

Clustering: how to and Applications

John Monroe

10/26/18

Clustering Objective

Complex Data



Clustering

Simple conclusion

Biology

Rich bacterial variability

Taxonomic classification

Machine Learning

Haphazardly collected data

Meaningful patterns

General Science

Many degrees of freedom

Essential underlying parameters

1. Data



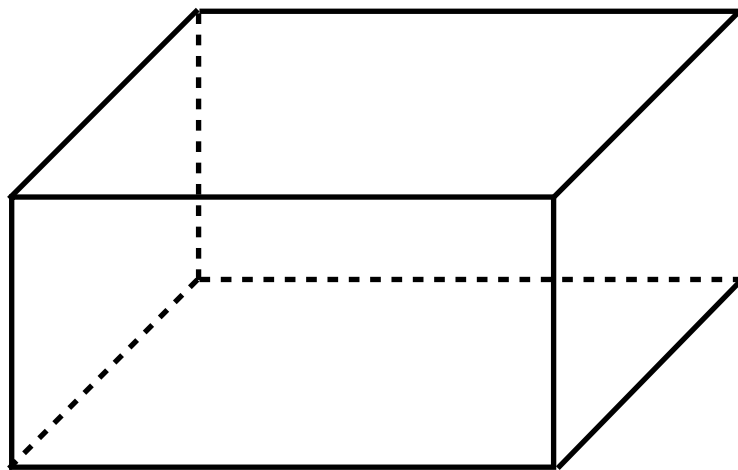
Shopping habits



- Similarity sets are not always immediately obvious
- High dimensional data is hard to visualize

Basic Ingredients

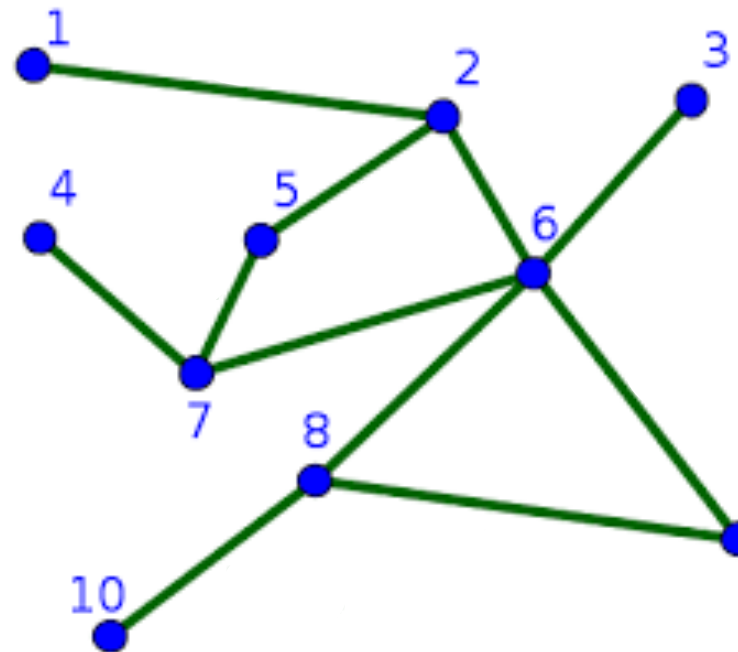
2. Distance Measure



Spatial Information

+

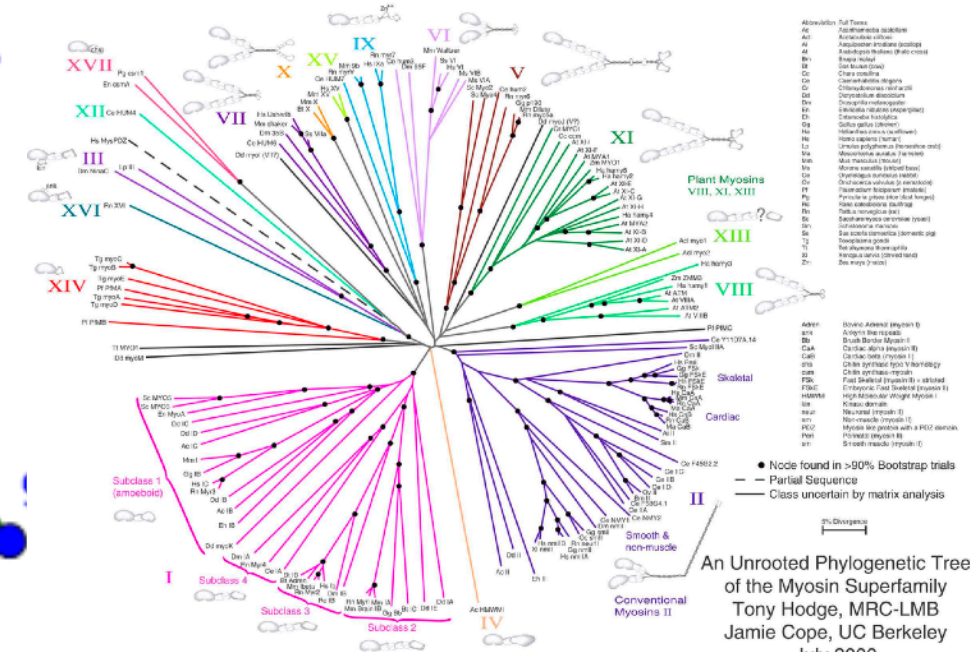
Euclidean Distance



Shopping habits

+

Similarity graph



Sequences

+

Genome counts

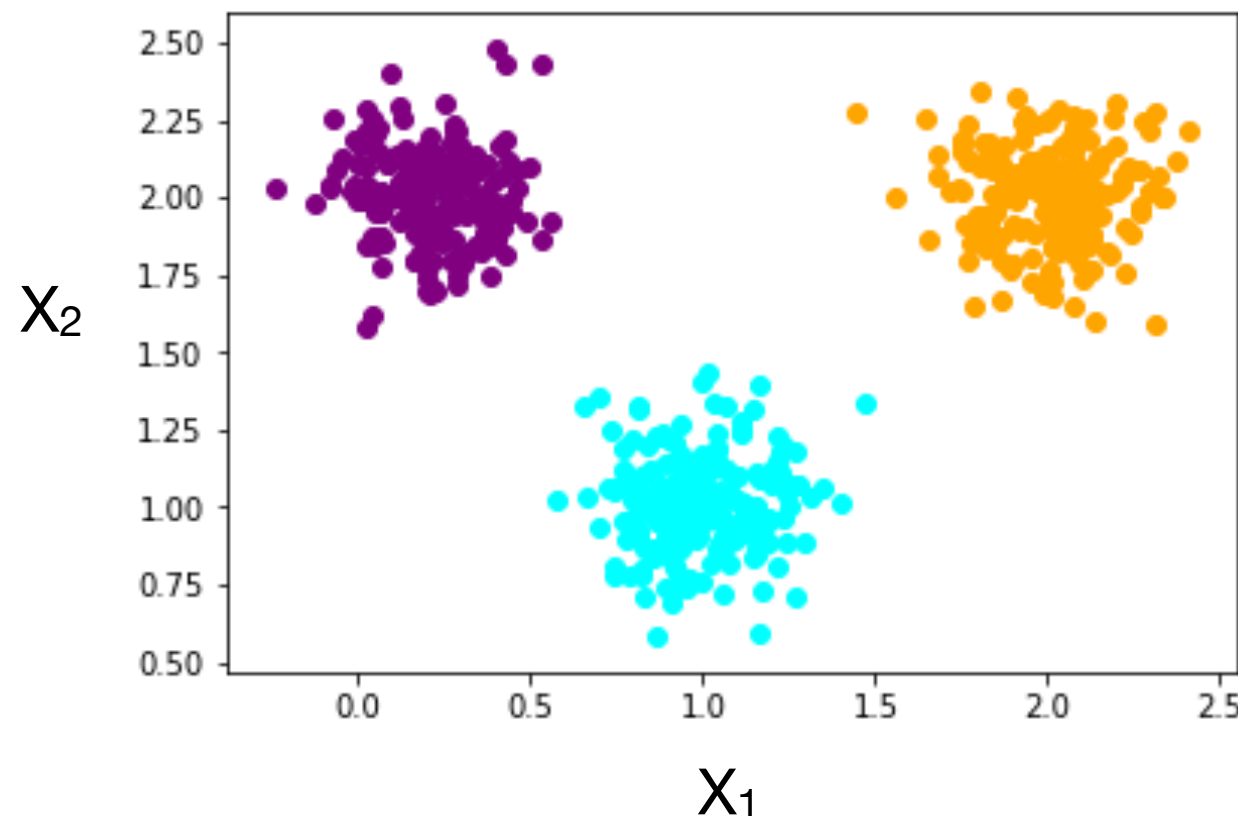
Distance measures are generally applicable to many datasets.

