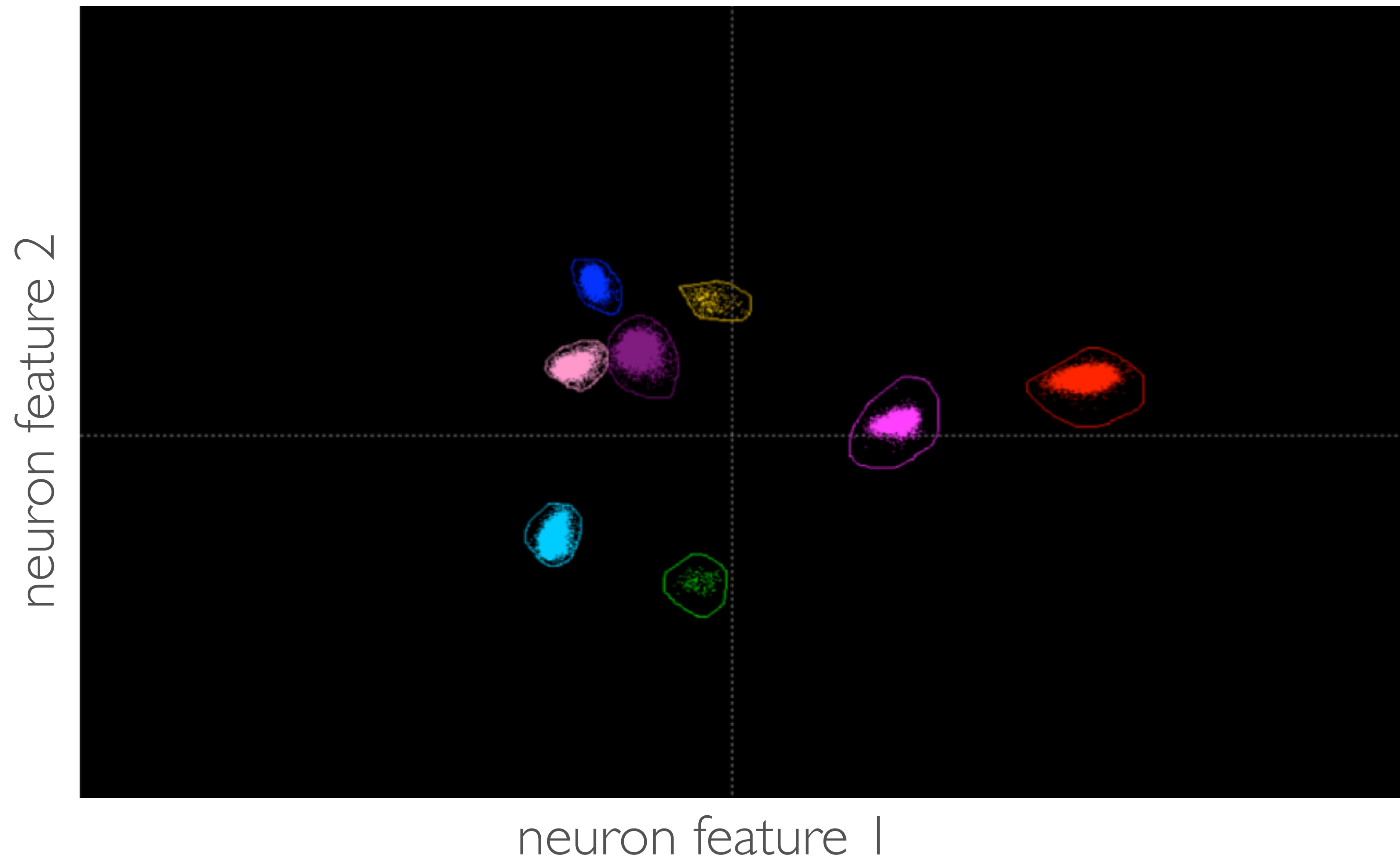## Data Science in Practice

*Dimensionality reduction and clustering*

# Clustering



neuron feature 1

neuron feature 2

UC San Diego

# Clustering



neuron feature 2

neuron feature 1

UC San Diego

# Why clustering?

- Unlabeled data – unsupervised learning
  - If it is labeled, we'll use classification!

- Want to find groups!

UC San Diego

# Why clustering?

- Data Compression – We could want to reduce the dimensionality in the data.
    - Dimensionality – An aspect or a feature of something is a dimension.

UC San Diego

# How do we cluster?

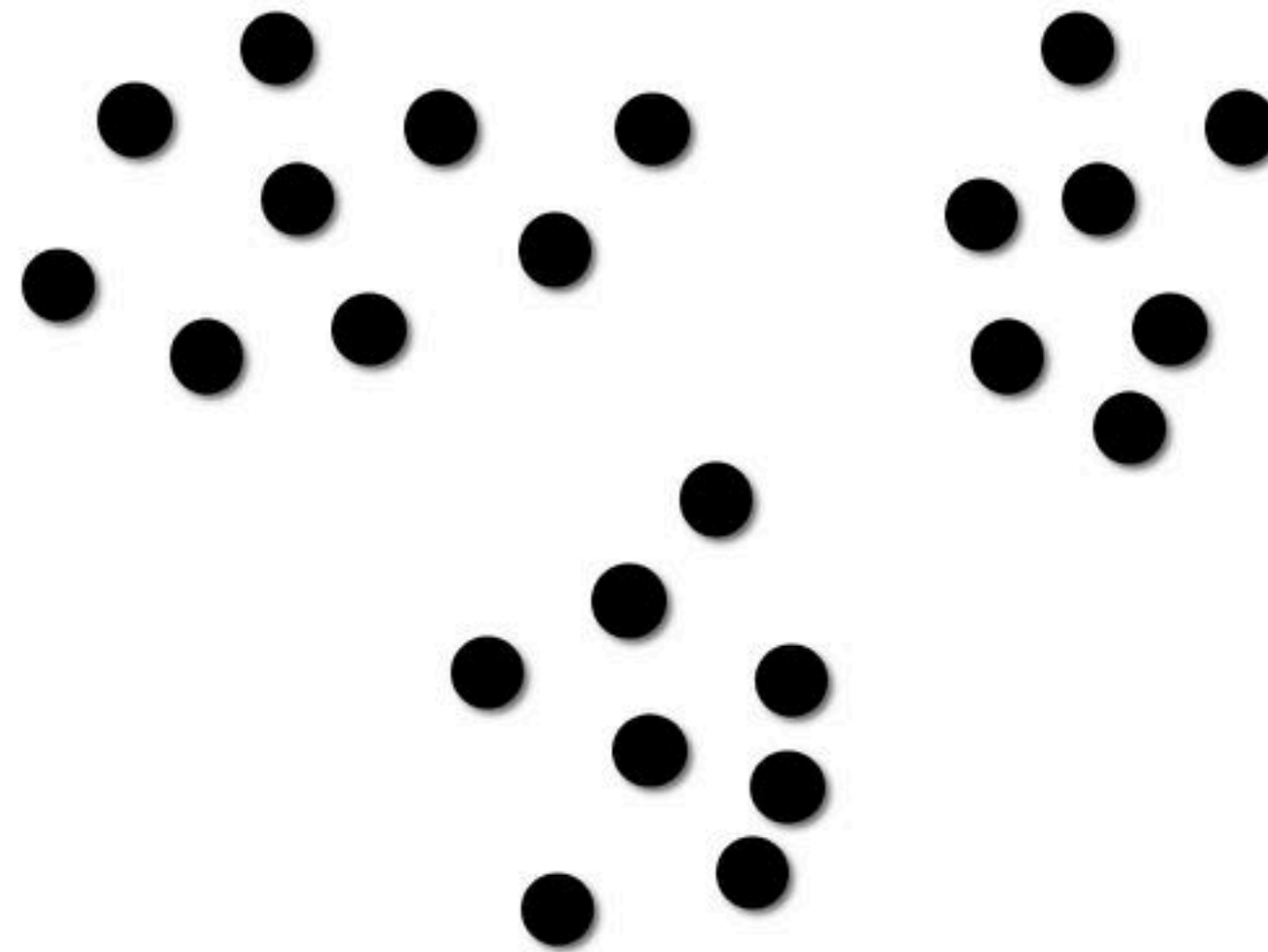- A lot of different ways we could, but we will start with K-Means.

UC San Diego

# k-means

- A simple but effective clustering algorithm

- Partitions the data in to K disjoint sets (clusters)
  - disjoint - no overlap

- We then run an iterative batch algorithm
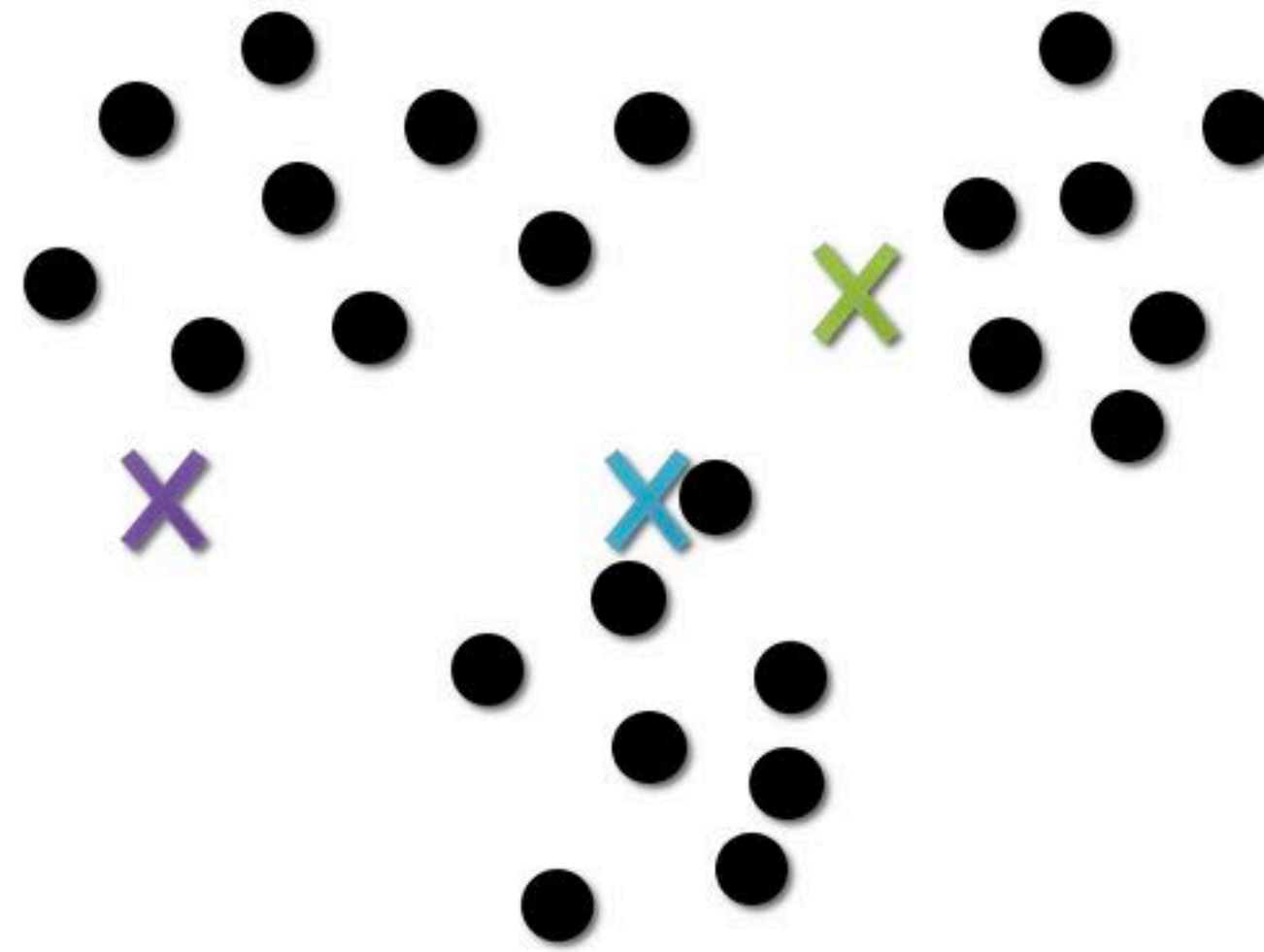  - batch – to do to all of the items (batch of cookies)

