

Lecture: K-Means Clustering

Rodrigo Canaan

rodrigocanaan.com

rocanaan@gmail.com

 @rocanaan

Cal Poly State University, May 26th

Slides and Code

Types of Machine Learning

- **Supervised Learning**

Learning to predict values or classify objects based on labeled data

- **Unsupervised Learning**

Learning patterns based on unlabeled data

- **Reinforcement Learning**

Learning to act intelligently based on interaction with an environment

Types of Machine Learning

- **Supervised Learning**

Learning to predict values or classify objects based on labeled data

- **Unsupervised Learning**

Learning patterns based on unlabeled data

- **Reinforcement Learning**

Learning to act intelligently based on interaction with an environment

Why do Unsupervised Learning?

- To gain insight about the data**
- To compress and/or visualize the data**
- To generate examples that are like the data**
- To label data for downstream tasks**

Examples in various domains

- **Recommender Systems**

Group similar products and/or users together

- **Visual Processing**

Group similar images even without labels

- **Games**

Build models of different player styles

- **Science**

Create a taxonomy of phenomena (e.g. stars) based on their observed properties (magnitude, spectrum, distance...)

Unsupervised Learning Tasks

- **Clustering**

Partitioning the data into “clusters” based on a measure of similarity

- **Dimensionality reduction**

Transforming the data from a high dimension to a low dimension

- **Generative models**

Learning a “model” that can be used to produce new examples that are similar to the data

Unsupervised Learning Tasks

- **Clustering**

Partitioning the data into “clusters” based on a measure of similarity

- **Dimensionality reduction**

Transforming the data from a high dimension to a low dimension

- **Generative models**

Learning a “model” that can be used to produce new examples that are similar to the data

Visual Example

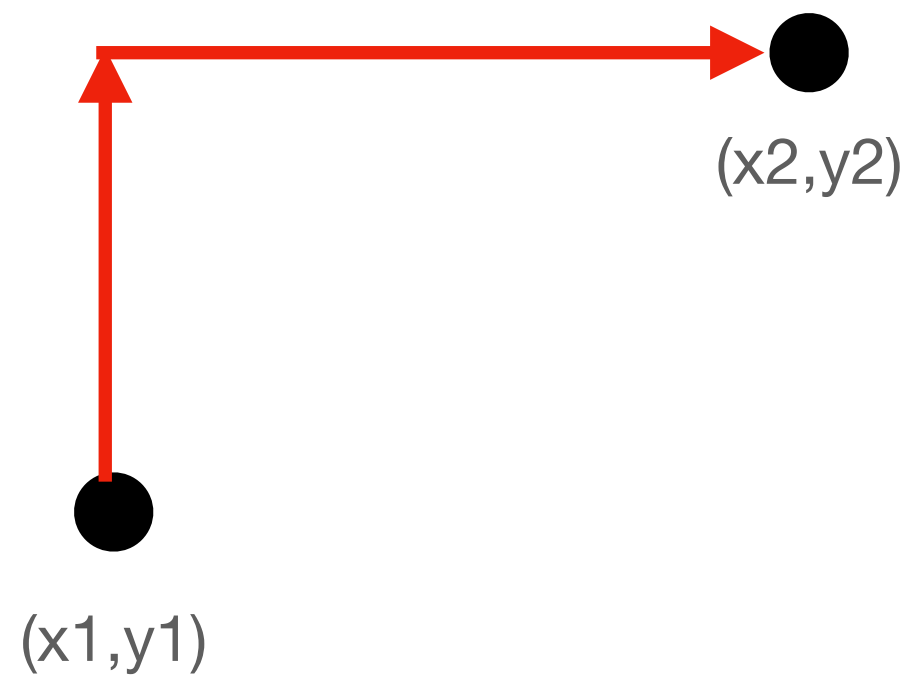
Clustering - Distances

- **What does it mean for two data points to be “similar”?**

Usually, it means to have a small “distance” according to some metric

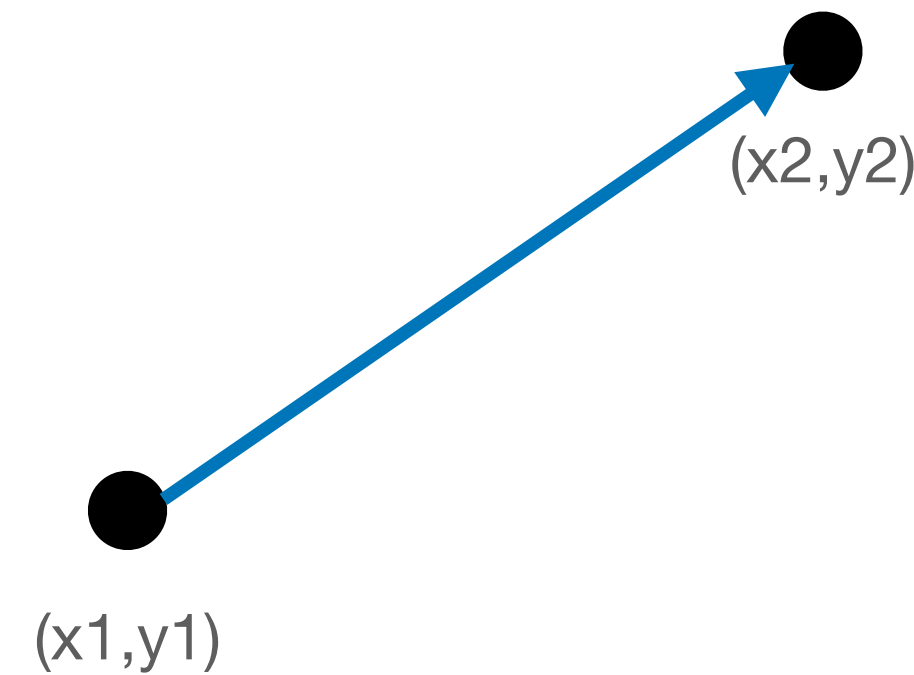
Examples (in 2-dimensions)