The k-means Algorithm

The algorithm:

```
> Initialize random cluster indices
> Repeat until convergence:
>   For each data point:
>     Distance( point, centerj )
>     new_index ← minj(distance)
>     total_distance += distance
```

The data:

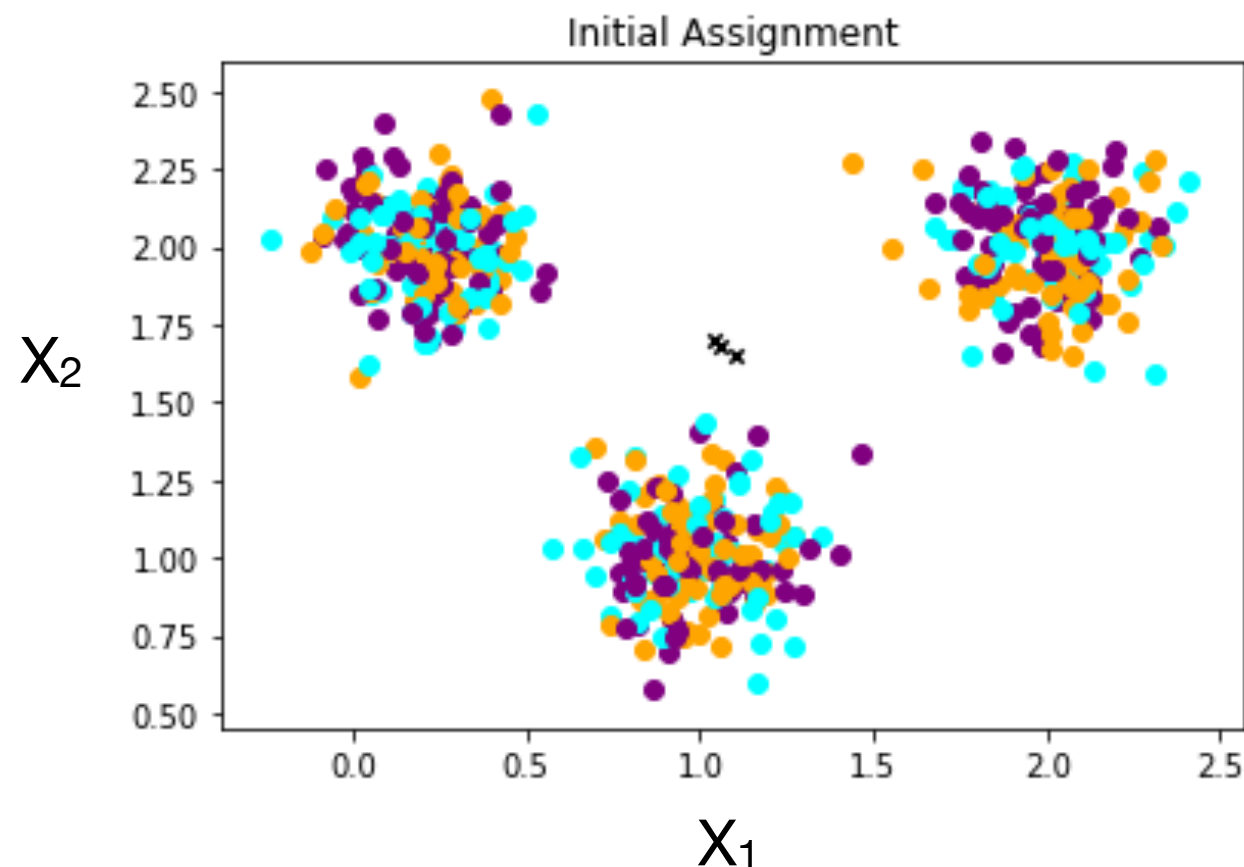


The k-means Algorithm

The algorithm:

```
> Initialize random cluster indices  
> Repeat until convergence:  
>   For each data point:  
>     Distance( point, centerj )  
>     new_index ← minj(distance)  
>     total_distance += distance
```

The data:

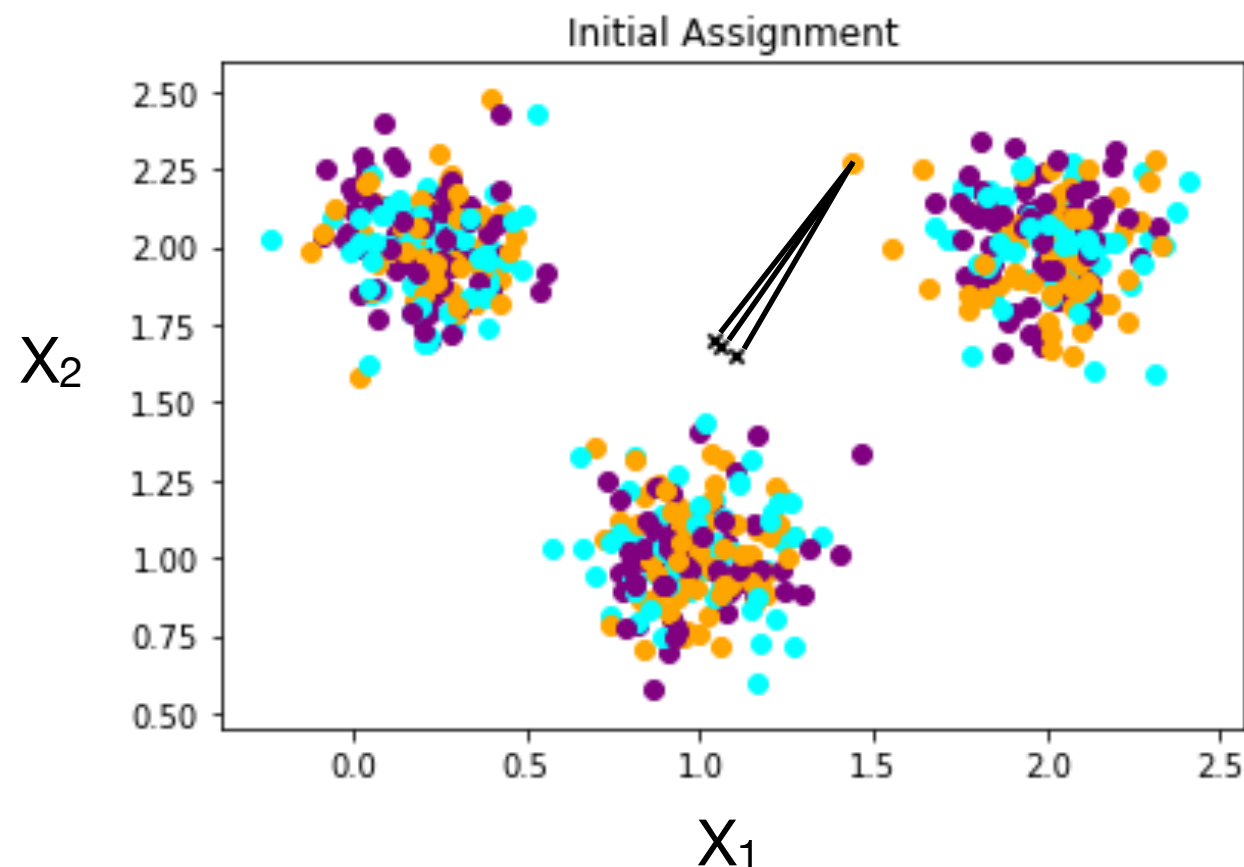


The k-means Algorithm

The algorithm:

```
> Initialize random cluster indices
> Repeat until convergence:
>   For each data point:
>     Distance( point, centerj )
>     new_index ← minj(distance)
>     total_distance += distance
```

The data:

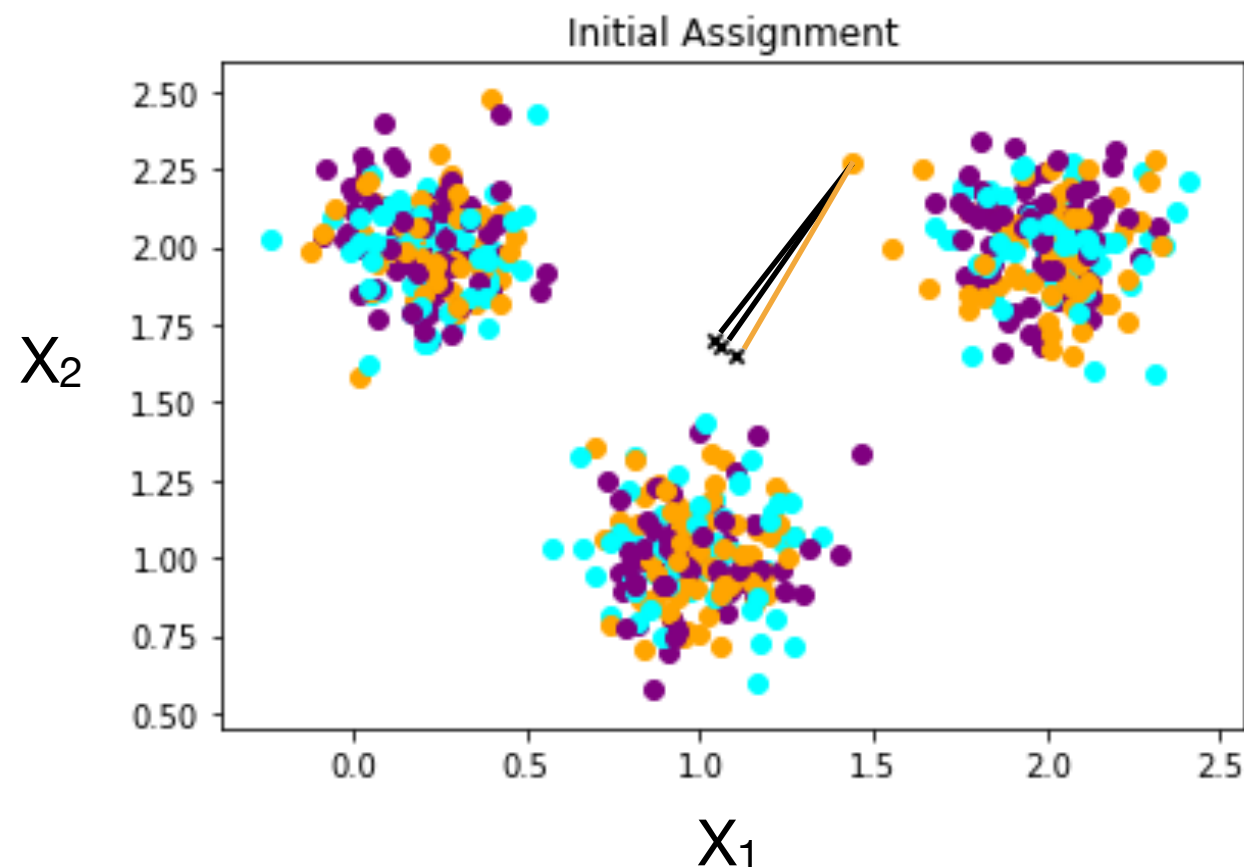


The k-means Algorithm

The algorithm:

```
> Initialize random cluster indices
> Repeat until convergence:
>   For each data point:
>     Distance( point, centerj )
>     new_index ← minj(distance)
>     total_distance += distance
```