UC San Diego

# k-means

- Start – Initialize a guess of k centers/"means"
  - notation – $j$ simply refers to which of the $k$ centers.

- $S^{(j)}$ is all points closest to $\boldsymbol{\mu}^{(j)}$
  - So each "mean" gets all of the points close to it in its set.

- Update the means:    $$\mu^{(j)} = \frac{1}{N_j} \sum_{n \in S^{(j)}} \mathbf{x}^{(n)}$$

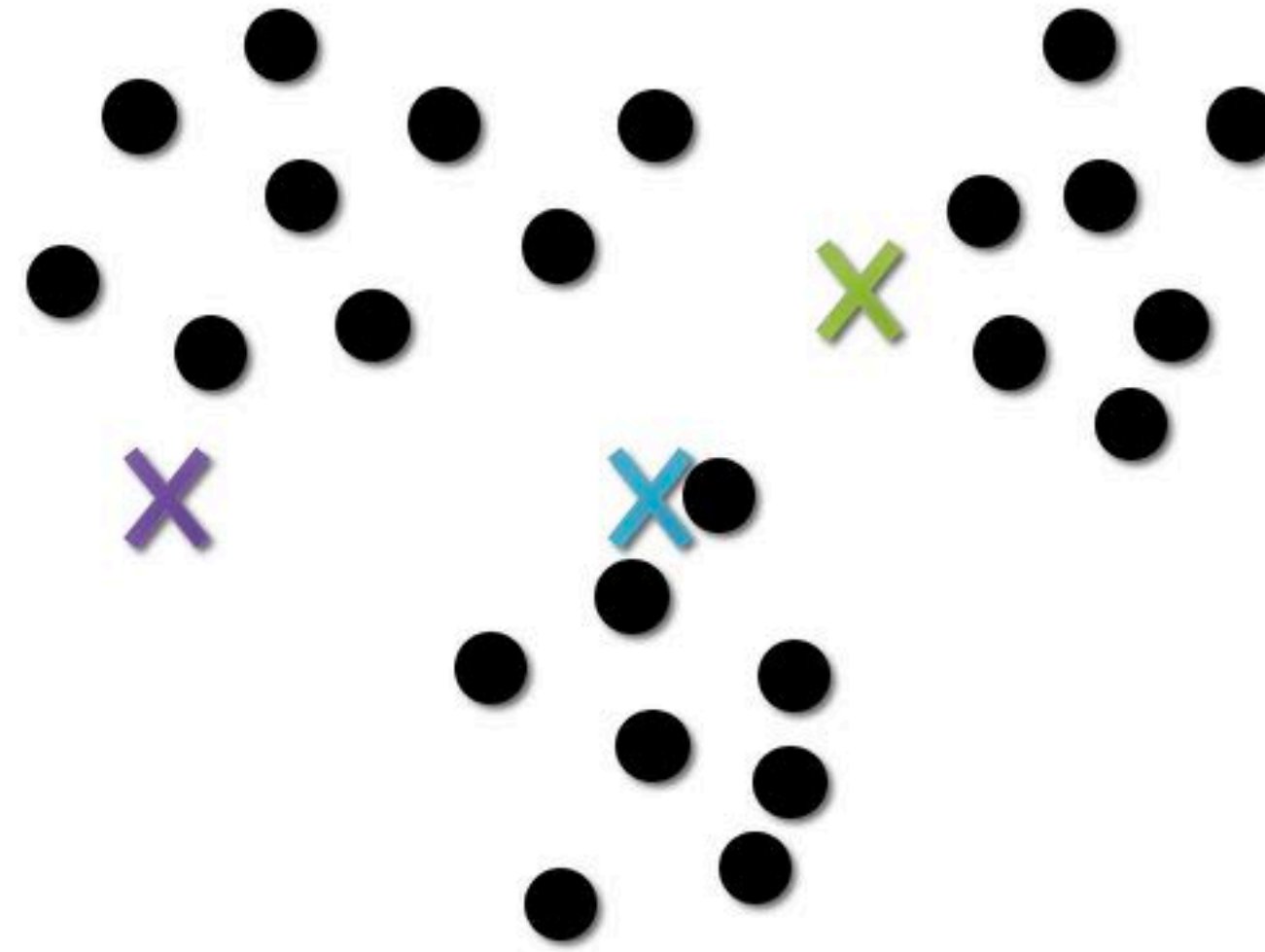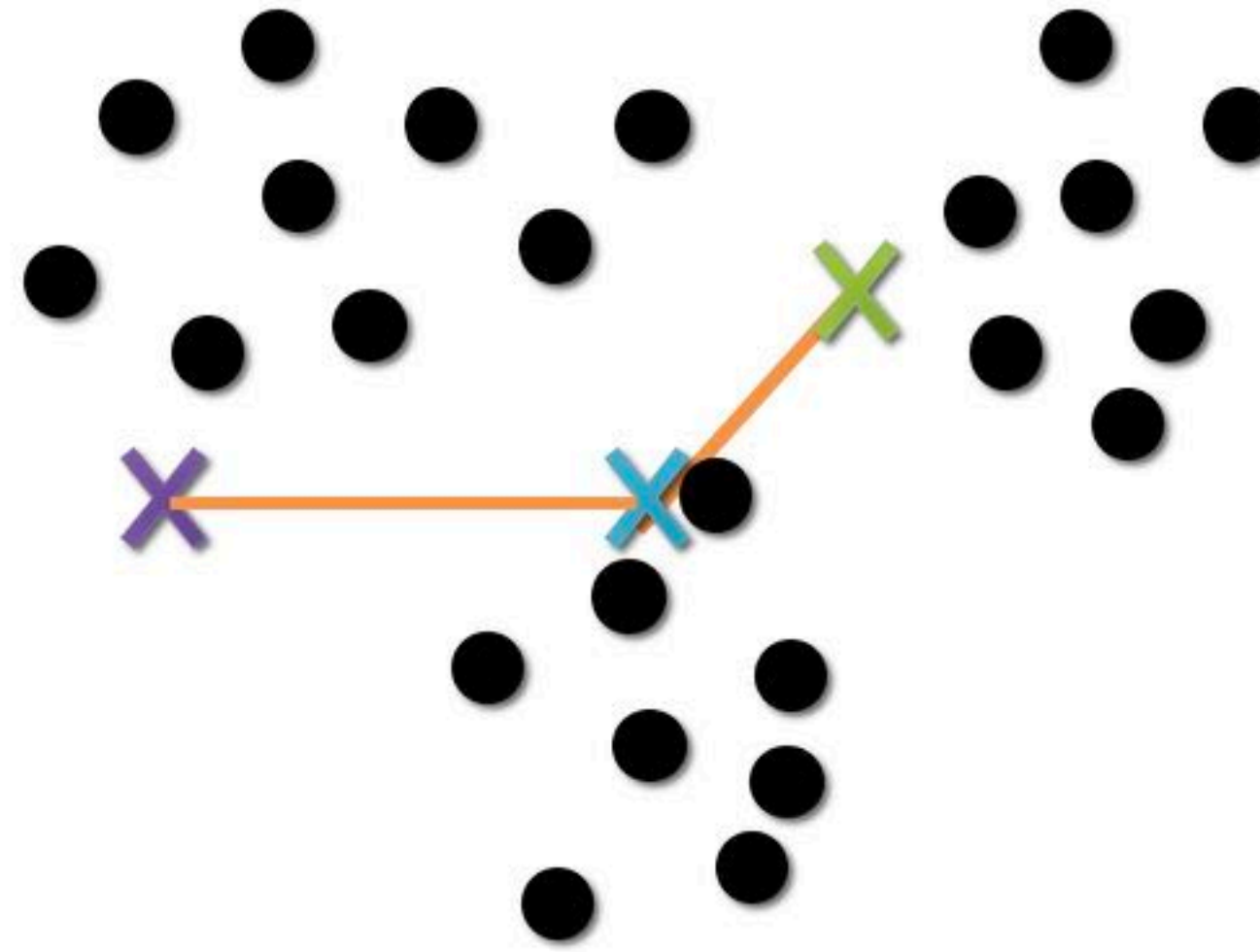- Continue until there is no change in the means or you decide to stop.
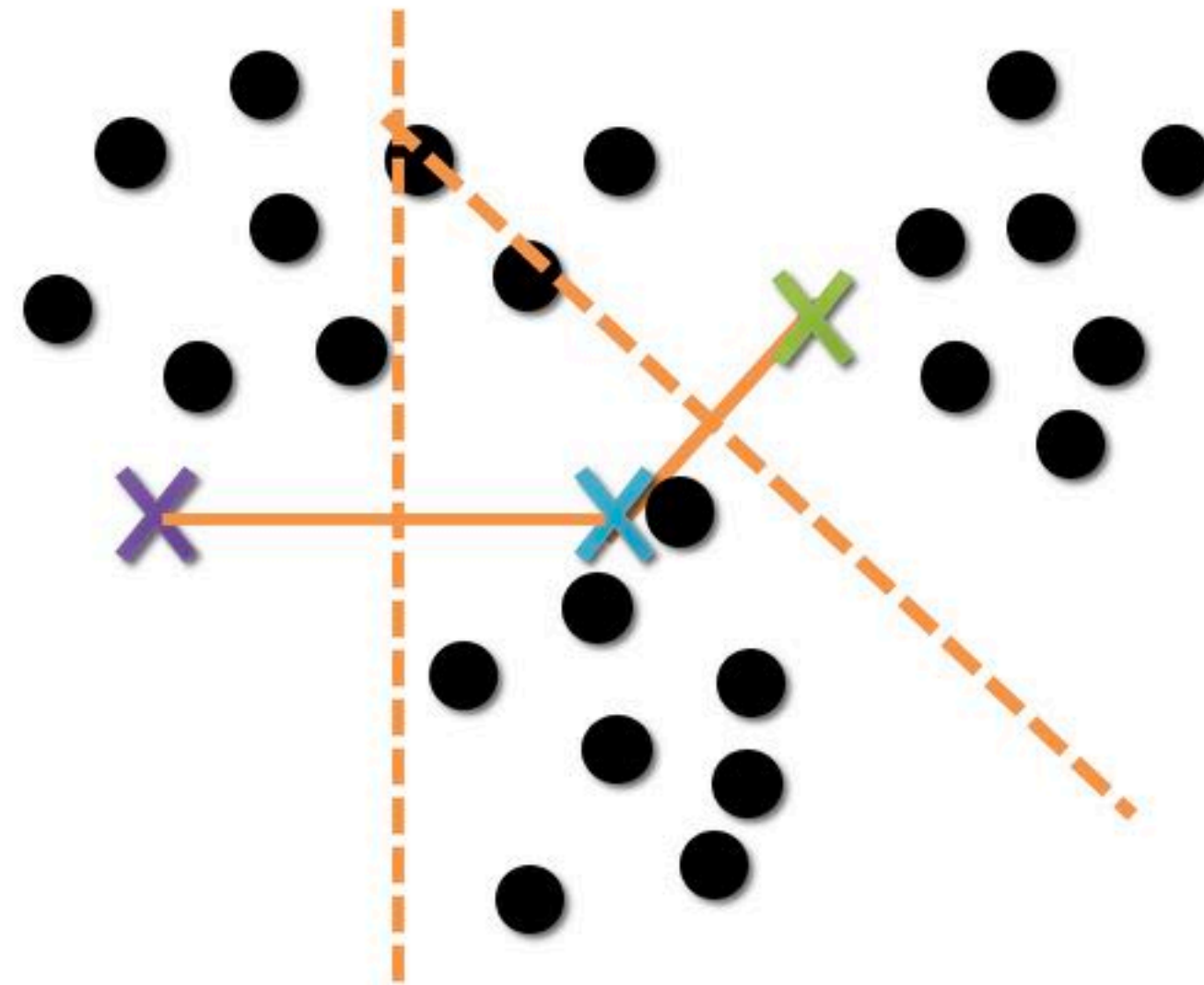
UC San Diego