Manhattan Distance



$$D_M = |x_2 - x_1| + |y_2 - y_1|$$

Euclidean Distance



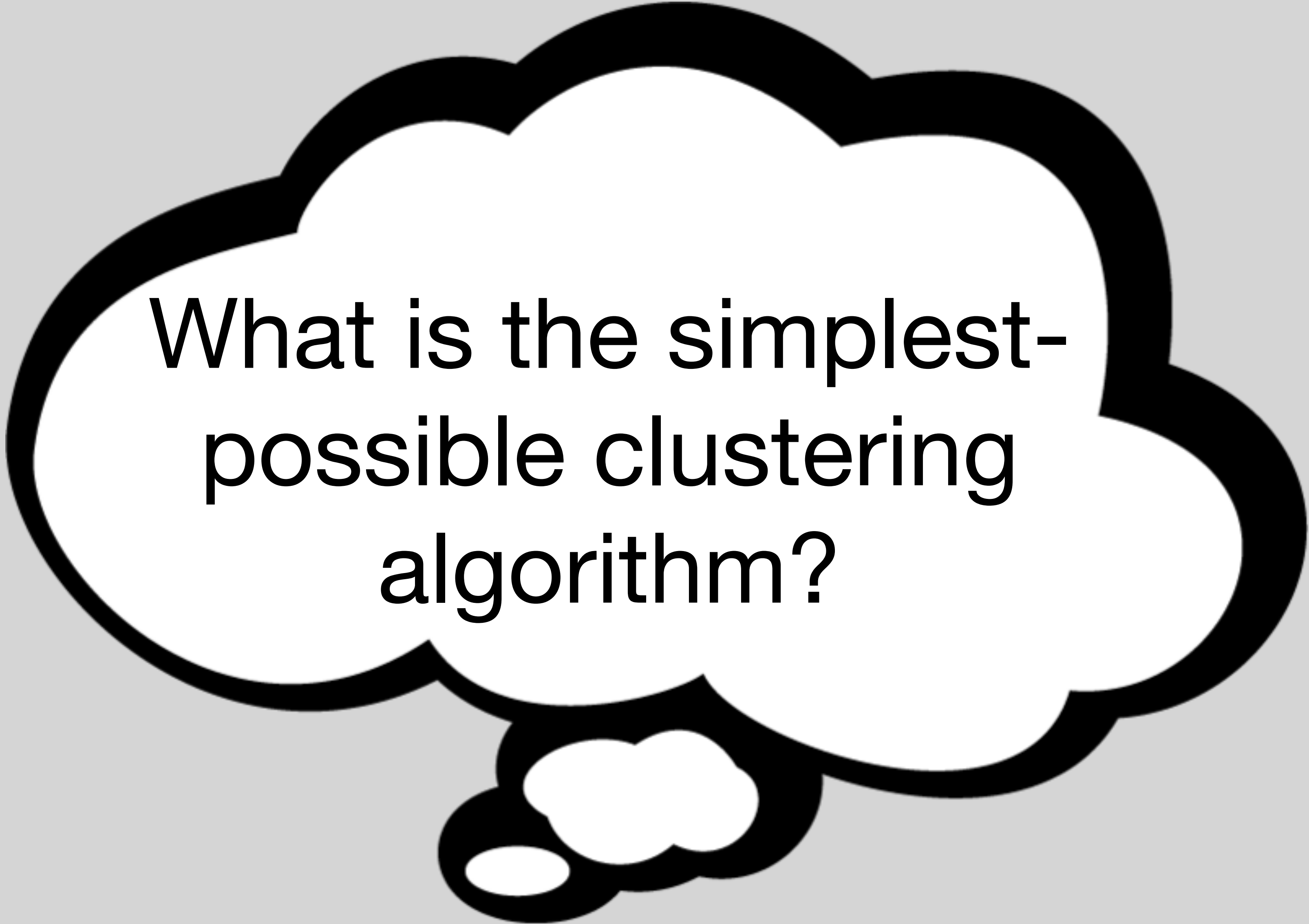
$$D_E = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Clustering - Partitions

What does it mean for a partition to be “good”?

Intuitively:

- Elements in the same cluster are similar.
- Elements in different clusters are dissimilar.



What is the simplest-
possible clustering
algorithm?

Clustering by exhaustive search

1. Generate all ____ possible partitions
2. Calculate clustering metric for each partition
3. Return the partition with best clustering metric

Clustering by exhaustive search

1. Generate all k^n possible partitions
2. Calculate WCSS for each partition
3. Return the partition with minimum WCSS

Clustering by exhaustive search

1. Generate all k^n possible partitions
2. Calculate WCSS for each partition
3. Return the partition with minimum WCSS



Very expensive!

Need **heuristic** for approximate (but faster) solution!

