

# Objects, Relationships, and Clustering

Scott C-H Pegg, Ph.D.

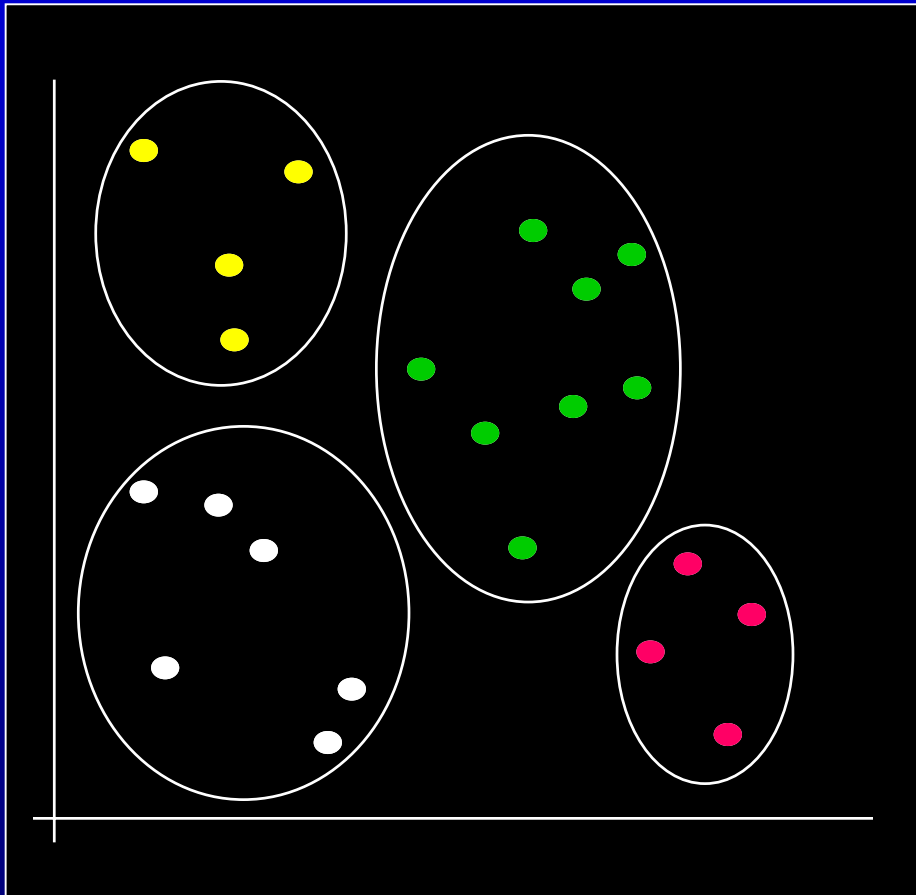
BMI203 April 20, 2004

# Overview

- What clustering means
- Relationships between objects
- Partitioning algorithms
- Hierarchical algorithms
- Choosing an algorithm
- Evaluating a clustering
- Clustering in Bioinformatics
- Homework instructions

# What is clustering?

Clustering is the grouping of objects together



A meaningful clustering groups objects such that there's a relationship between objects in a group.

# What's a Relationship?

## Binary relations

Relationships between two objects

Binary relations can have all sorts of properties.  
Here are three simple ones:

Given objects  $a_0, a_1, \dots, a_n$  in set  $A$ , and relation  $R$ ,

1.  $a_i R a_i$  for all  $a_i$  in  $A$

Reflexivity

2.  $a_i R a_j \implies a_j R a_i$

Symmetry

3.  $a_i R a_j$  and  $a_j R a_k \implies a_i R a_k$

Transitivity

# Equivalence Relations

A binary relation that satisfies all three of these properties is an **equivalence relation**.

What does this have to do with clustering?

An equivalence relation on a set of objects defines a partitioning of that set (ie. a clustering).

Any partitioning (clustering) of a set of objects defines an equivalence relation on that set.

# Equivalence Relations

You don't actually need an equivalence relation in order to perform a clustering.

So why am I showing you this?