The data:

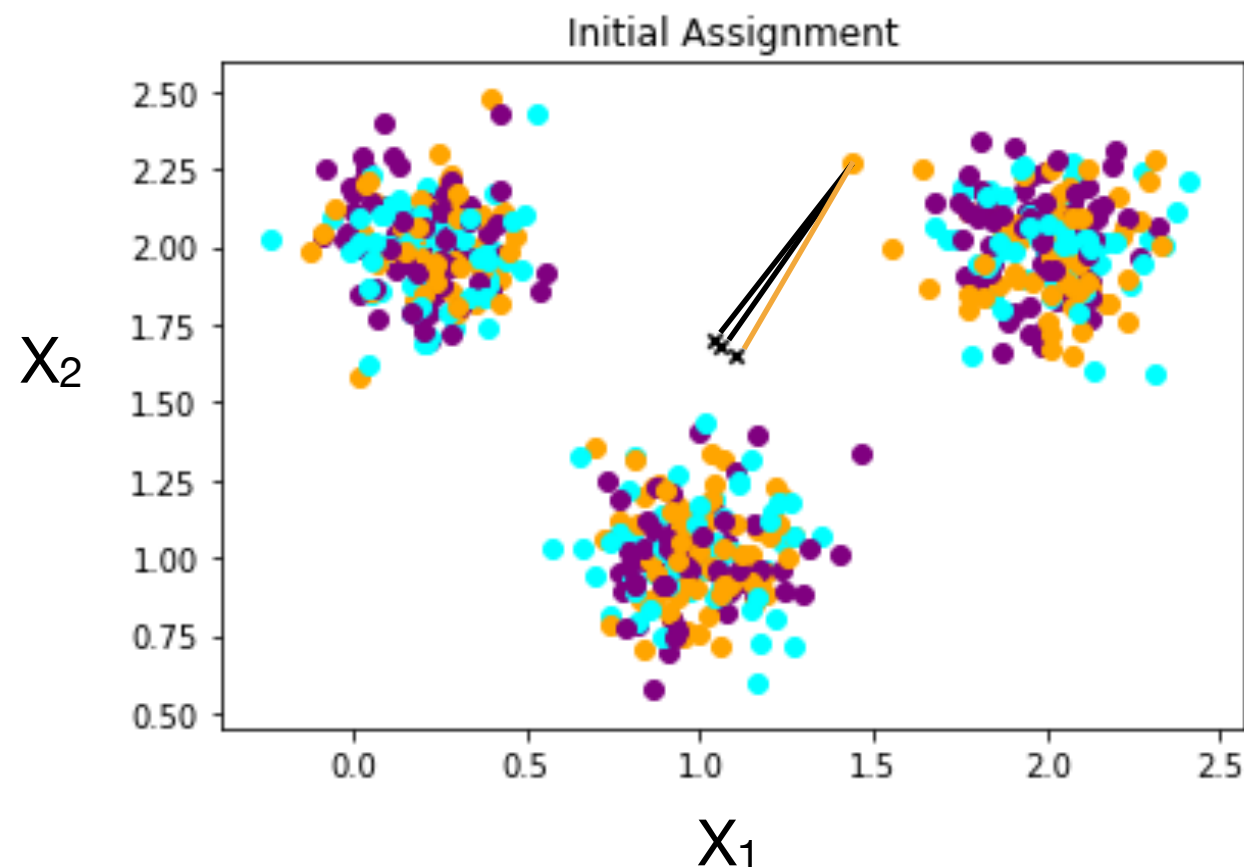


The k-means Algorithm

The algorithm:

```
> Initialize random cluster indices  
> Repeat until convergence:  
>   For each data point:  
>     Distance( point, centerj )  
>     new_index ← minj(distance)  
>     total_distance += distance
```

The data:

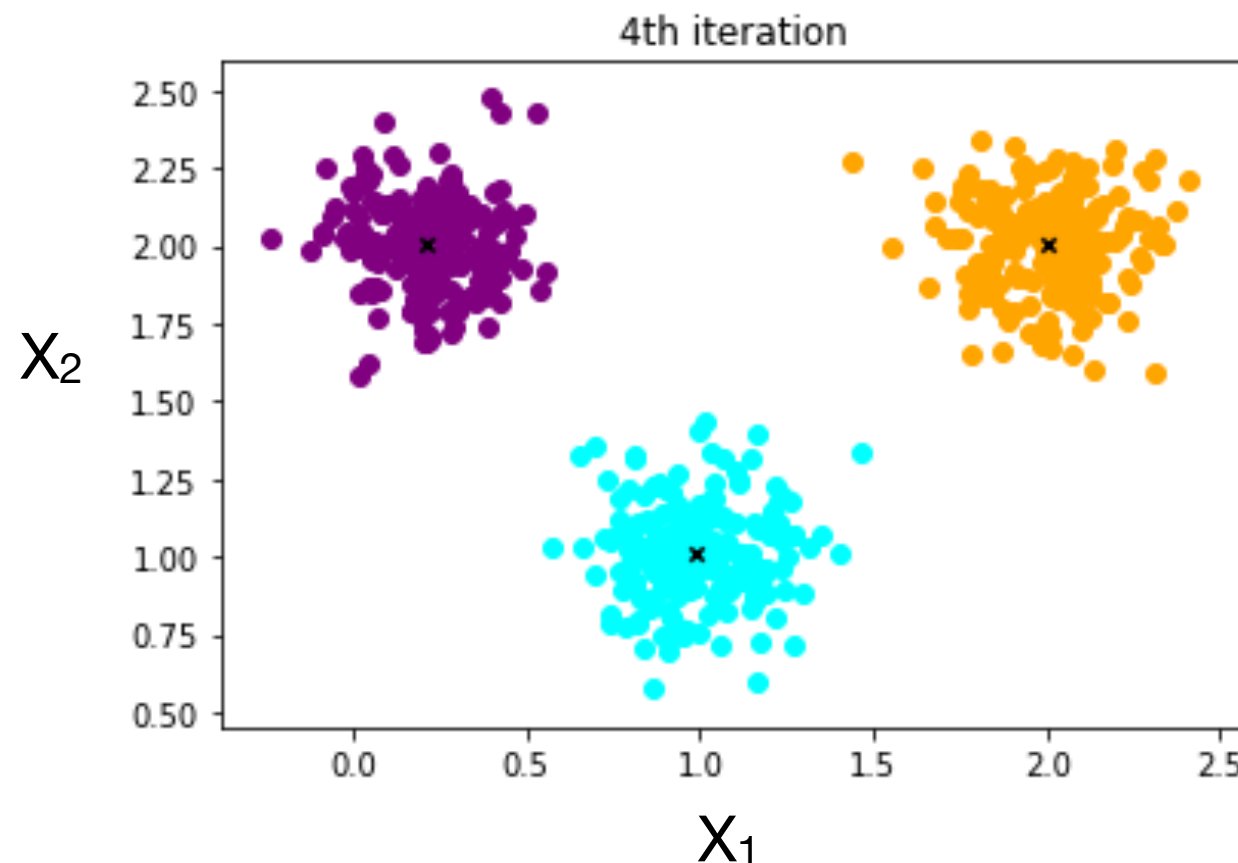


The k-means Algorithm

The algorithm:

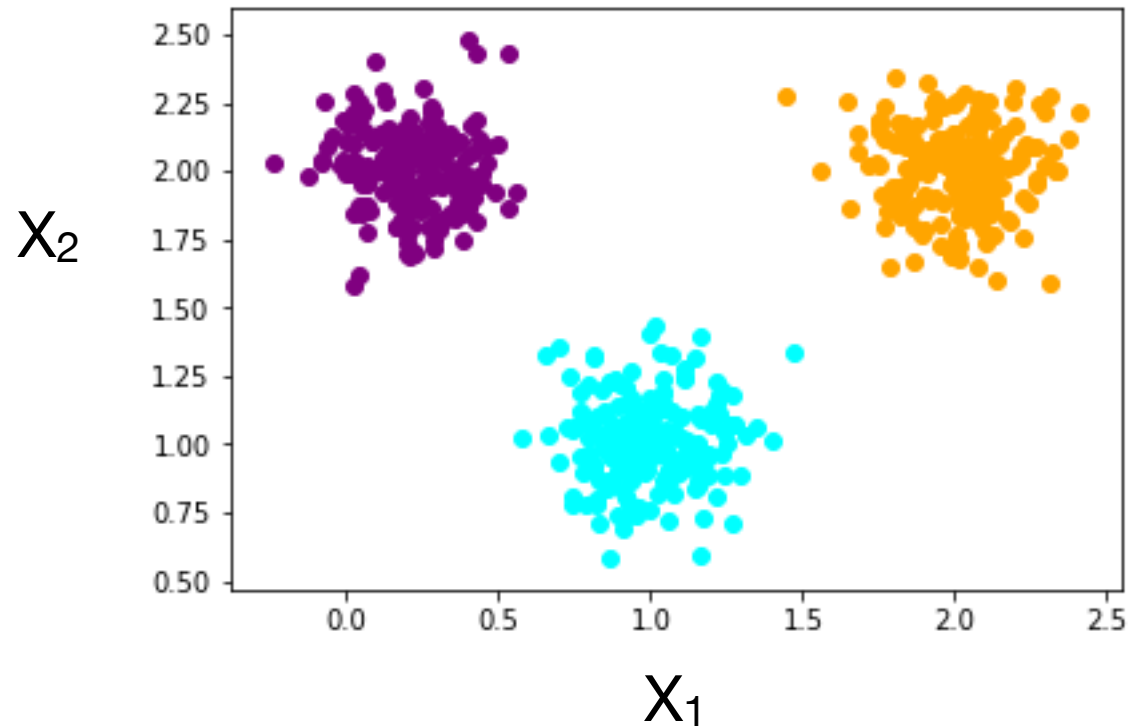
```
> Initialize random cluster indices
> Repeat until convergence:
>   For each data point:
>     Distance( point, centerj )
>     new_index ← minj(distance)
>     total_distance += distance
```