Clustering by exhaustive search

1. Generate all k^n possible partitions
2. Calculate WCSS for each partition
3. Return the partition with minimum WCSS

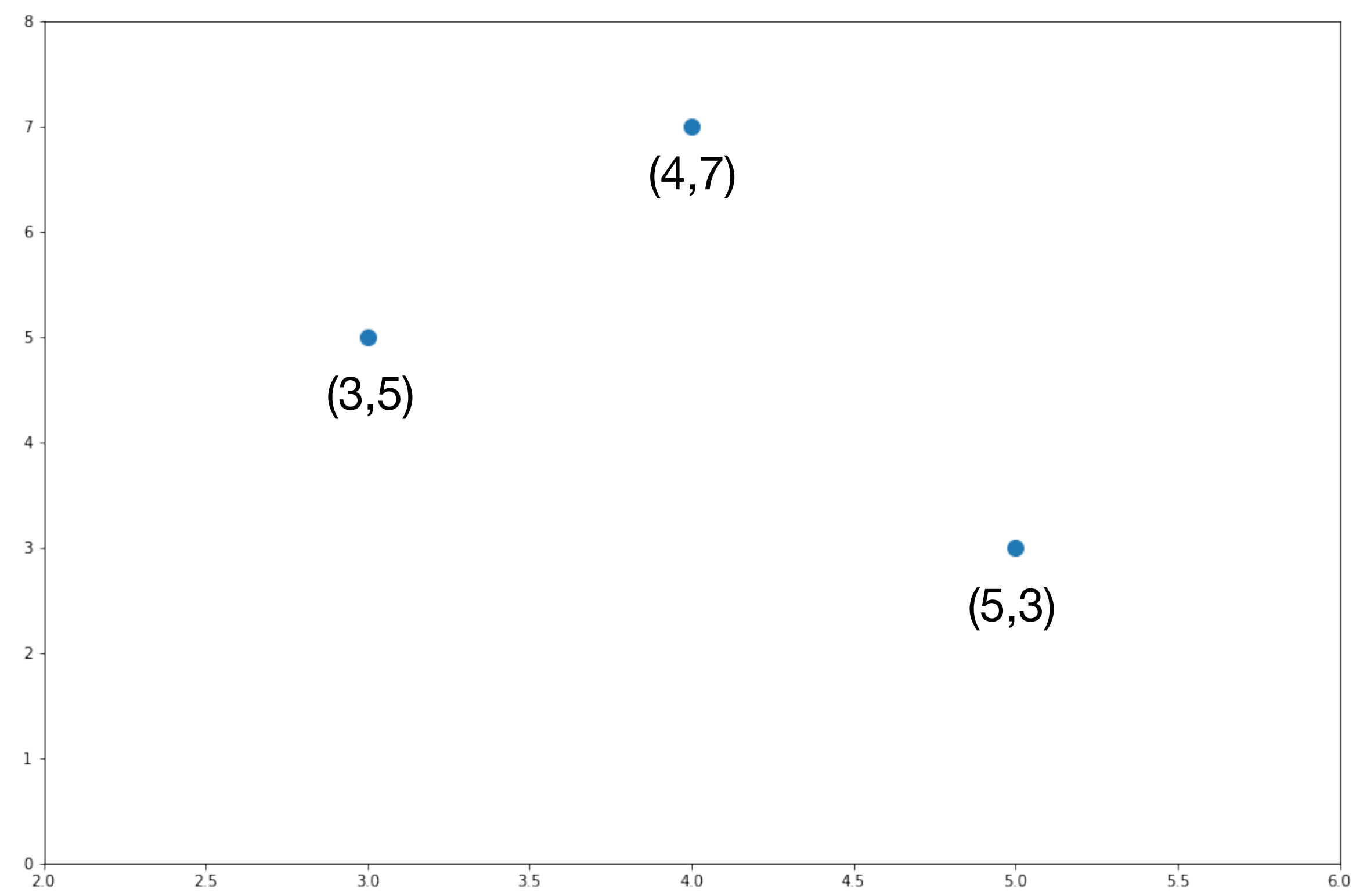


Very expensive!

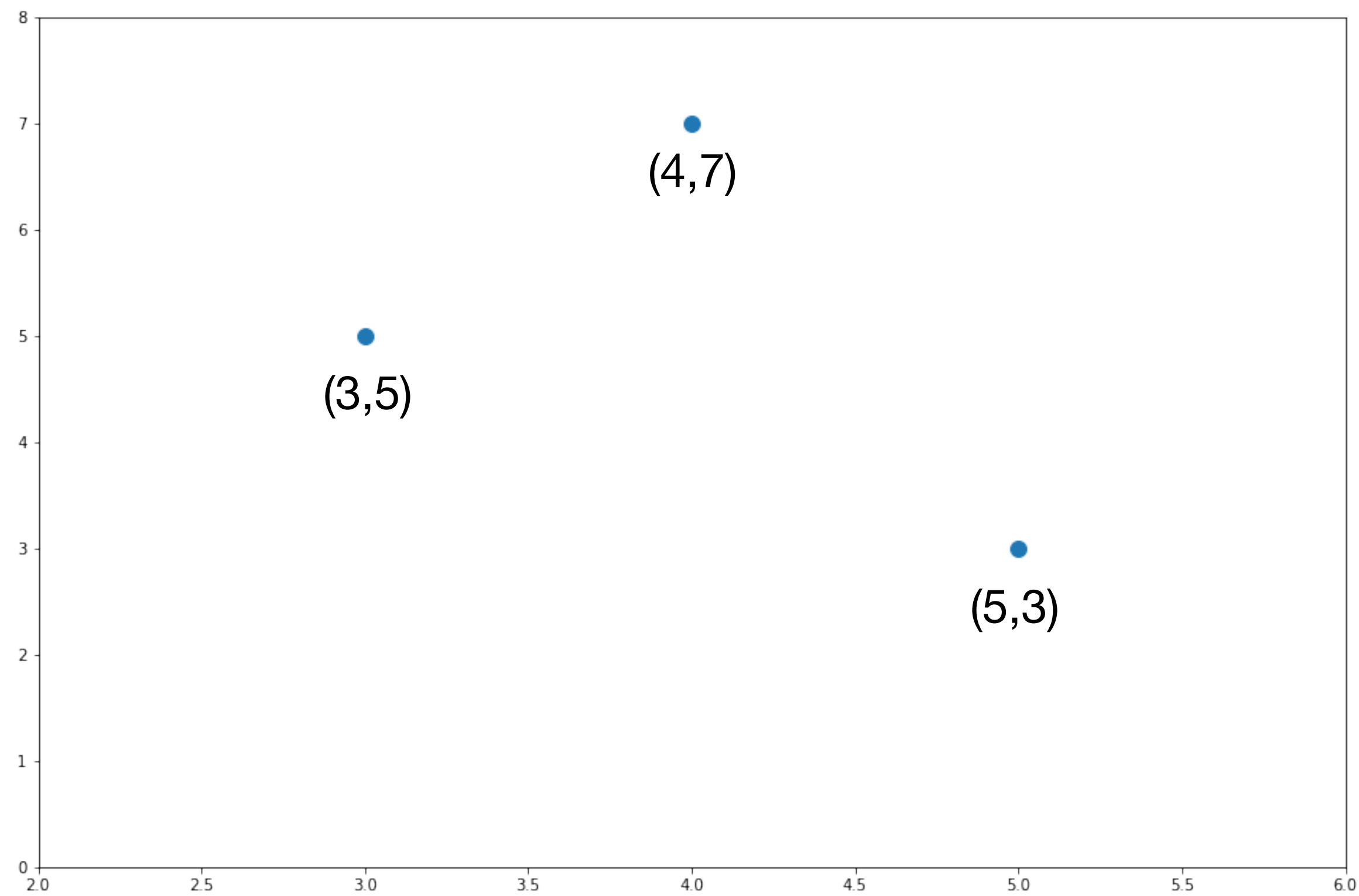
Need **heuristic** for approximate (but faster) solution!