- An introduction to Relational Algebra. This is what Relational Database Theory is based upon.
- Clustering is really about relationships between objects.
- It can provide a method of evaluating your clustering (we'll come back to this).

# Equality

This one is simple. Just substitute the  $R$  for  $=$ .

Is this an equivalence relation?

Yes

1.  $a_i = a_i$  for all  $a_i$  in  $\mathbf{A}$

2.  $a_i = a_j \implies a_j = a_i$

3.  $a_i = a_j$  and  $a_j = a_k \implies a_i = a_k$

Reflexive

Symmetric

Transitive

# Sequence Similarity

Let my objects in  $\mathbf{A}$  be a set of protein sequences.

Let  $R$  be “is similar to” (using a threshold score from pairwise sequence alignment).

Is this an equivalence relation?

1.  $a_i R a_i$  for all  $a_i$  in  $\mathbf{A}$  **Reflexive**
2.  $a_i R a_j \implies a_j R a_i$  **Symmetric**
3.  $a_i R a_j$  and  $a_j R a_k \not\implies a_i R a_k$  **NOT Transitive!**

No

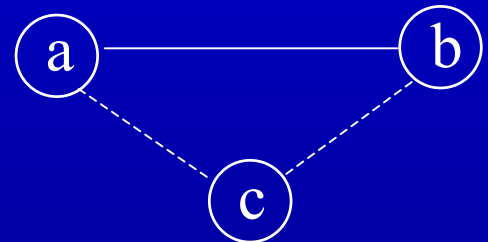


# Distance Relations

These relations (functions, really) consider the distance between two objects in a geometric sense.

Distance relations must satisfy the following properties

1.  $D(a,a) = 0$
2.  $D(a,b) \geq 0$
3.  $D(a,b) = D(b,a)$
4.  $D(a,b) \leq D(a,c) + D(c,b)$



triangle inequality

Is sequence identity a distance relation?

# Distance Relations

Here are a couple of examples:

## Euclidean distance

$$D(a,b) = \sqrt{(x_{a0} - x_{b0})^2 + (x_{a1} - x_{b1})^2 + \dots + (x_{an} - x_{bn})^2}$$

( $x_{a0}$  = the 0<sup>th</sup> property of object  $x_a$ )

## Manhattan (or city-block) distance

$$D(a,b) = |x_{a0} - x_{b0}| + |x_{a1} - x_{b1}| + \dots + |x_{an} - x_{bn}|$$

# Similarity Relations

There are lots of ways to define the similarity between objects.

A typical similarity function requires that

1.  $0 \leq S(a,b) \leq 1$
2.  $S(a,a) = 1$
3.  $S(a,b) = S(b,a)$

Note that we can define **dissimilarity** as

$$D(a,b) = 1 - S(a,b)$$

# Similarity Relations

Here are a couple of examples:

**% identity of strings**

$$S(a,b) = \frac{\# \text{ of symbols in common}}{\# \text{ of symbols in longer string}}$$

**Tanimoto similarity of bitstrings**

$$S(a,b) = \frac{\# \text{ of 1's in common}}{\# \text{ of 1's in both strings}}$$

# Missing Values

Note that all of the functions we've looked at require a complete set of values for every object.

But we may not have values for every variable of every object.

So what do we do about missing values?

# Missing Values

Here are a few options:

- Remove the objects with missing values from the overall clustering
- Replace the missing values of the variable with the average value of the variable over all the objects.
- Use prior knowledge of the variable's distribution and replace the missing value with the most likely value.

# Normalizing Values

Let's go back to the Euclidean distance between two objects

$$D(a,b) = \sqrt{(x_{a0} - x_{b0})^2 + (x_{a1} - x_{b1})^2 + \dots + (x_{an} - x_{bn})^2}$$

In some cases we're basing the distance on more than one variable. If these variables are not on the same scale (eg. Daltons vs. % identity), the variables will contribute unequally to the distance.

# Normalizing Values

We can alleviate this by normalizing the values of the variables.