# Initialize means

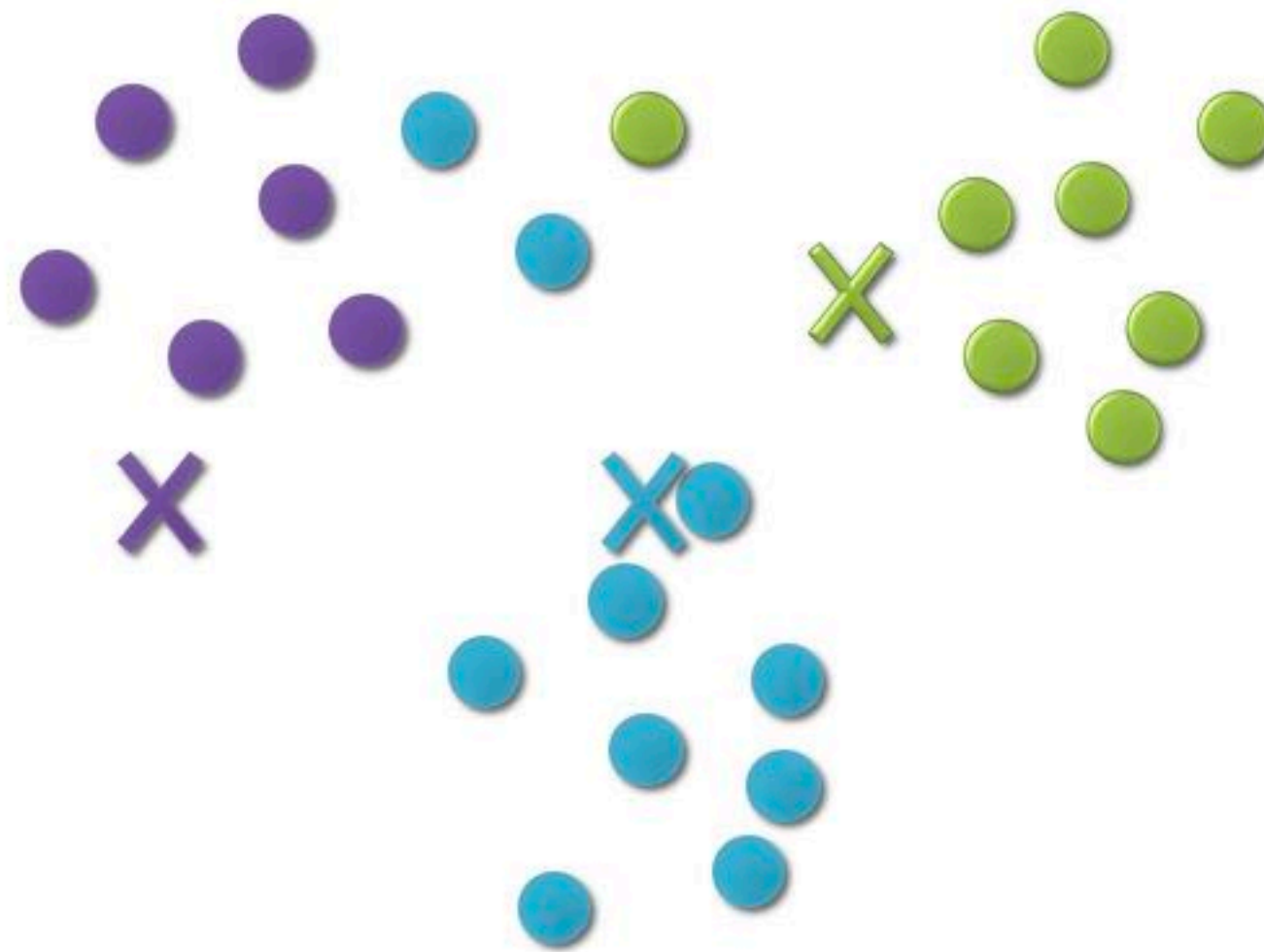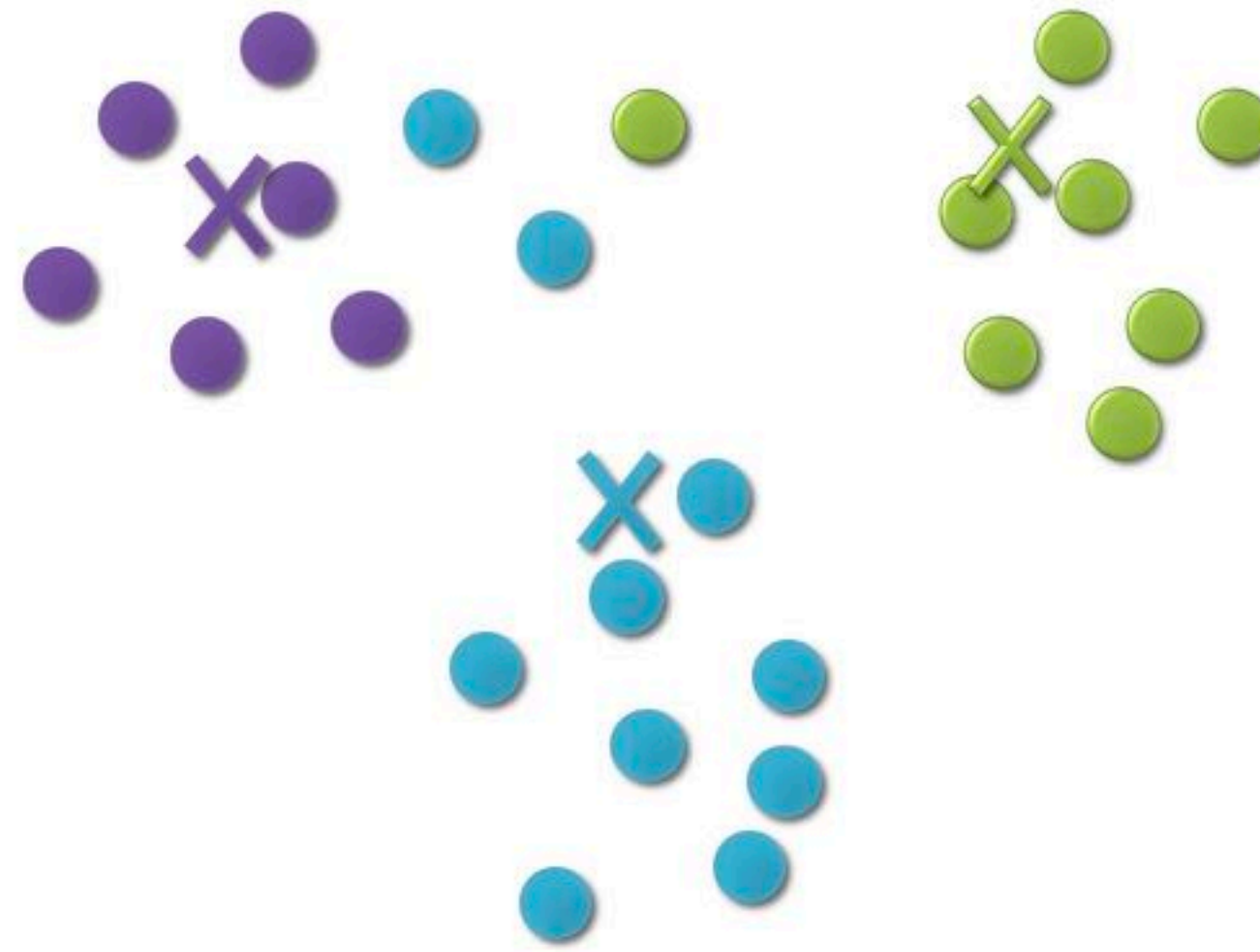# Initialize means