The data:

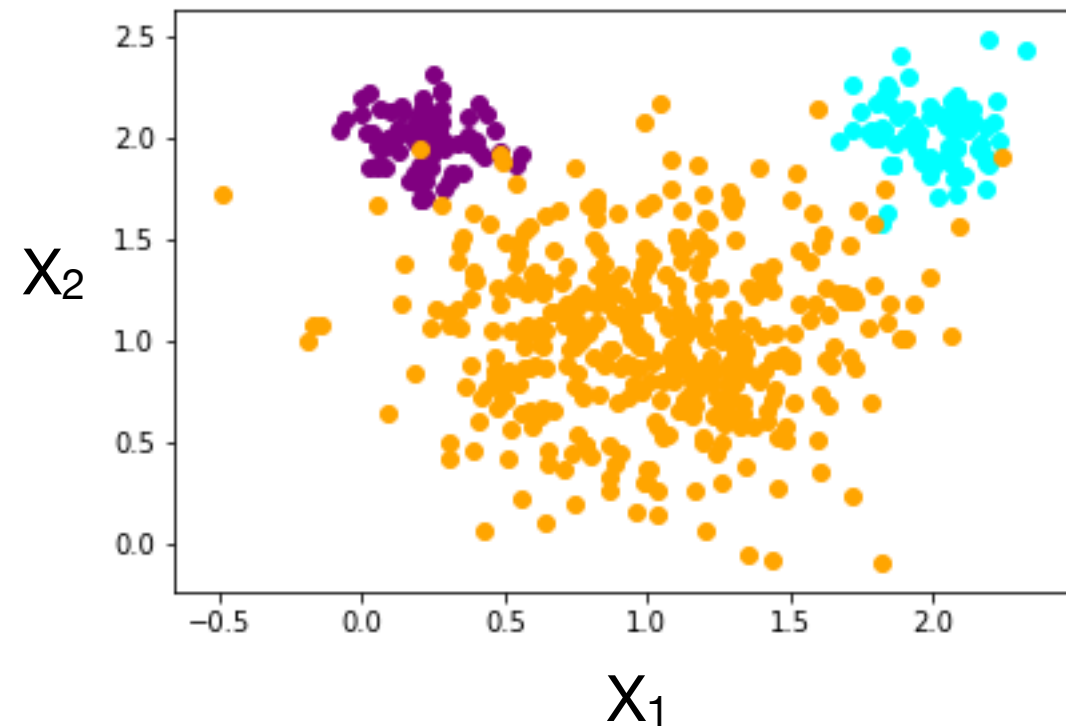


Problems with k-means

Previous data:



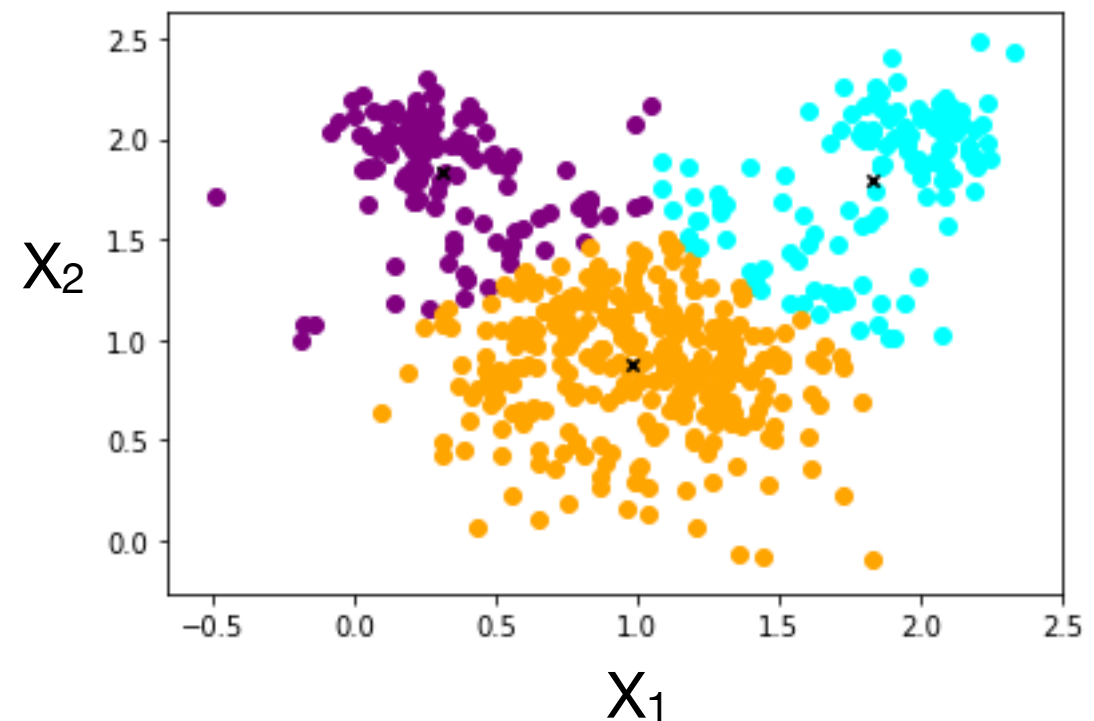
New data:



Poor fits to dissimilar clusters

Solution:
Soft k-means

4th iteration



Soft k-means Algorithm

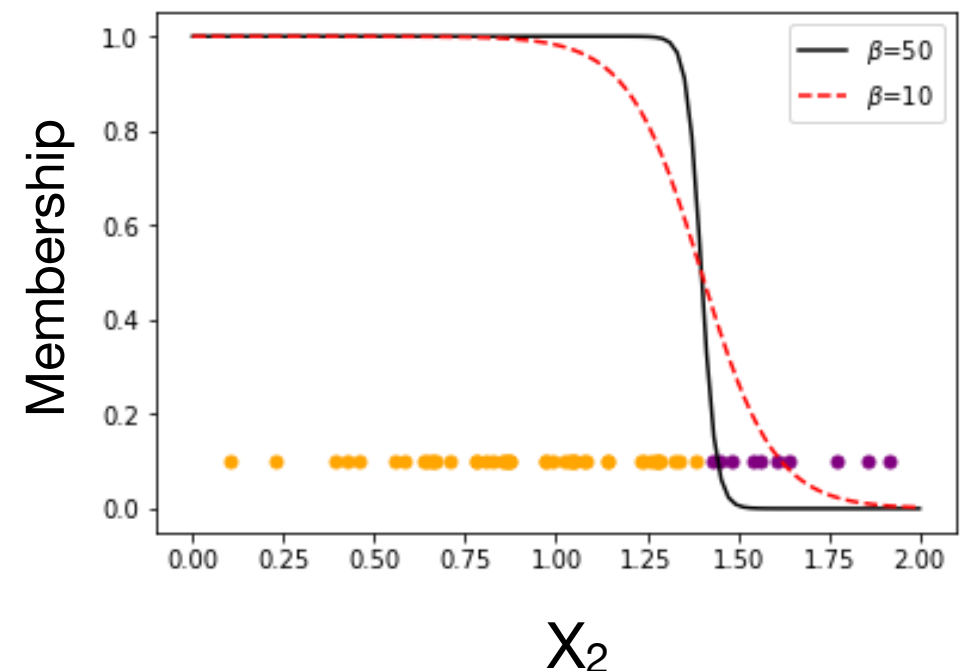
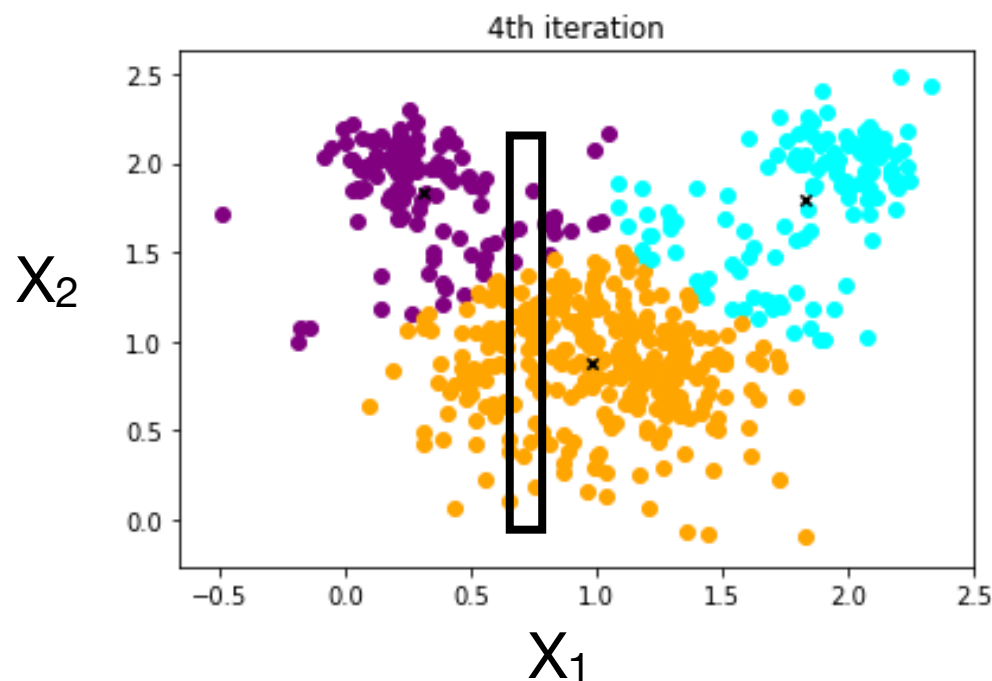
Hard k-means:

```
> new_index ← minj(distance)
```