For example, we can divide by the mean absolute deviation

Change each  $x_i$  to  $z_i$ , where

$$z_i = \frac{x_i - m}{S}$$

$$m = 1/n \{x_0 + x_1 + \dots + x_n\}$$

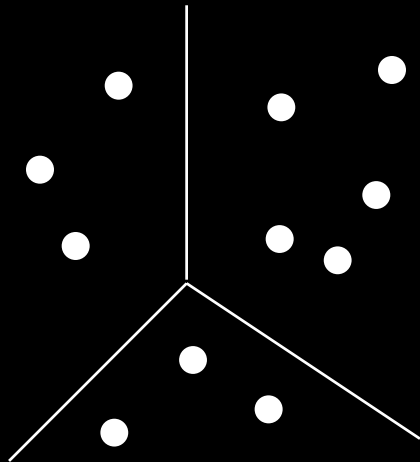
$$S = 1/n \{|x_0 - m| + |x_1 - m| + \dots + |x_n - m|\}$$



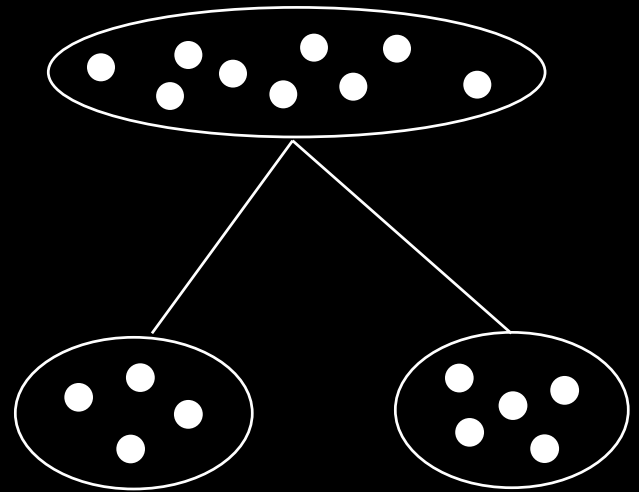
# Clustering Algorithms

Clustering algorithms come in two basic flavors

## Partitioning



## Hierarchical



# Partitioning Algorithms

Partitioning algorithms divide a set of objects into a **given** number of smaller sets.

You always get the number of clusters you asked for.

**But what if I don't know the number of clusters I want?**

Run the program several times over a range of cluster numbers.

# Partitioning Algorithms

The most common partitioning algorithms follow a 2-step process

1. Choose a set of “representative” objects.
2. Assign each remaining object to its nearest representative.

The critical part of these algorithms lies in the first step.

# The K-means Algorithm

One of the most widely used partitioning algorithms.