Example: assign each point to cluster with nearest centroid

Centroid Calculation



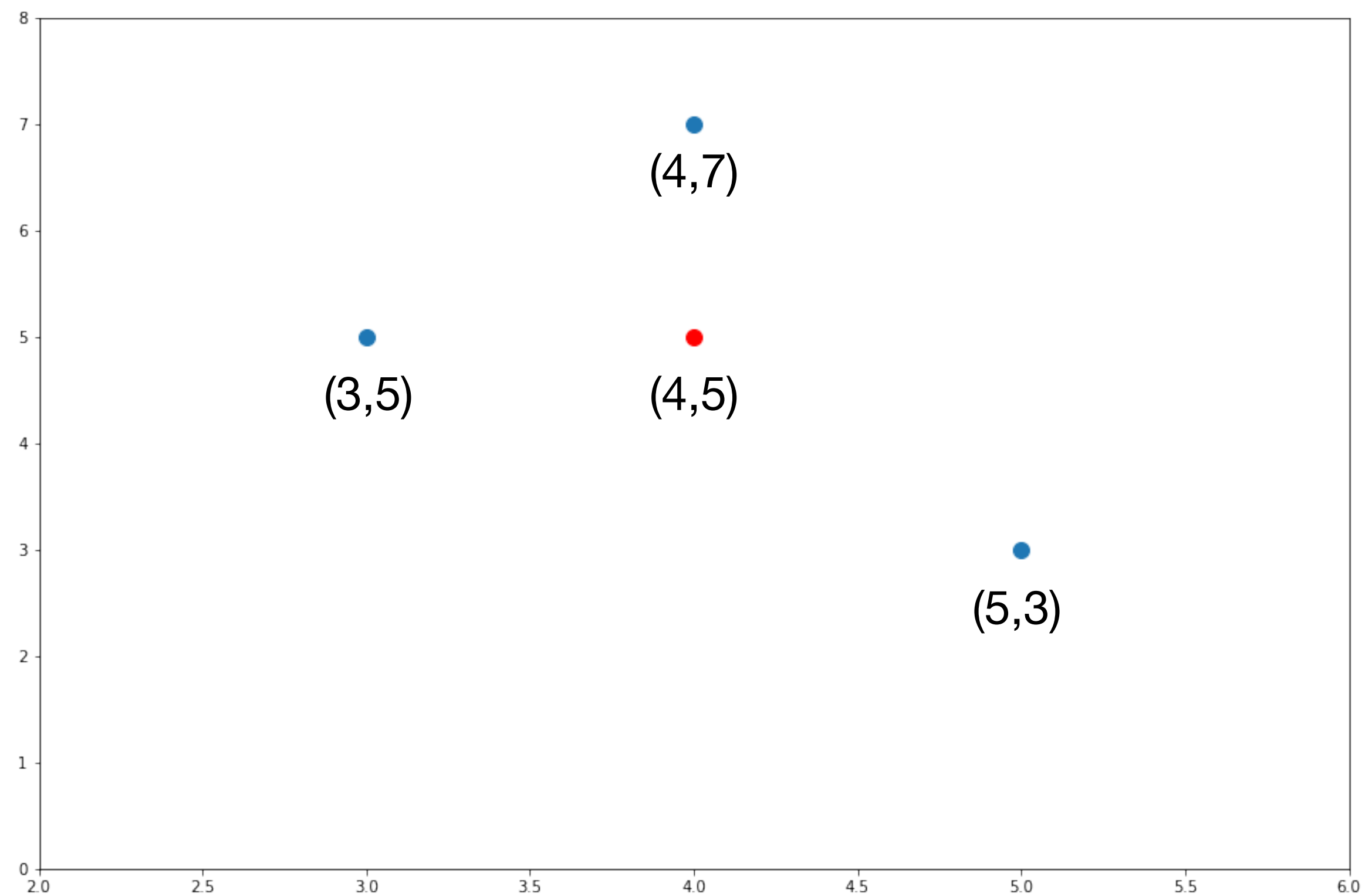
Centroid Calculation



$$x_mean = (3+4+5)/3 = 4$$

$$y_mean = (5+7+3)/3 = 5$$

Centroid Calculation



$$x_mean = (3+4+5)/3 = 4$$

$$y_mean = (5+7+3)/3 = 5$$

$$\text{Centroid} = (4,5)$$

Outline of K-Means Clustering Algorithm

- 1. Generate k initial centroids**
- 2. Re-assign points based on heuristic**
Each point is assigned to cluster with nearest centroid
- 3. Recalculate centroids**
- 4. Repeat (2, 3) until no re-assignments**
Alternatively, up to some maximum number of repetitions

(Naive) K-Means Pseudo-Code

```
1 # Pseudo-code of K-Means algorithm
2
3 Function k-means (data,k)
4     P ← initialize_partition(data,k)
5     stop ← False
6
7     while not stop # runs until no new assignments
8         # optionally, until some max_iterations is reached
9         stop = True
10        P_new ← empty_partition(k)
11        for d in data
12            new_label ← P.find_min_sq_distance_cluster(d)
13            P_new.add_element(d,new_label)
14            if new_cluster <> P.get_label(d) # new assignment, don't stop!
15                stop = False
16            Endif
17        Endfor
18        P ← P_new
19        P.compute_centroids()
20    Endwhile
21
22    return P
23 Endfunction
```

K-Means Initialization

How to initialize clusters?