# Label points

# Connect adjacent means

UC San Diego

# Bisect

UC San Diego

# Label points

# Move means to centroids

UC San Diego

# Move means to centroids

UC San Diego

# If means moved, repeat!

UC San Diego

# Bisect

UC San Diego

# Label points

UC San Diego

# Connect adjacent

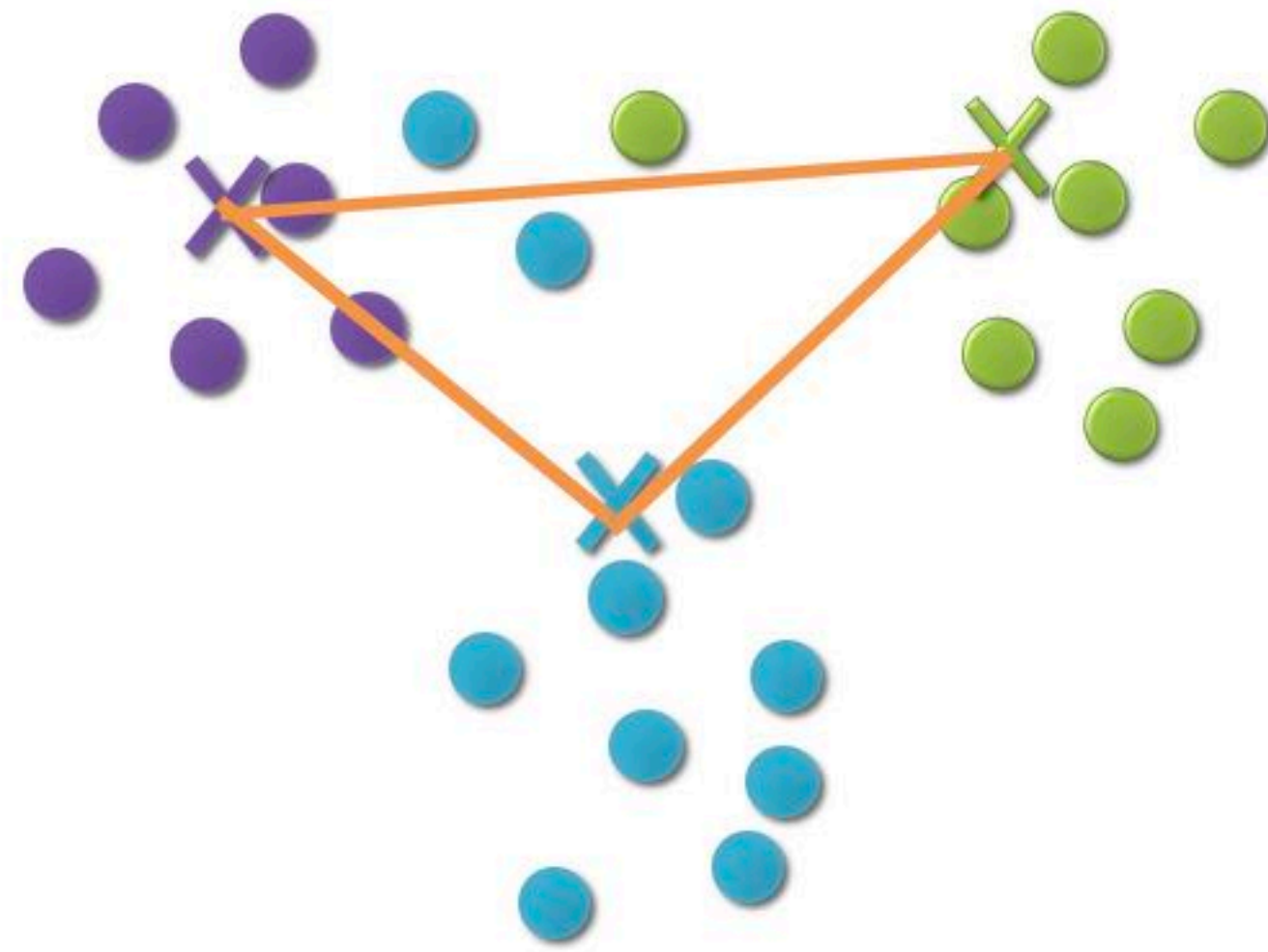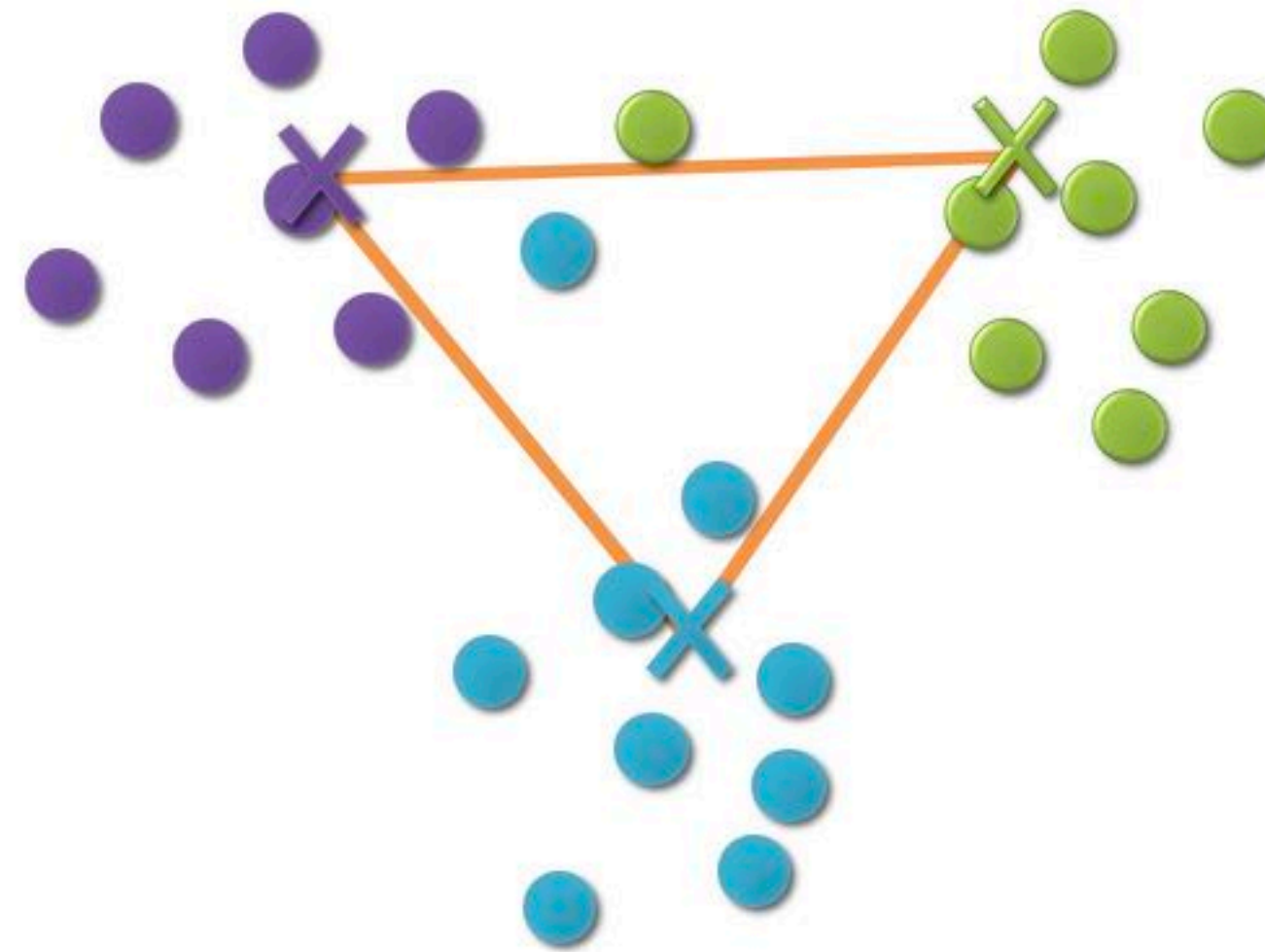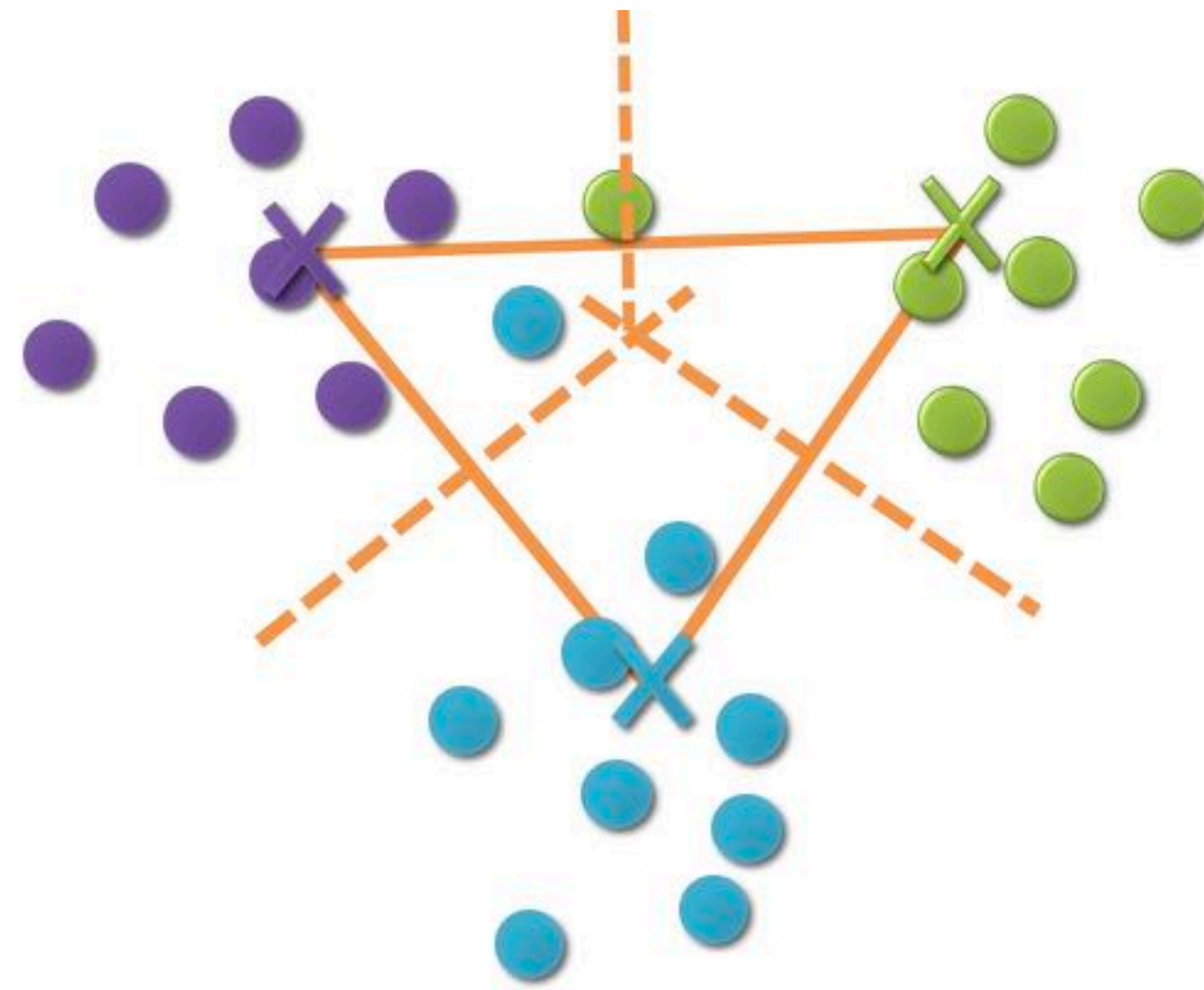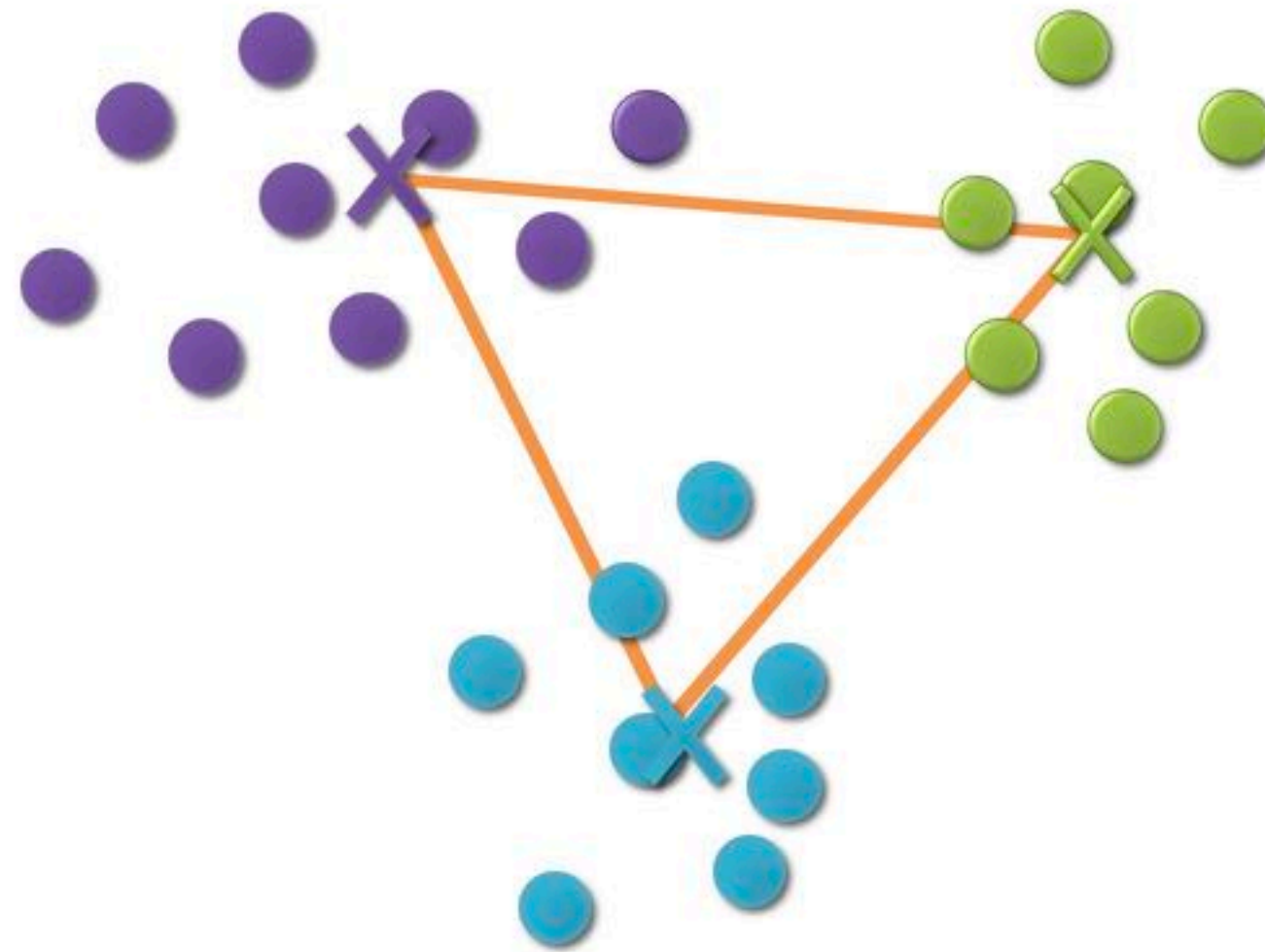UC San Diego