Soft k-means:

```
> membership ← rj(distance)
```

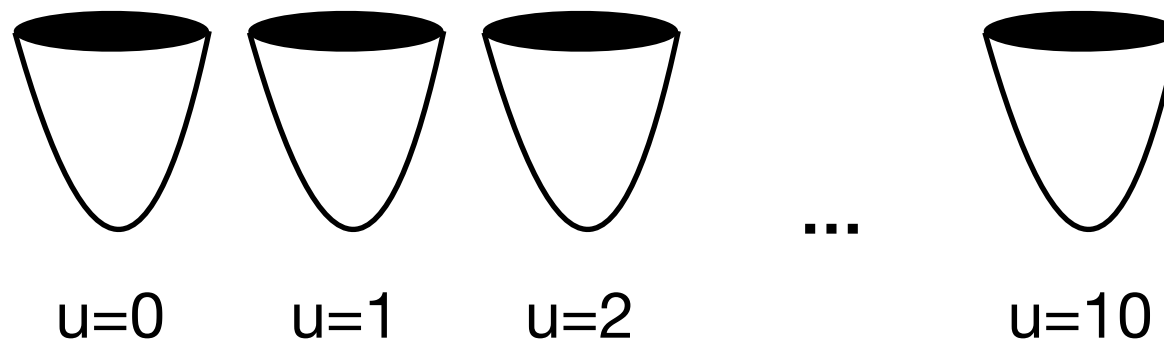
$$r_i = \frac{\exp[-\beta \text{dist}(x_i, c_j)]}{Z}$$



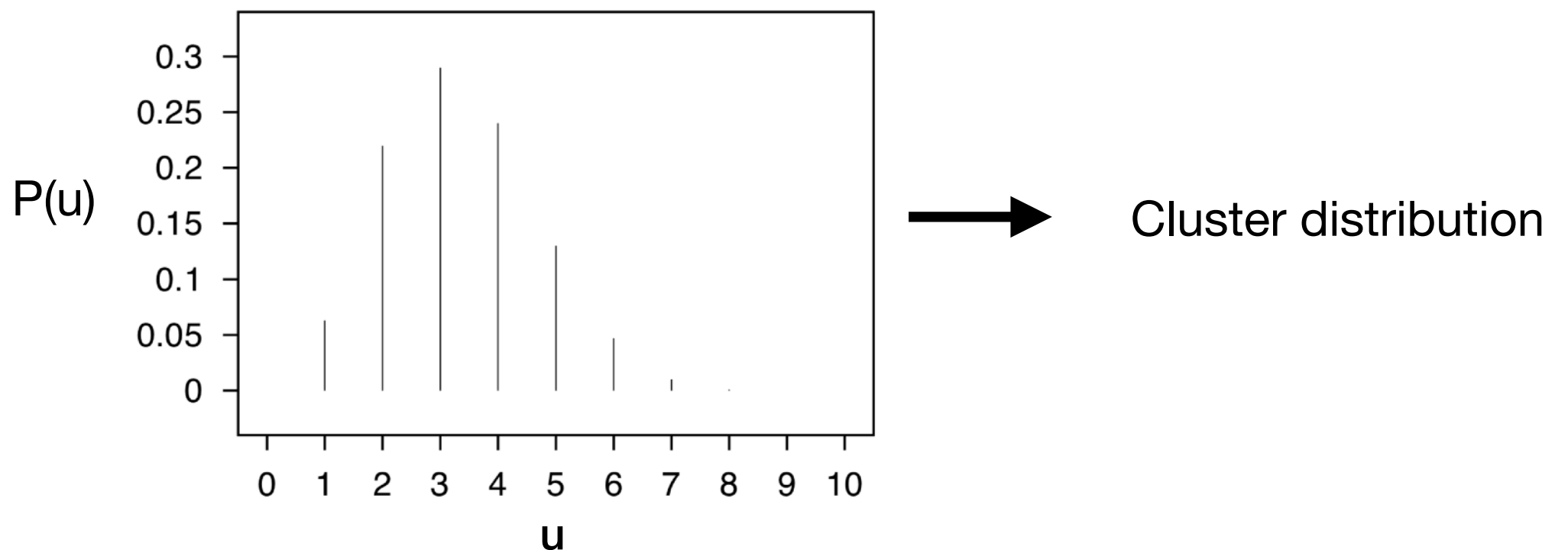
Soft k-means Algorithm

Utilizes full posterior distribution rather than approximate solution

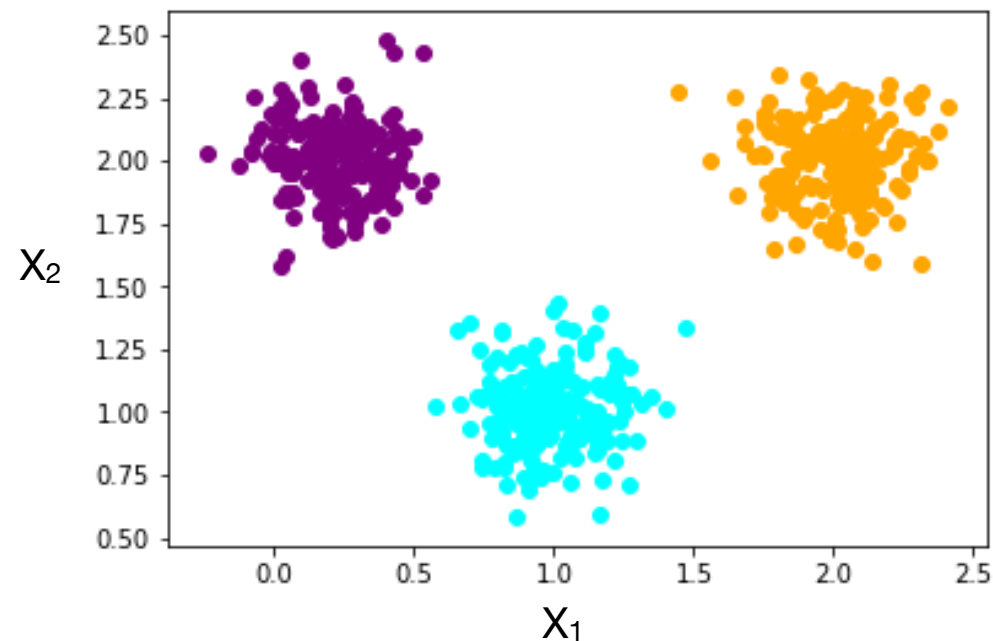
Recall the urn problem:



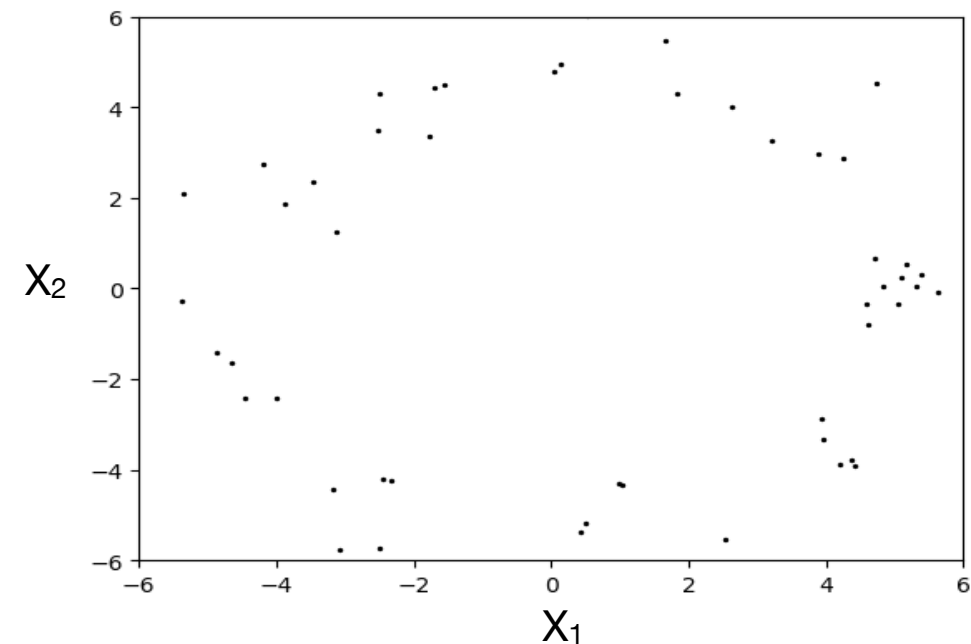
N=10 draws ● ● ● ○ ○ ... ○ → Cluster = urn 3



Another Problem: Choosing k



k is “obvious”



What is k?