- **Forgy method:** Choose k initial centroids randomly (from the data), assign other points according to distance to centroids
- **Random Partition:** Assign each datapoint to a random cluster label, then compute centroids
- Other options: see comparative study by ([Celebi et al., 2013](#))

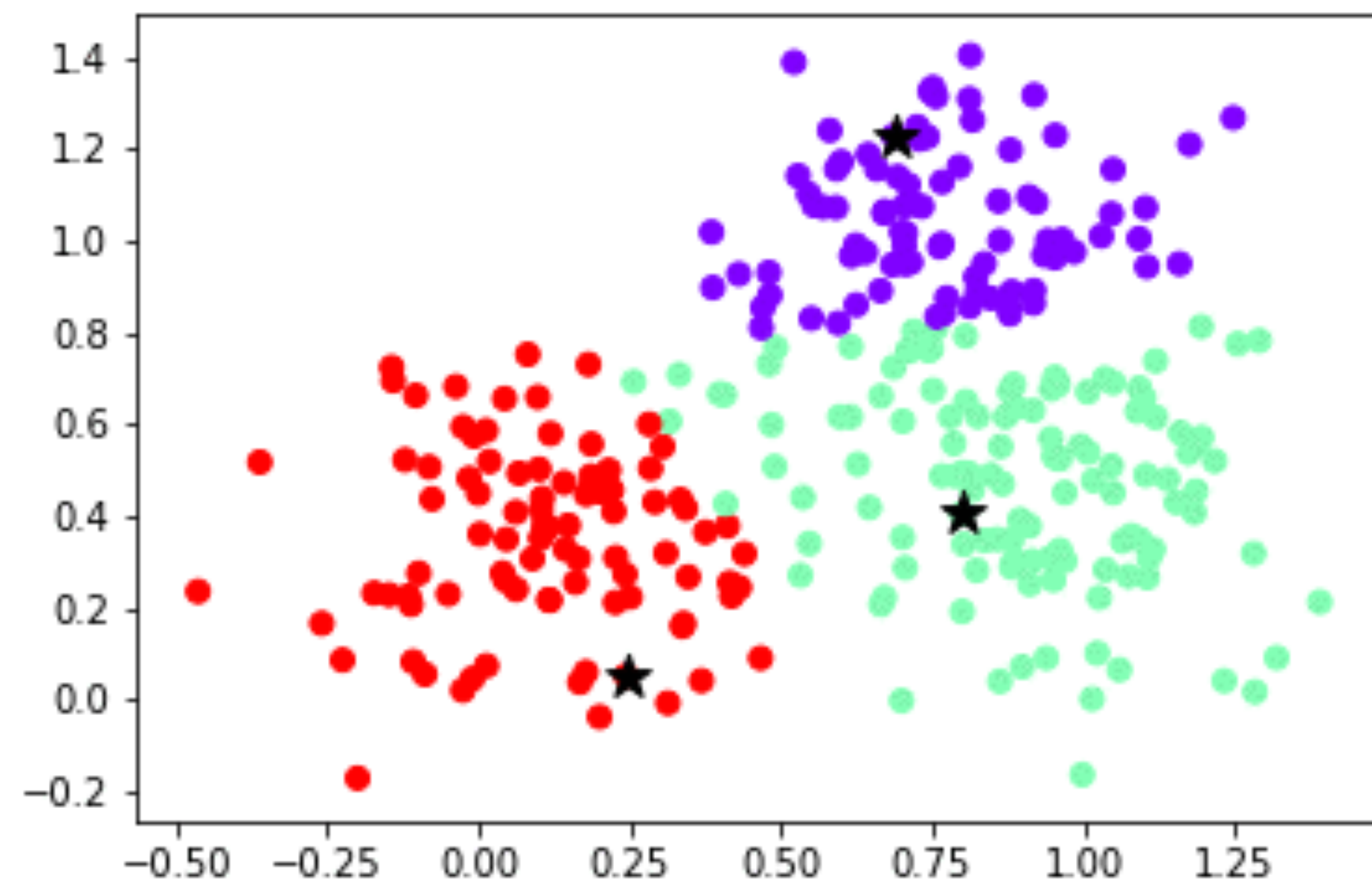
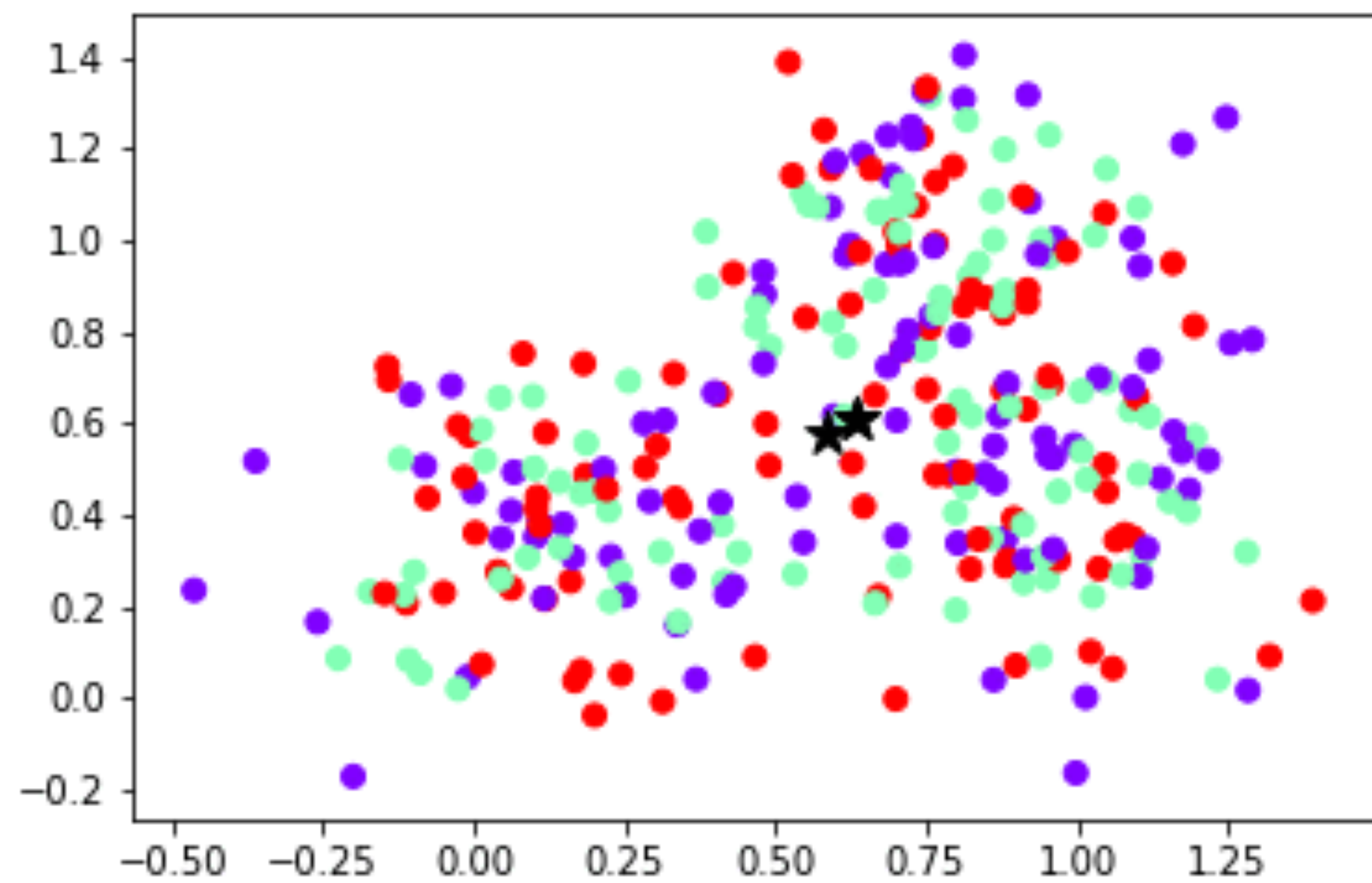
Forgy Method

```
1 # Forgy initialization
2 # Select k random points as initial centroids,
3 # then assign remaining elements to closest centroid
4
5 Function initialize_partitions(data,k)
6     P <- empty_partition(k)
7     medoids <- sample_without_replacement(data,k) # get k points from data
8     for i in 0..k-1
9         P.add_element(medoids[i],i)
10    Endfor
11    P.compute_centroids()
12
13    for d in data
14        label <- P.find_min_sq_distance_cluster(d)
15        P.add_element(d,label)
16    Endfor
17    P.compute_centroids
18    return P
19
20 Endfunction
```