# Bisect

# Label points

# Connect adjacent

# Bisect

# Label points

UC San Diego

# Finished!

UC San Diego
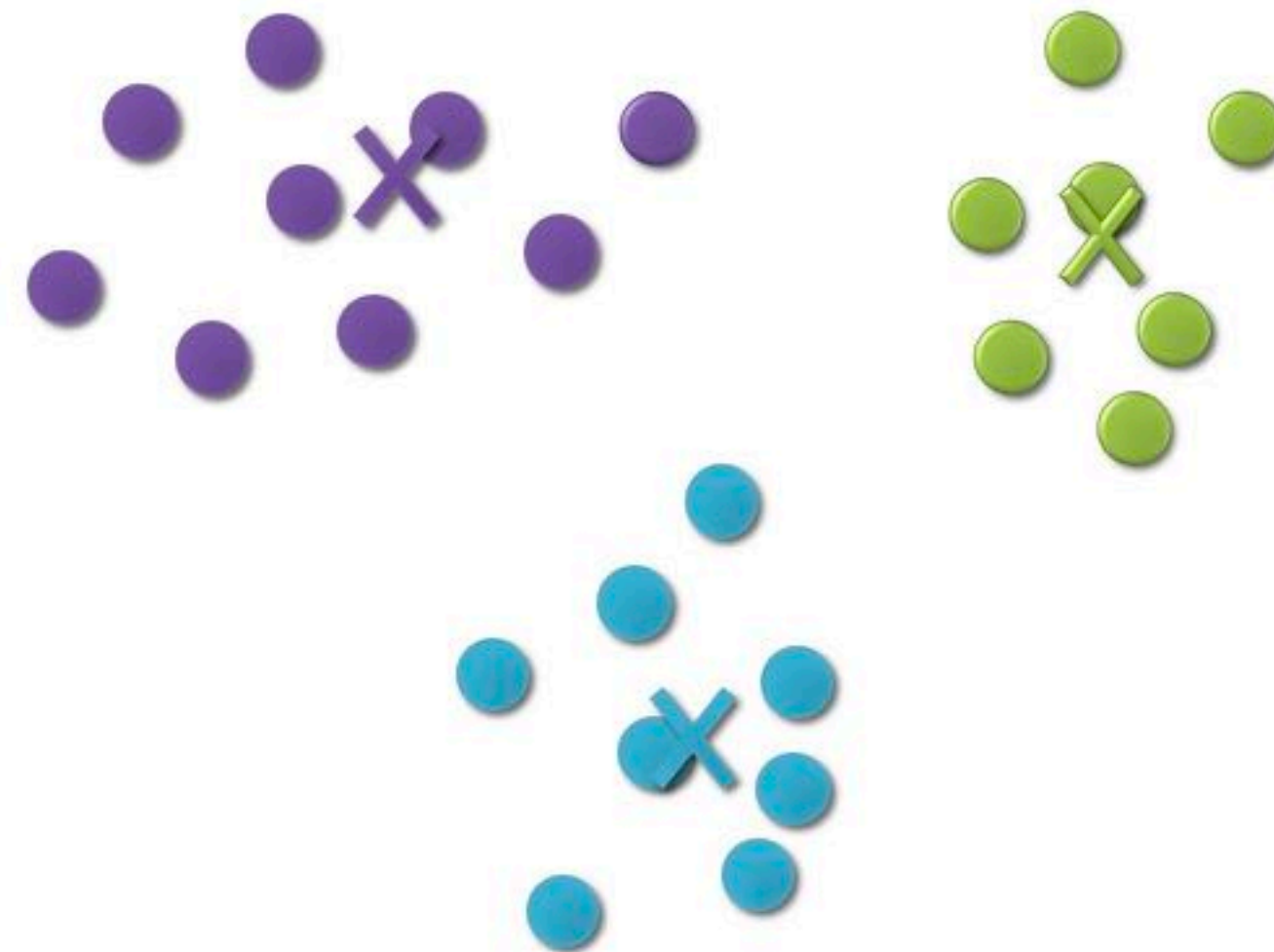
# When do we reduce dimensionality?

Alicia is researching the resistance of plants in different soil to bugs. She has 10 each of 5 different types of plants. She has counted the number of bugs on each once per week for 10 weeks. Should Alicia analyze her bug count data with a clustering analysis? Why?

# When do we reduce dimensionality?

Pat is studying cell types and morhpologies in the basal forebrain of mice. They recorded the waveform amplitude, shapes, and firing patterns of 400 neurons. Should Pat run a clustering analysis?

UC San Diego

# When do we reduce dimensionality?

Linh is working at Apple. They have data on the accounts for all of their users, including what programs they use, how often, and for how long. They've been tasked with mining the data to find patterns that could affect marketing strategies. Should Linh run a clustering analysis?