**Input:** A set of objects  $X = \{x_0, x_1, \dots, x_n\}$ ,  
and an integer  $K$ .

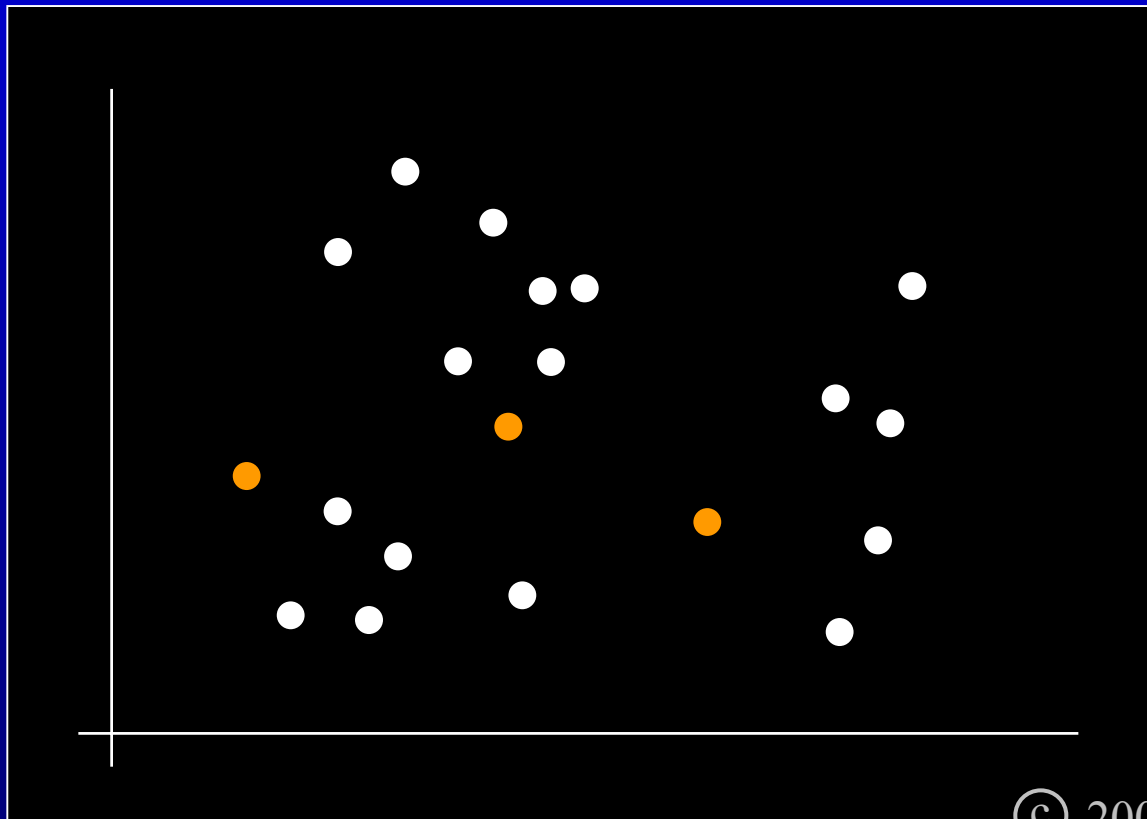
**Output:** A partition  $S$  of the objects that minimizes  
the sum of squared distance to the center of  
the cluster:

$$\sum_{j=1}^K \sum_{n \text{ in } S_j} D(x_n - \mu_n)^2$$

Where  $\mu_j$  is the  
mean of cluster  $S_j$ .

**Step 1:** Arbitrarily choose from the set of objects  
K initial cluster centers,

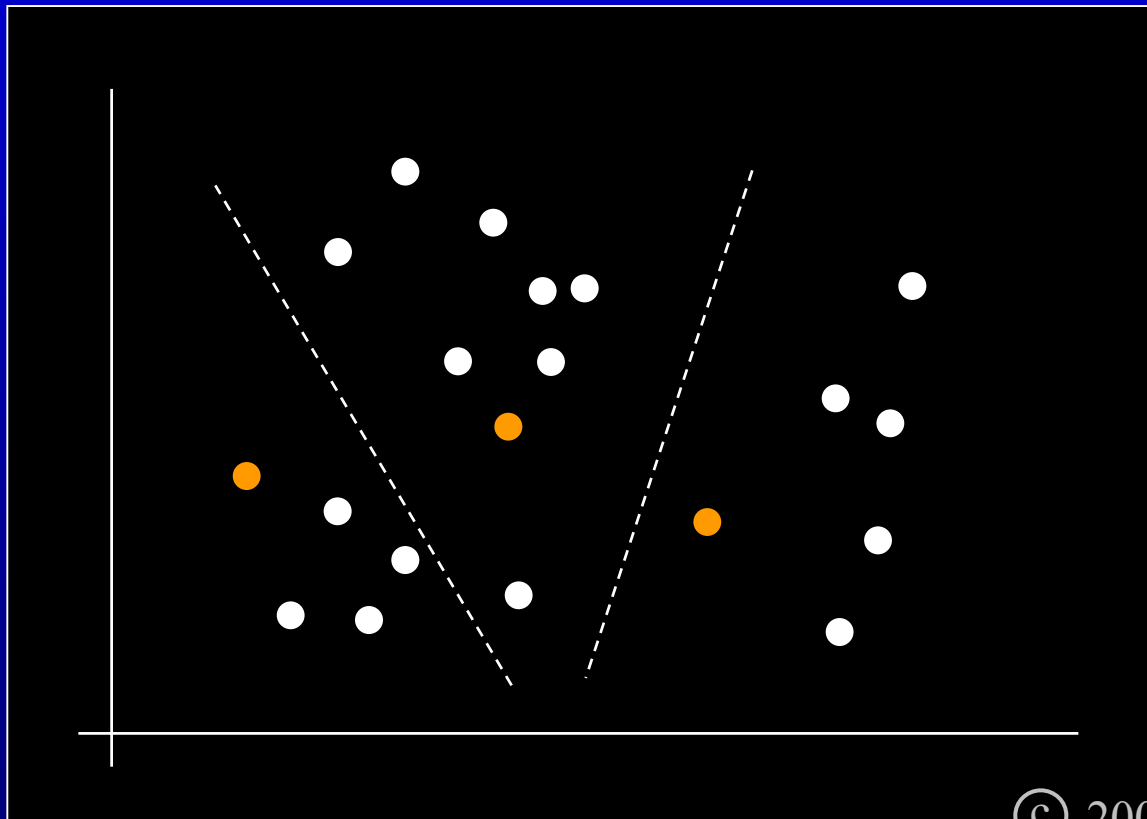
$$M_1^{(0)}, M_2^{(0)}, \dots, M_K^{(0)}$$