Few clusters

Many clusters

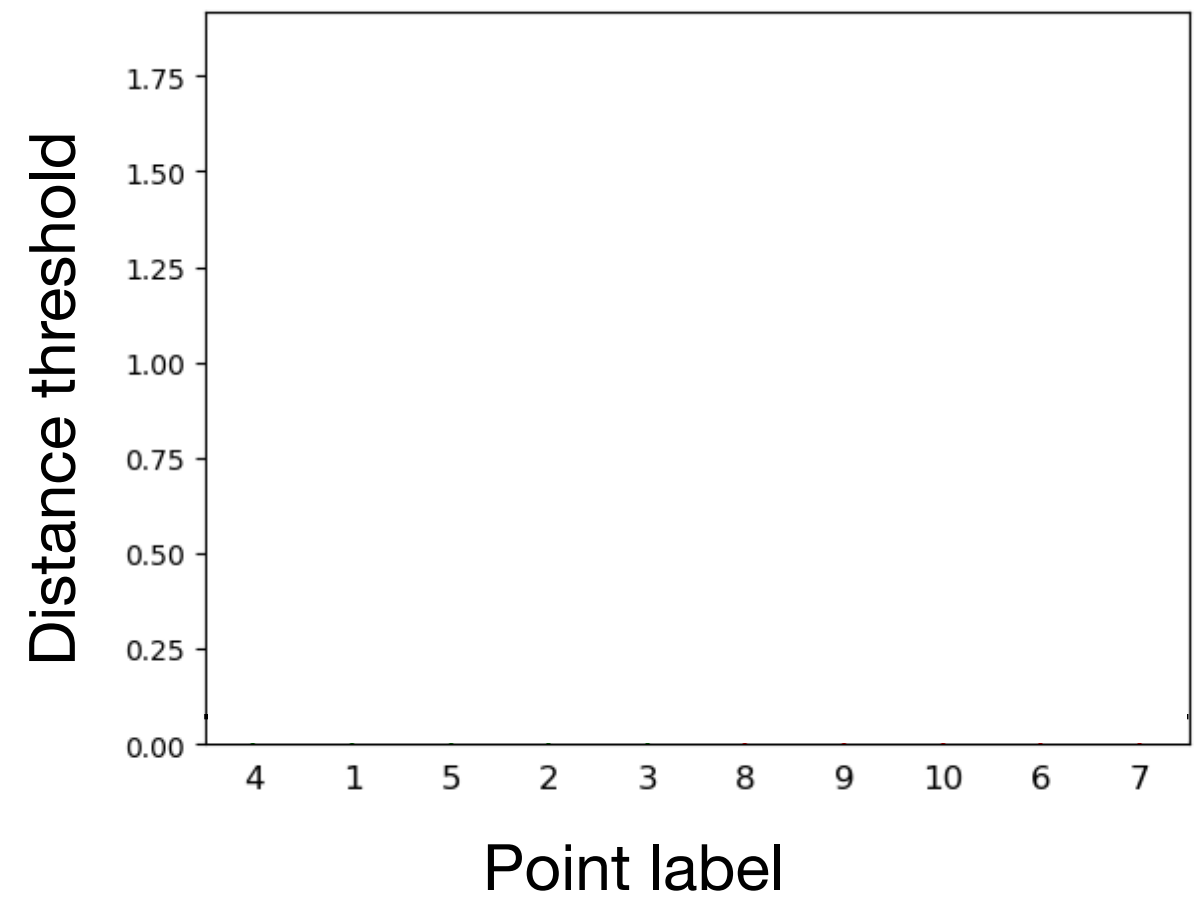
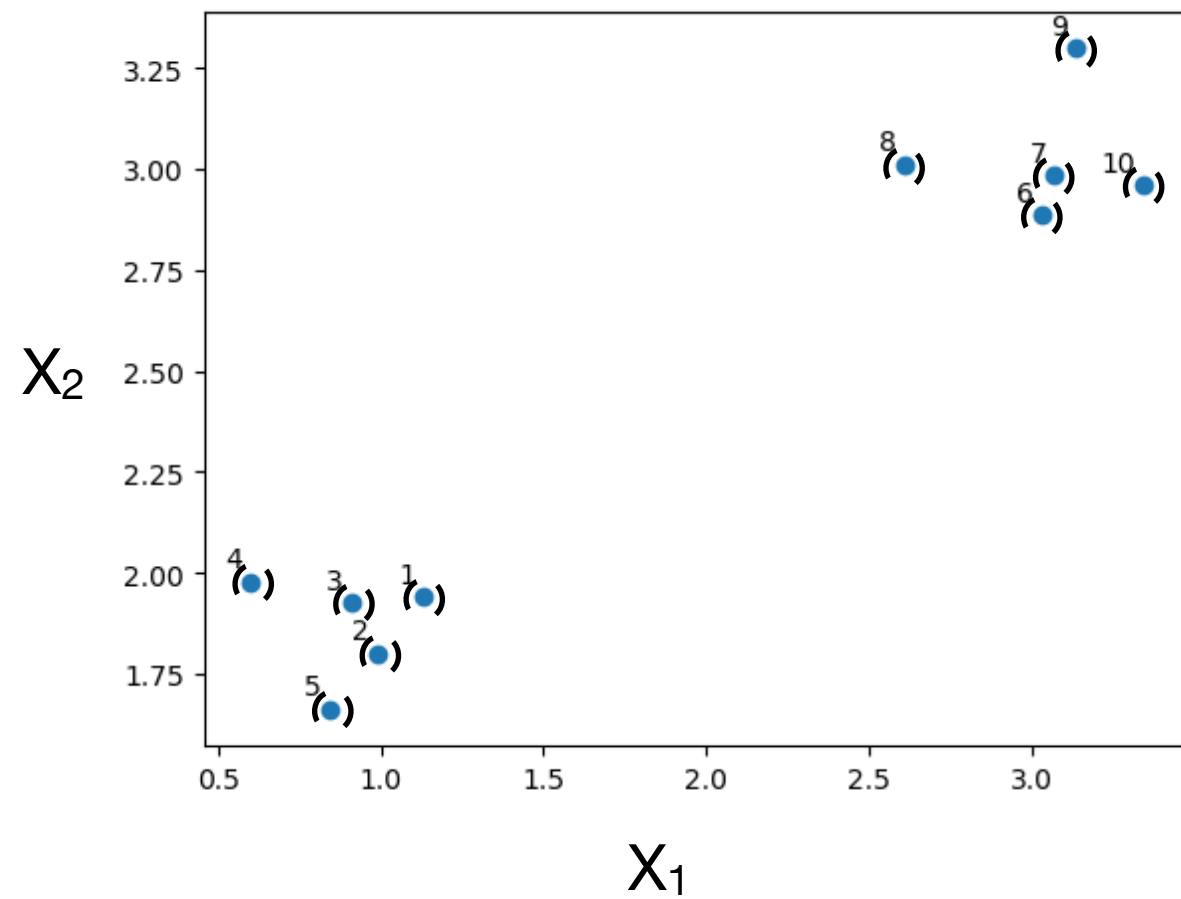


Compression

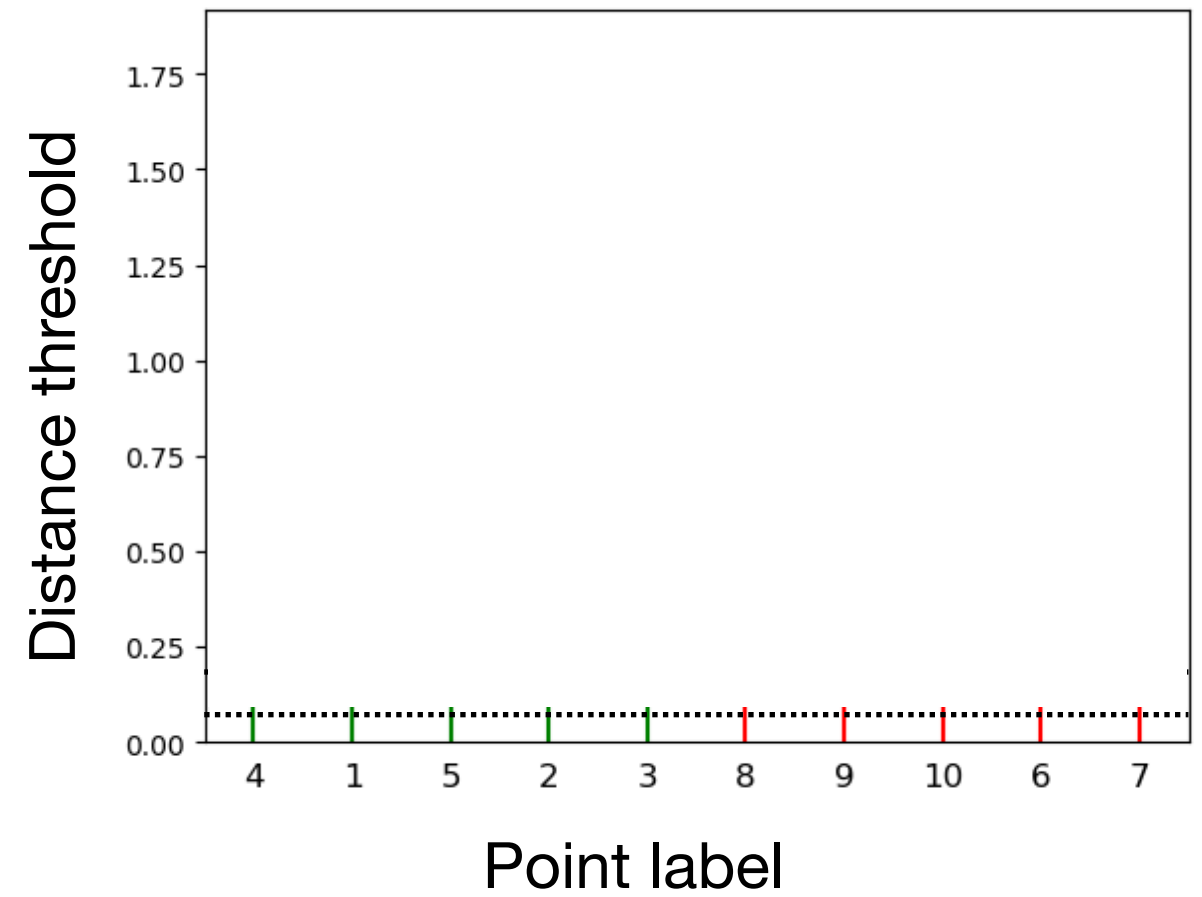
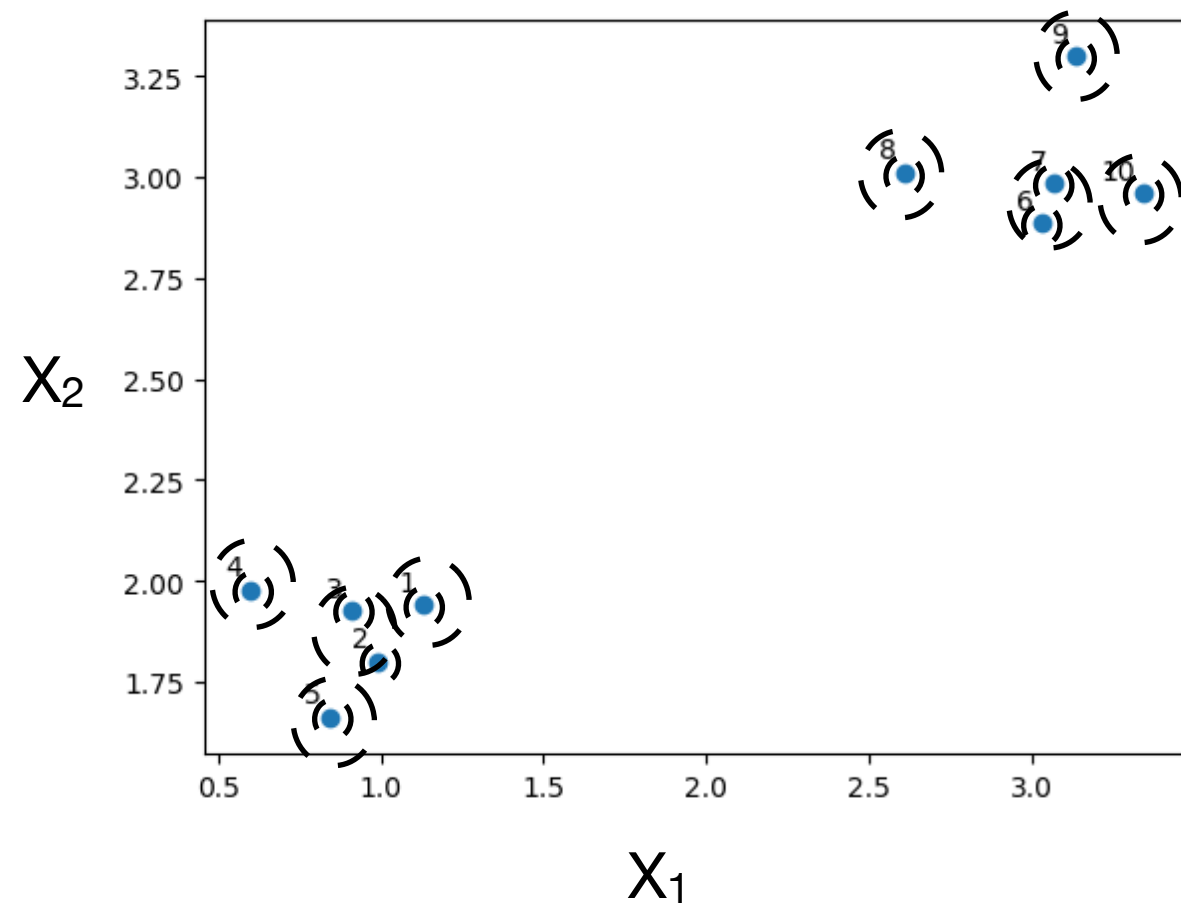
Accuracy