Source: Jake Olson

# COGS 108
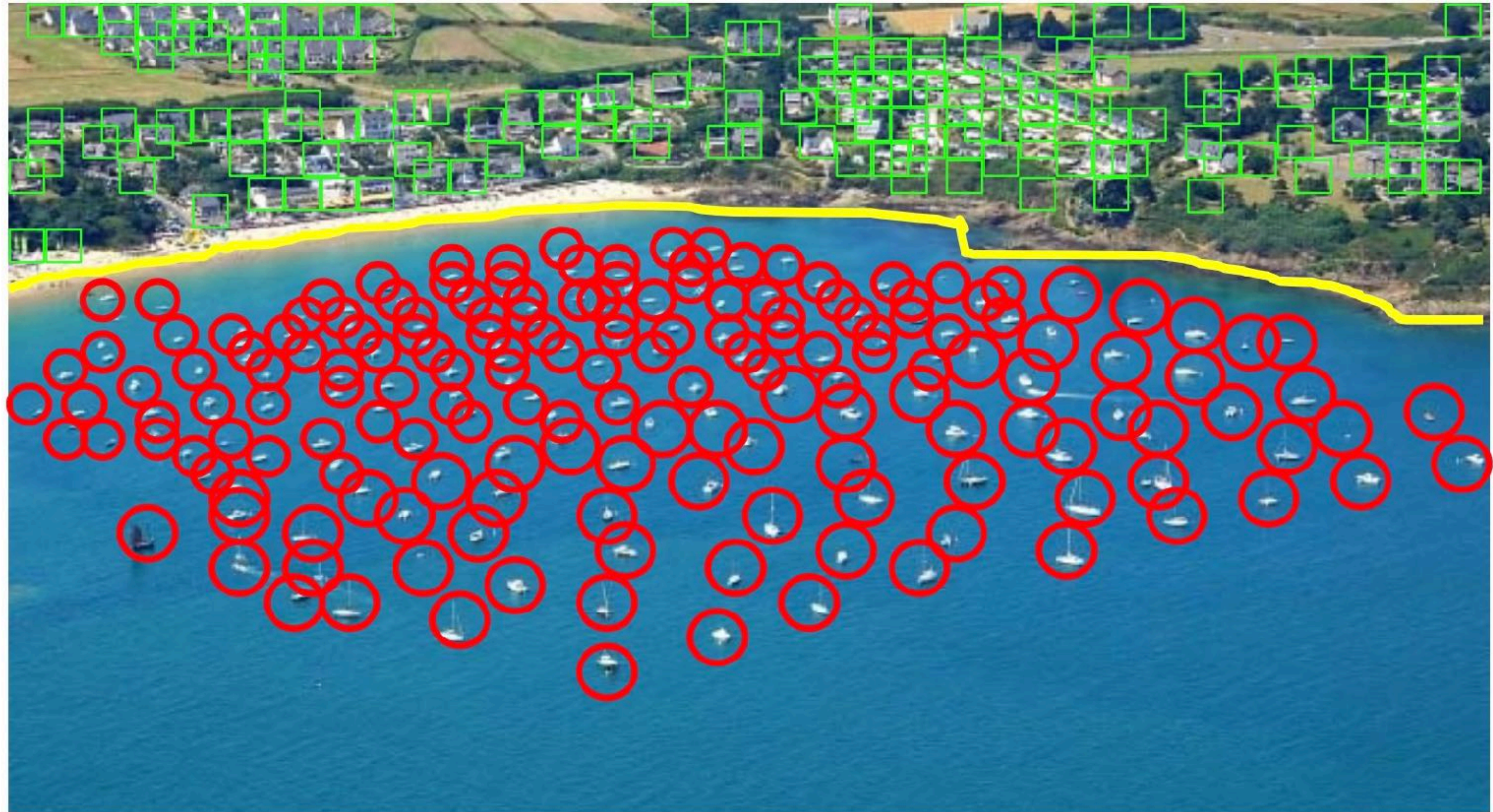# Data Science in Practice

## *Support Vector Machines*

# SVM

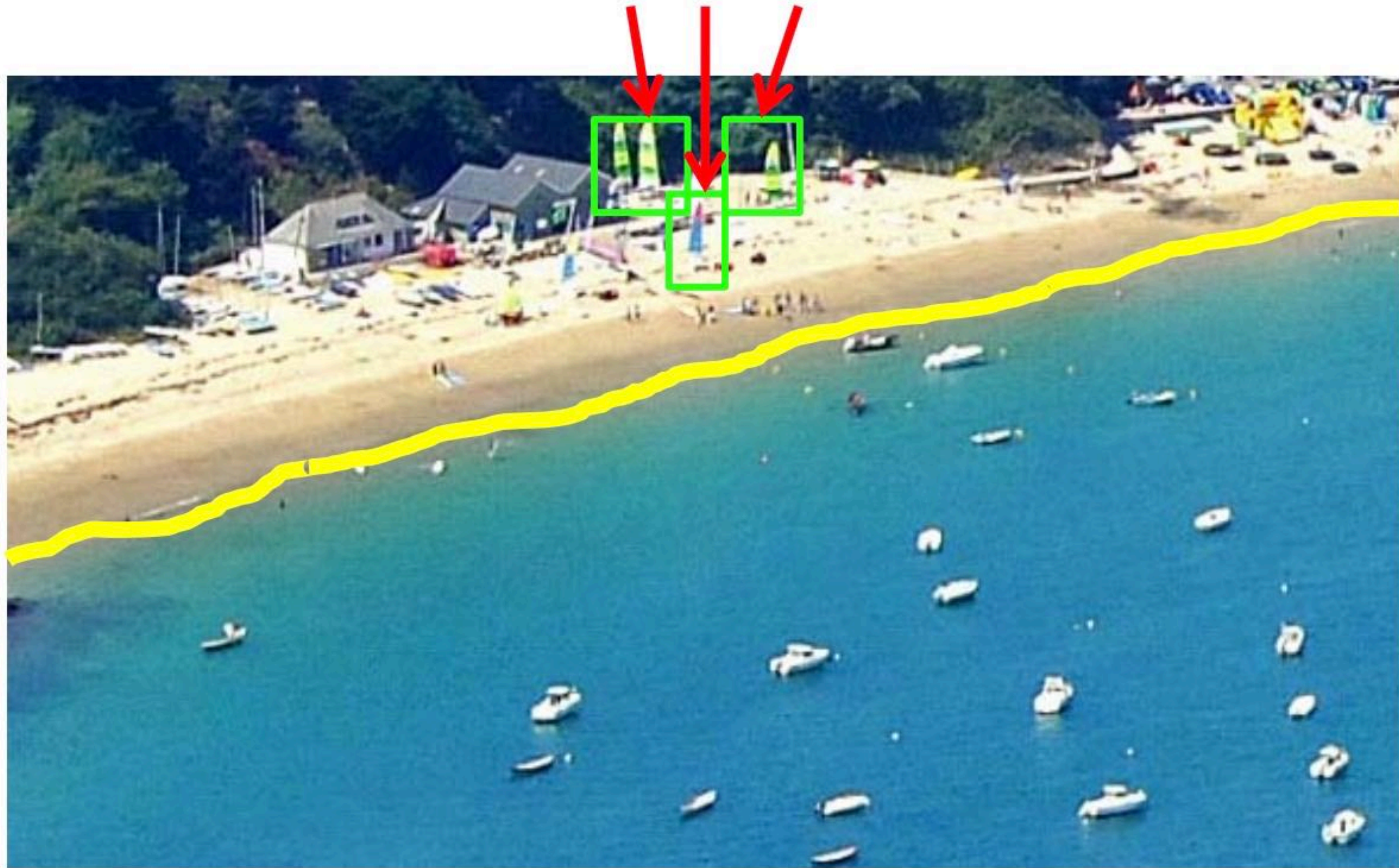**When classifying isn't easily linear,**
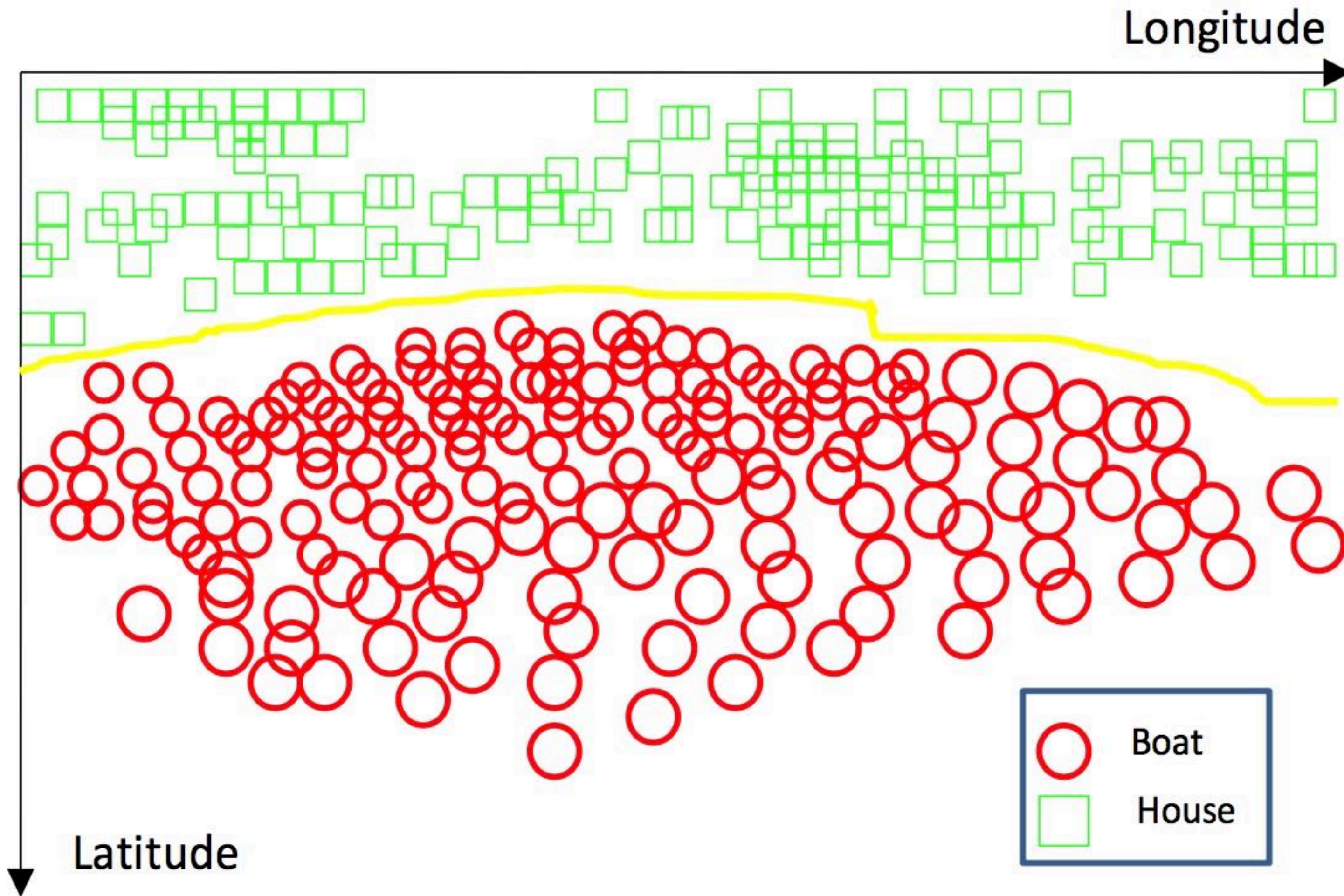*make it linear.*

# SVM

# SVM

# SVM



These boats will be misclassified as houses

# SVM



Shoreline is your decision boundary

Source: Statnikov, Hardin, Guyon, Aliferis *AMIA* 2009

# SVM



NOT LINEAR

Source: Statnikov, Hardin, Guyon, Aliferis *AMIA* 2009

# SVM

# SVM



Classify!