**Step 2:** Assign each of the objects in  $X$  to the cluster with the nearest cluster center.

$X_i$  in  $S_j$  if

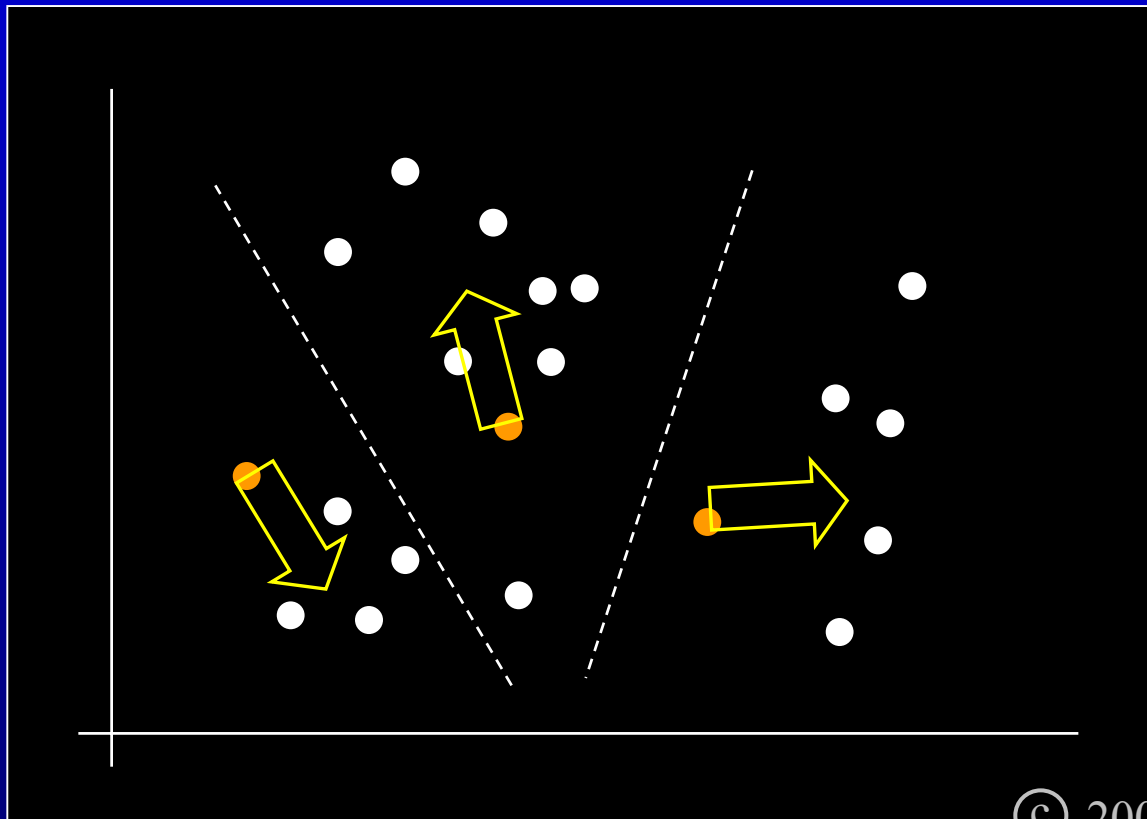
$$D(X_i, M_j^{(1)}) = \min \{ D(X_i, M_h^{(1)}, h = 1, \dots, k) \}$$