Recall hypothesis comparison with many parameters

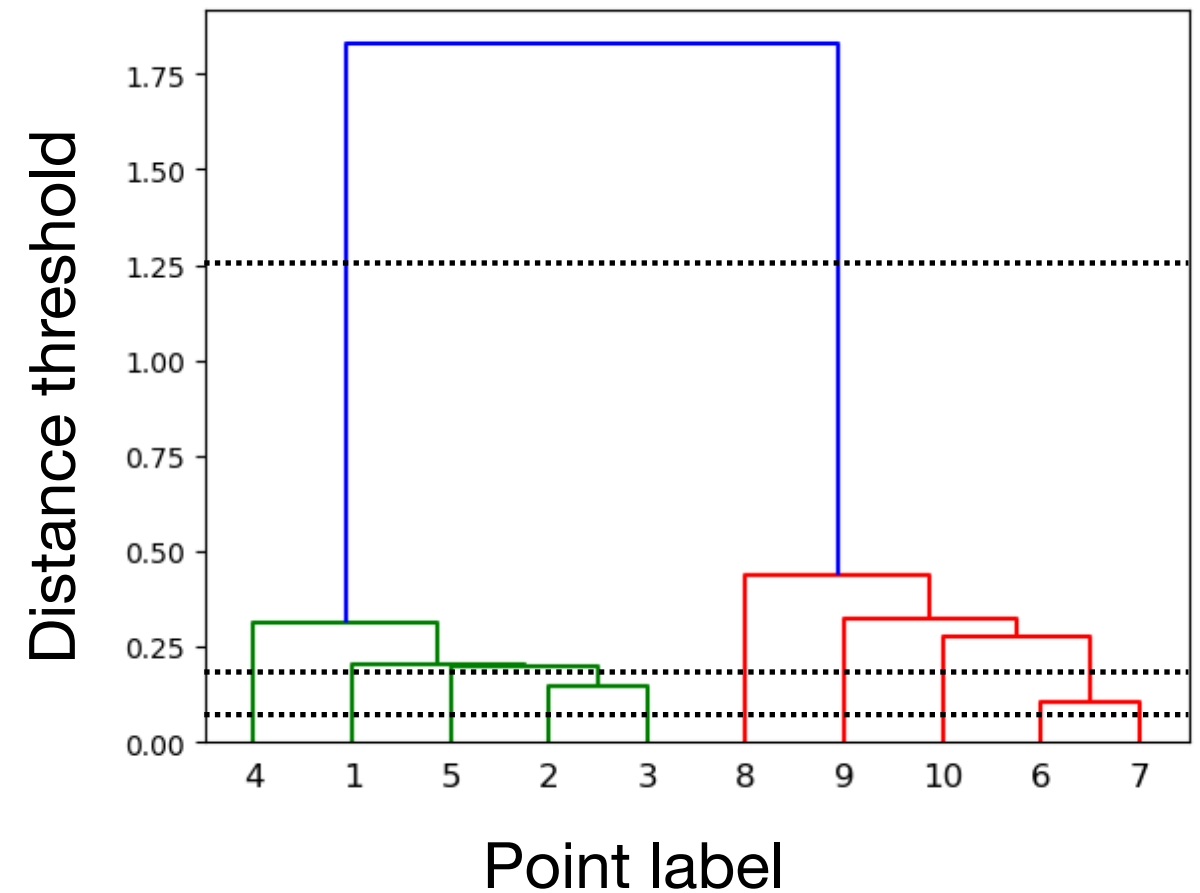
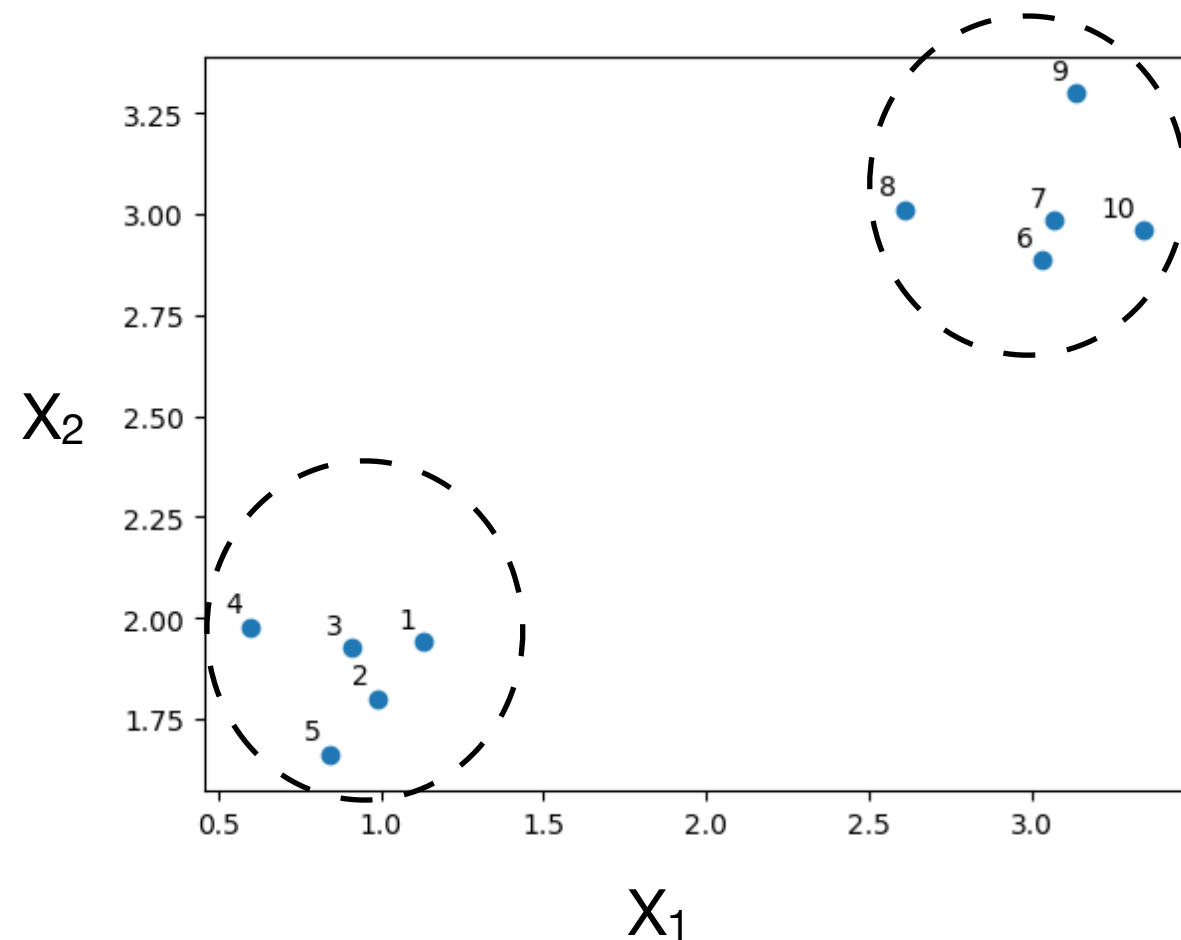
Solution: Dendrogram