# SVM

Classify!

# SVM



Classify!

**DONE**

UC San Diego

# SVM



Classify!

~~DONE~~

# SVM



This is not the optimal separator

UC San Diego

# SVM



How do we find the "optimal"?

# SVM



These are your
"support vectors"

# SVM



Support vectors are the points nearest to the plane

UC San Diego

# SVM



margins

These are your margins

# SVM



**Goal:** maximize the margins ($m$) such that the decision boundary is as far away from the data of both classes as possible
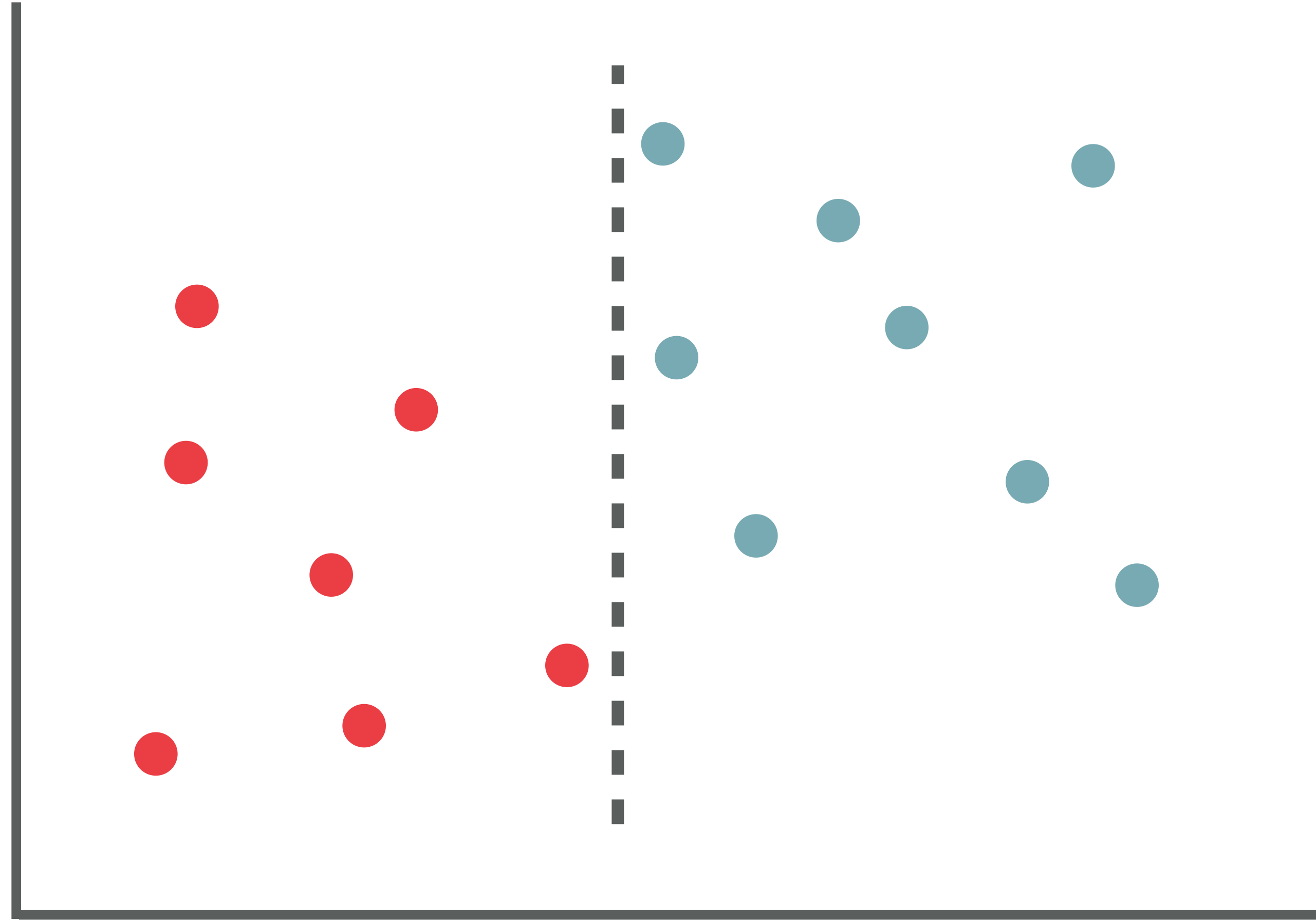
# SVM



Classify!

**DONE**

# SVM

Okay… but…

# SVM
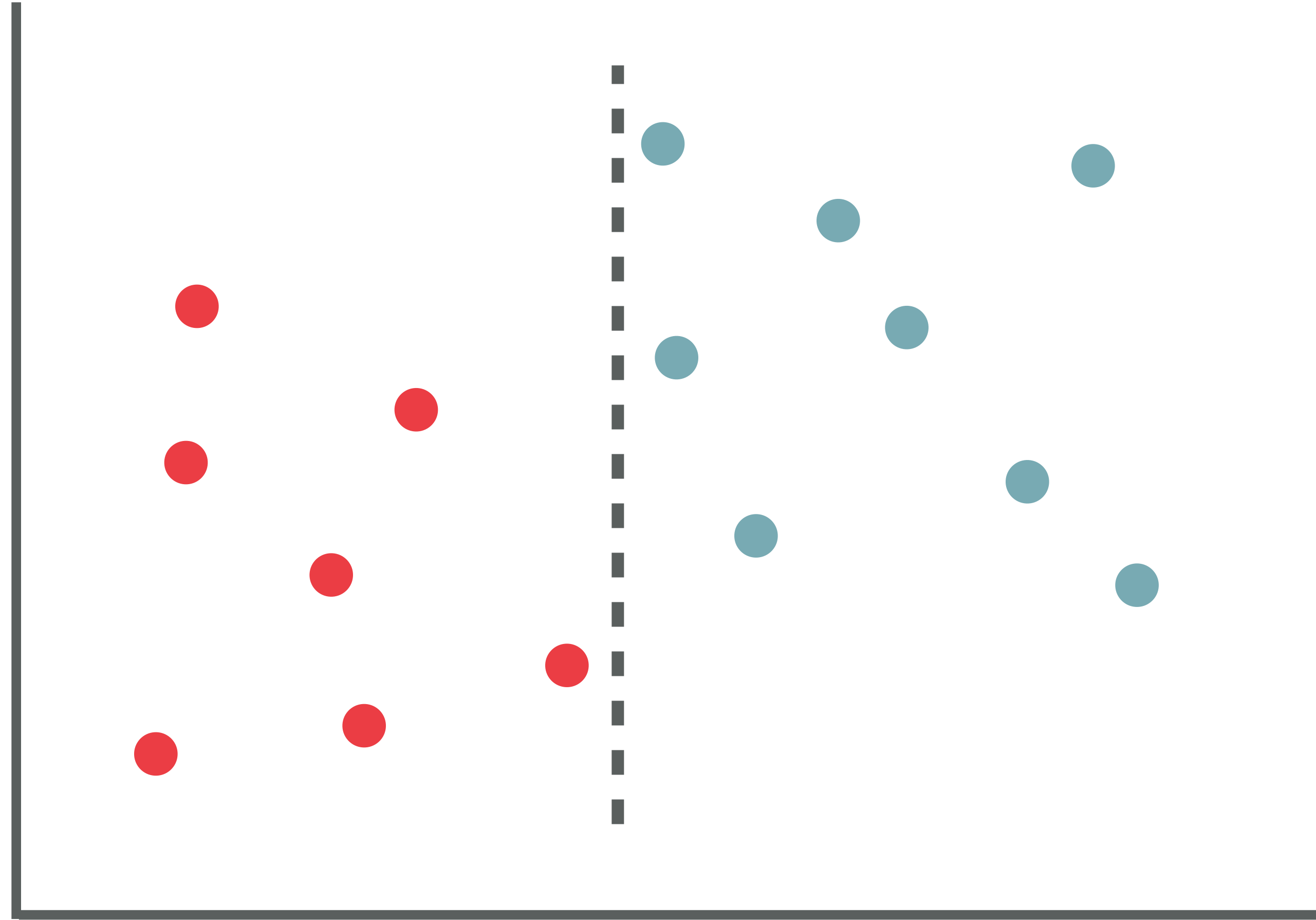


Classify!

# SVM



Classify!

# SVM



Classify!

# SVM



Classify!