Random Partition

```
1 # Random Initialization
2 # Assign each point to a random label, then compute centroids
3
4 Function initialize_partitions(data,k)
5     P <- empty_partition(k)
6
7     for d in data
8         label <- random_uniform(k) # select a random int uniformly in 0..k-1
9         P.add_element(d,label)
10    Endfor
11
12    P.compute_centroids
13    return P
14 Endfunction
```

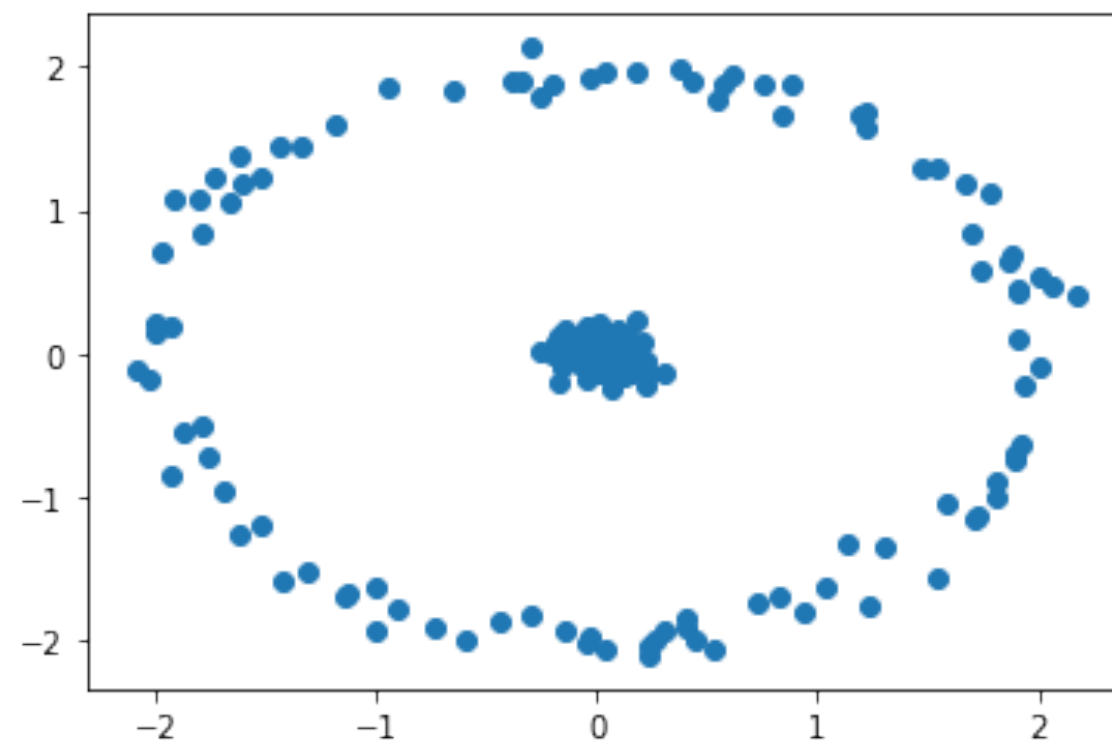

Demo



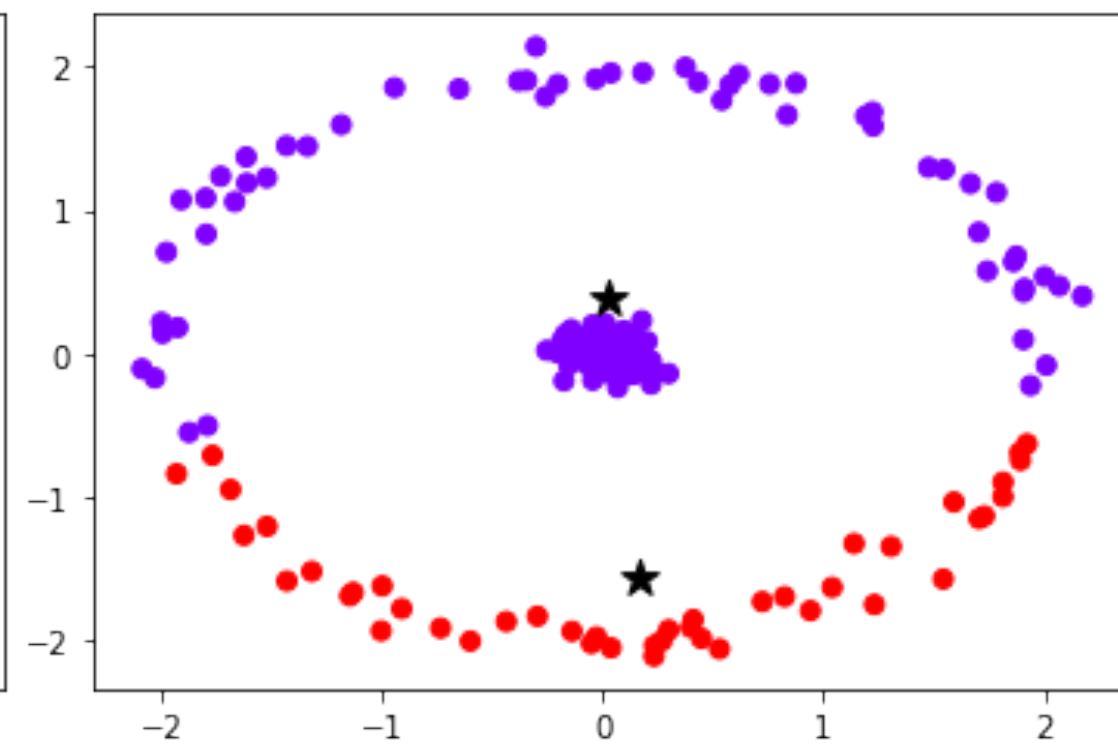
(Naive) K-Means Disadvantages

- **Struggles with some cluster shapes, sizes and outliers**

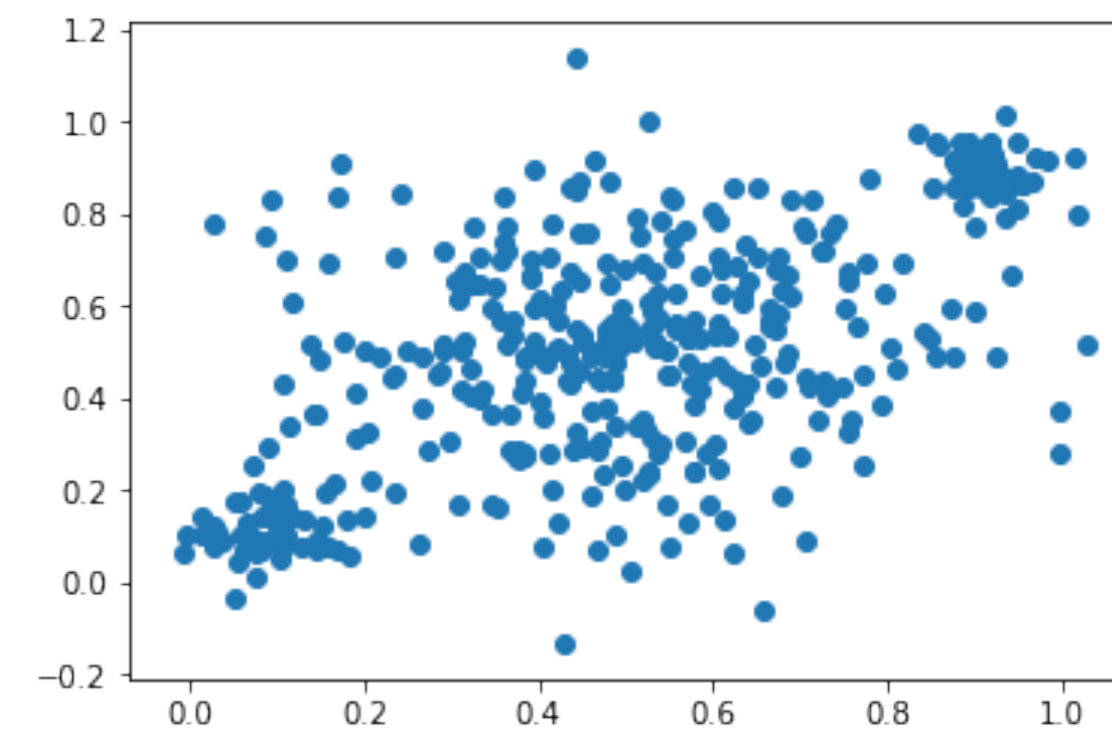
Non-Convex Clusters



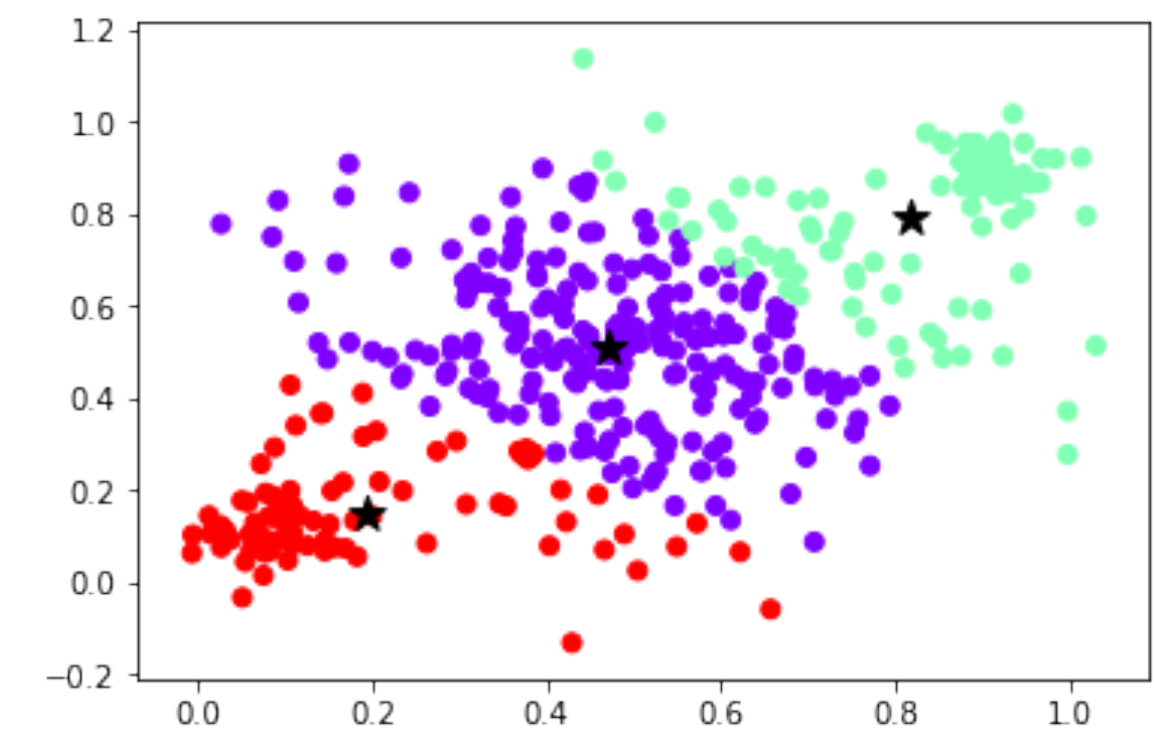
Partition with k=2



Clusters of different sizes



Partition with k=3



(Naive) K-Means Disadvantages

- **Struggles with some cluster shapes, sizes and outliers**
- **Struggles with categorical data**
- **Guaranteed convergence only to local (not global optimum)**
- **Need to specify k in advance**
- **“Curse of dimensionality”**
- **Worst-case performance?**

Time Complexity

$O(i)$
Outer-loop
iterations

$O(n)$
Inner-loop
iterations

```
1 # Pseudo-code of K-Means algorithm
2
3 Function k-means (data,k)
4     P <- initialize_partition(data,k)
5     stop <- False
6
7     while not stop # runs until no new assignments
8                   # optionally, until some max_iterations is reached
9         stop = True
10        P_new <- empty_partition(k)
11        for d in data
12            new_label <- P.find_min_sq_distance_cluster(d)
13            P_new.add_element(d,new_label)