Solution: Dendrogram



Solution: Dendrogram



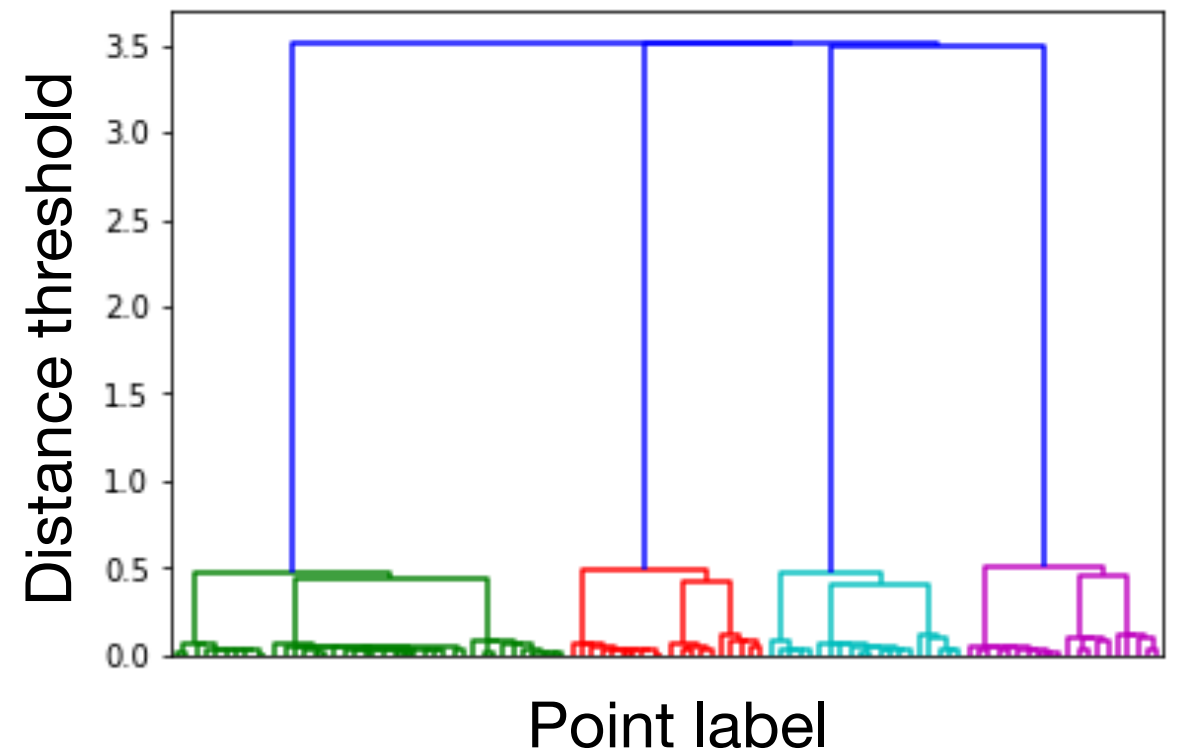
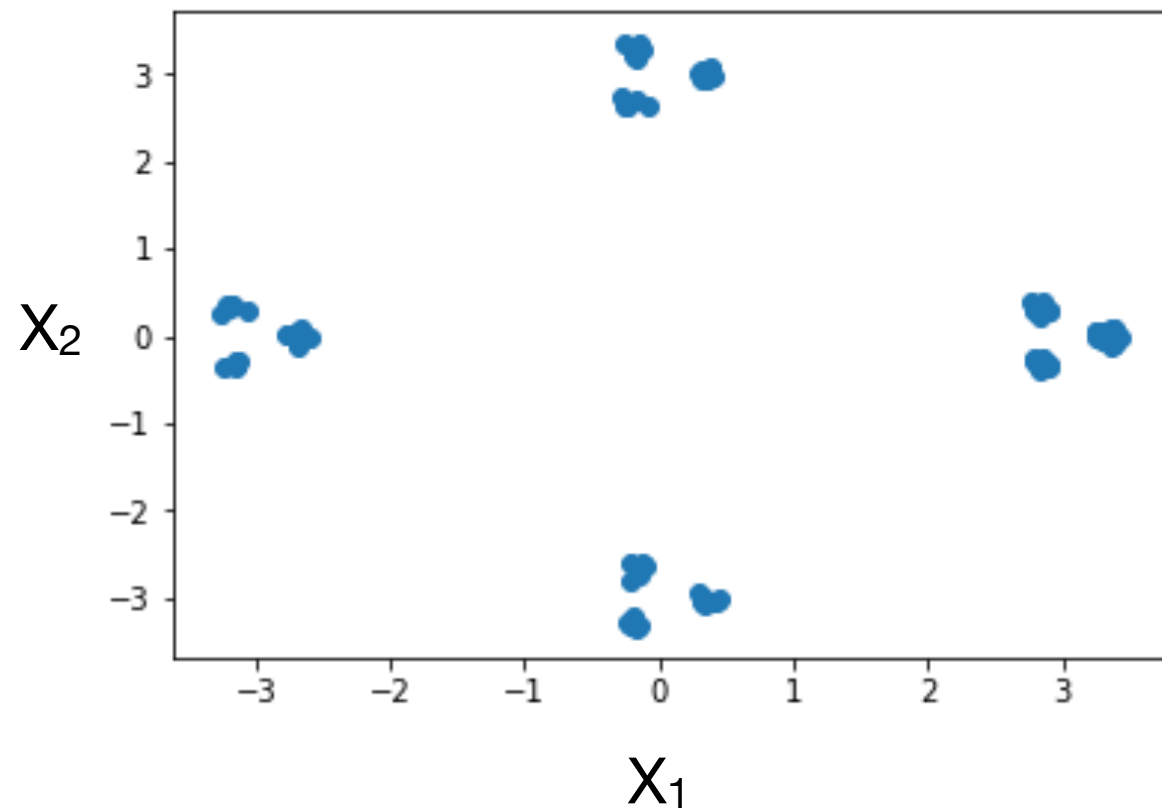
Definite criteria for k:

Choose maximum change in distance threshold

Hypothesis comparison:

Allow the data to speak for itself

No Silver Bullet



Which threshold is “correct?”

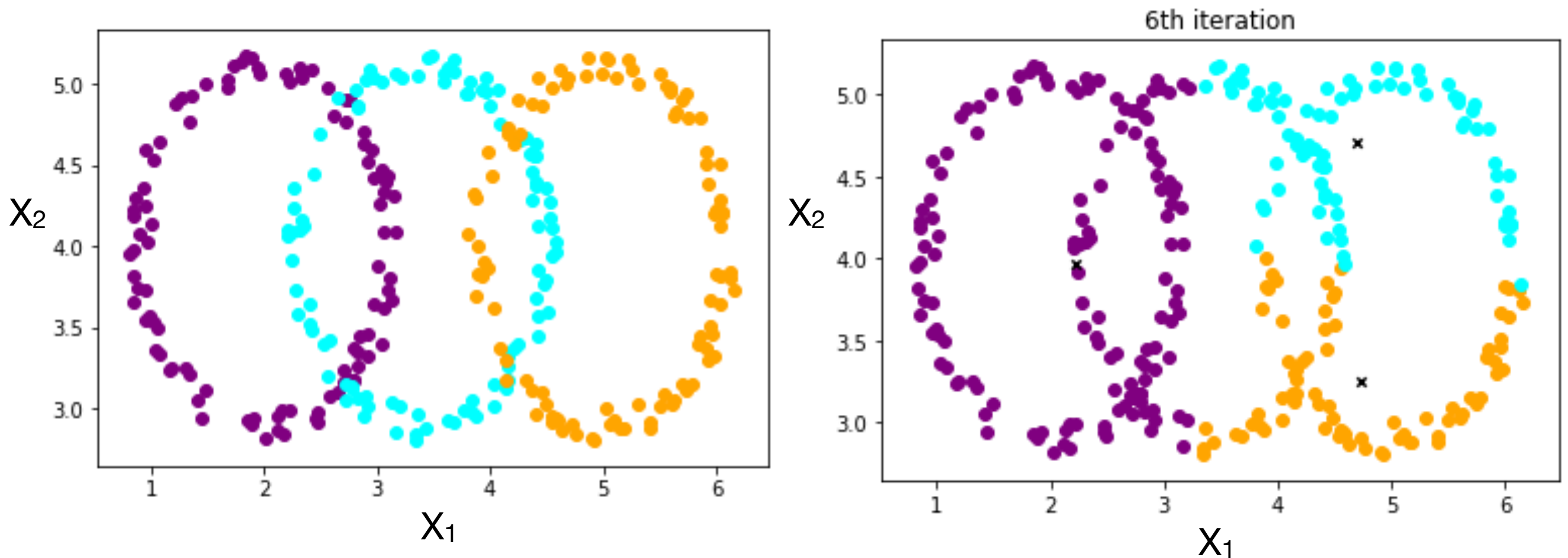
$k=4$?

$k=12$?

The question is not “which k should I choose?”
but rather “how is my data structured?”

No Silver Bullet Continued

Some data is very difficult to cluster



Clustering is

Takeaways

Clustering is a valuable tool for dimensional reduction.

Modified k-means is a good tool for clustering.

Clustering strategy is ultimately sensitive to data's structure.