# SVM



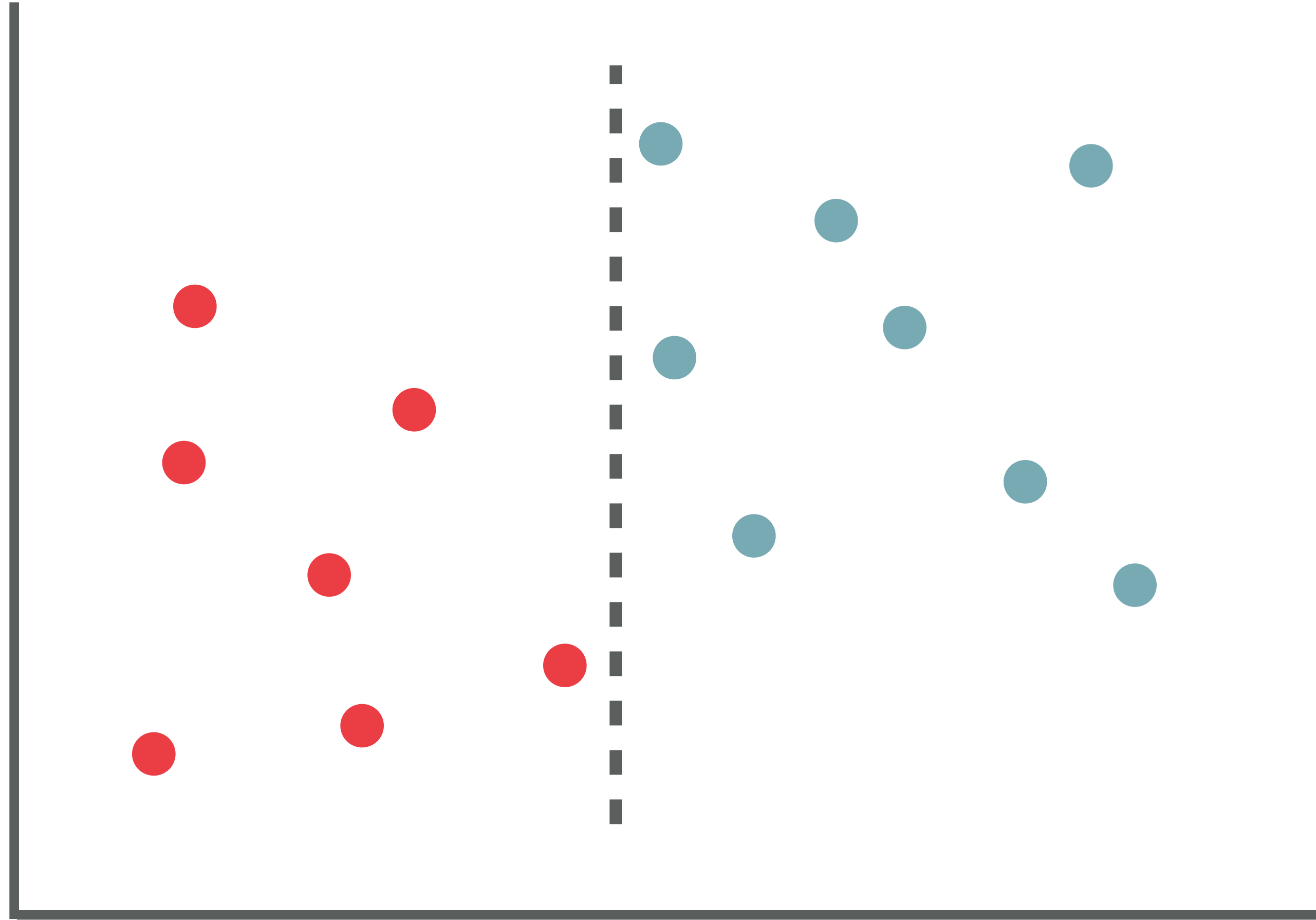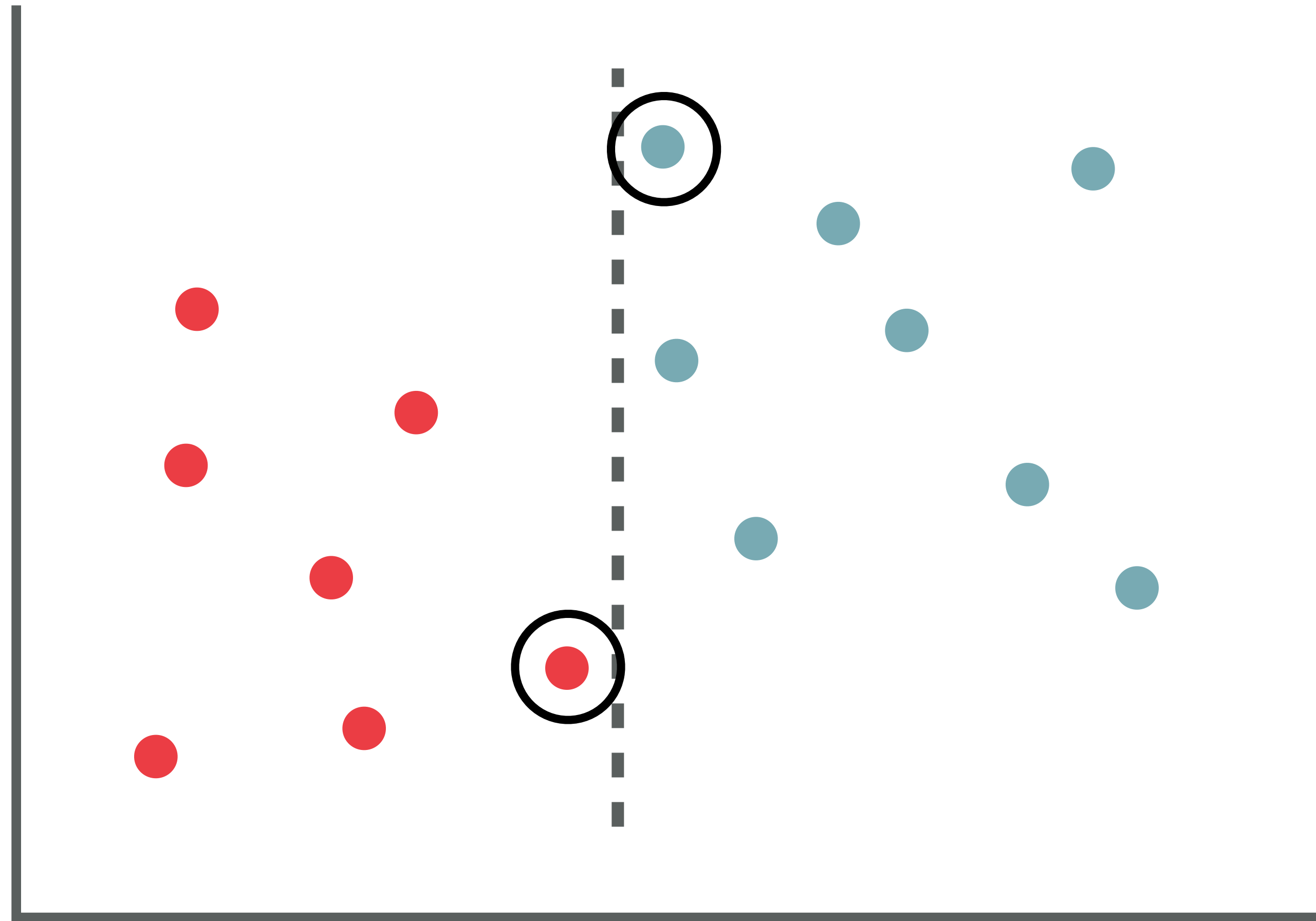fy!

# SVM

# SVM



If we have 3 training points in 2D space, no matter how the points are labeled we can perfectly classify them using a linear algorithm

UC San Diego

# SVM



But add a fourth point, and this can fail

# SVM



**SVM:** Total hack to turn data that are not optimally suited for linear classification into a linearly classifiable problem by warping them into higher dimensions

UC San Diego

# SVM



If we have *N* training points in (*N-1*)D space, no matter how the points are labeled we can perfectly classify them using a linear algorithm

# SVM



But rarely do you need to go to $N$-1 space to do this

# SVM



Classify!

# Common SVM kernels

**Examples:**

$$K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$
Linear kernel

$$K(\vec{x}_i, \vec{x}_j) = \exp(-\gamma \|\vec{x}_i - \vec{x}_j\|^2)$$
Gaussian kernel

$$K(\vec{x}_i, \vec{x}_j) = \exp(-\gamma \|\vec{x}_i - \vec{x}_j\|)$$
Exponential kernel

$$K(\vec{x}_i, \vec{x}_j) = (p + \vec{x}_i \cdot \vec{x}_j)^q$$
Polynomial kernel

$$K(\vec{x}_i, \vec{x}_j) = (p + \vec{x}_i \cdot \vec{x}_j)^q \exp(-\gamma \|\vec{x}_i - \vec{x}_j\|^2)$$
Hybrid kernel

$$K(\vec{x}_i, \vec{x}_j) = \tanh(k\vec{x}_i \cdot \vec{x}_j - \delta)$$
Sigmoidal

# COGS 108
## Data Science in Practice

*Geospatial Analyses*

# Geocoding

**Geocoding** is the process of converting addresses (like a street address) into geographic coordinates (like latitude and longitude), which you can use to place markers on a map, or position the map.

Source: Google

# Google Maps API

**Client-side geocoding**, which is executed in the browser, generally in response to user action. The Google Maps JavaScript API provides classes that make the requests for you.

UC San Diego

# Google Maps API

**When to use client-side geocoding**

The short answer is "almost always." The reasons are:

- Client-side request and response provide a faster, more interactive experience for users.
- A client-side request can include information that improves geocoding quality: user language, region, and viewport.

Source: Google

# Google Maps API

**HTTP server-side geocoding**, which allows your server to directly query Google's servers for geocodes. The Google Maps Geocoding API is the web service that provides this functionality. Typically, you integrate this service with other code that is running server-side.

UC San Diego

# Google Maps API

**When to use server-side geocoding**
Server-side geocoding is best used for applications that require you to geocode addresses without input from a client. A common example is when you get a dataset that comes independently of user input, for instance if you have a fixed, finite, and known set of addresses that need geocoding. Server-side geocoding can also be useful as a backup for when client-side geocoding fails.

UC San Diego

# Google Maps API

The Google Maps Geocoding API has the following limits in place:

**Standard Usage Limits**

Users of the standard API:

- 2,500 free requests per day, calculated as the sum of client-side and server-side queries.

- 50 requests per second, calculated as the sum of client-side and server-side queries.

UC San Diego

# Geocoding

**Reverse geocoding** is the process of converting geographic coordinates into a human-readable address. The Google Maps Geocoding API's reverse geocoding service also lets you find the address for a given place ID.

UC San Diego

Source: Google

# Geocoding

**Jupyter Demo!**

# Administrative stuff

- **http://cape.ucsd.edu/**

- **COGS: https://goo.gl/dmujYA**
- The cognitive science department runs it's own evaluations. You can submit evaluations for the instructor, your section TA, and/or any TA you had enough interaction with to provide feedback. Please do fill these out - they are very useful for us to get feedback on our teaching, and provide feedback to the department, and this is especially useful for working on this new course.

# Bradley Voytek, Ph.D.
UC San Diego
Cognitive and Neural Dynamics Laboratory

Department of Cognitive Science
Neurosciences Graduate Program
Halıcıoğlu Data Science Institute

bvoytek@ucsd.edu
@bradleyvoytek

UC San Diego