14            if new_cluster <> P.get_label(d) # new assignment, don't stop!
15                stop = False
16        Endif
17    Endfor
18    P <- P_new
19    P.compute_centroids()
20 Endwhile
21
22 return P
23 Endfunction
```

$O(k)$

Time Complexity = $O(nki)$

$O(i)$
Outer-loop
iterations

$O(n)$
Inner-loop
iterations

```
1 # Pseudo-code of K-Means algorithm
2
3 Function k-means (data,k)
4     P <- initialize_partition(data,k)
5     stop <- False
6
7     while not stop # runs until no new assignments
8                   # optionally, until some max_iterations is reached
9         stop = True
10        P_new <- empty_partition(k)
11        for d in data
12            new_label <- P.find_min_sq_distance_cluster(d)
13            P_new.add_element(d,new_label)
14            if new_cluster <> P.get_label(d) # new assignment, don't stop!
15                stop = False
16        Endif
17    Endfor
18    P <- P_new
19    P.compute_centroids()
20 Endwhile
21
22 return P
23 Endfunction
```

$O(k)$

(Naive) K-Means Disadvantages

- **Struggles with some cluster shapes, sizes and outliers**
- **Guaranteed convergence only to local (not global optimum)**
- **Need to specify k in advance**
- **“Curse of dimensionality”**
- **High worst-case performance**

K-Means Advantages

- **Simple to implement**
- **Good performance in many practical scenarios**
- **Adaptations can handle outliers, different shapes and sizes, higher dimensions, categorical data...**
- **Performance can be improved by non-naive implementations**

Thank you!



www.rodrigocanaan.com

rodrigo.canaan@nyu.edu

Scholar Page