**Step 3:** Recalculate the cluster centers to get  $M_j^{(1+1)}$

$$M_j^{(1+1)} = 1/N_j^{(1)} \sum_{X \text{ in } S_j} X$$

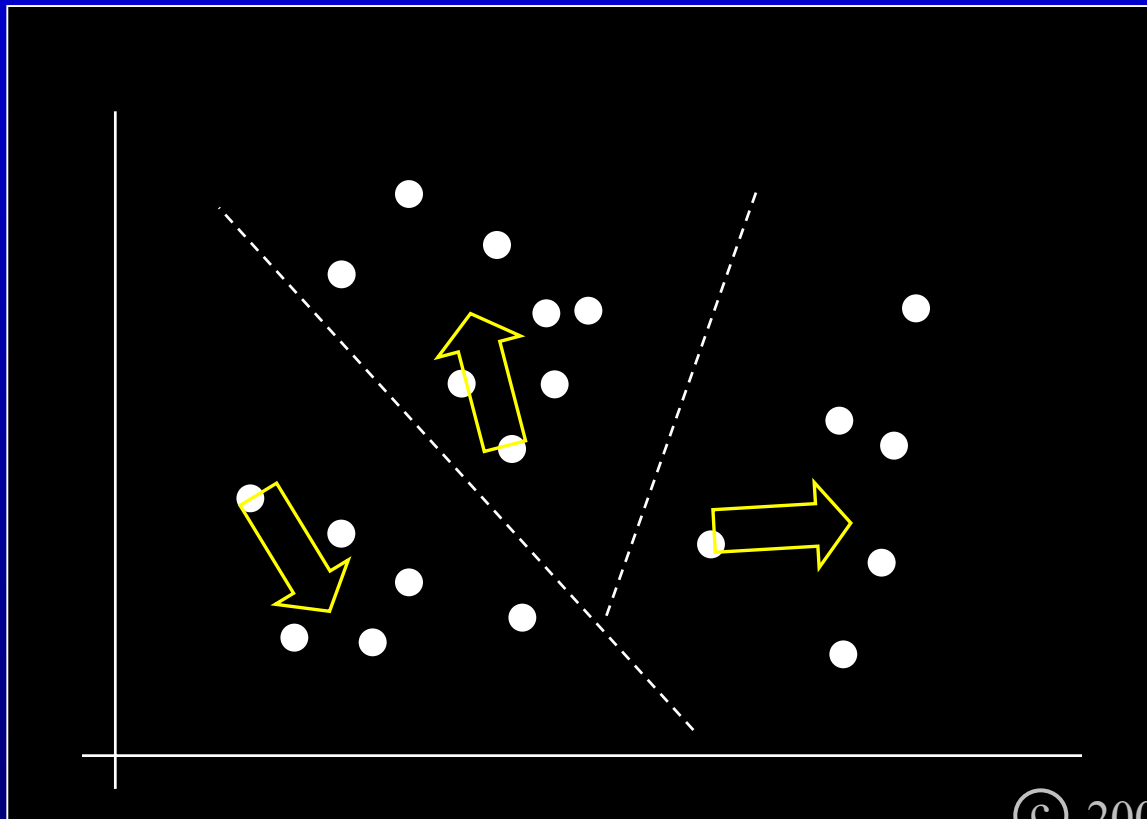
where  $N_j^{(1)}$  is the number of samples in cluster  $S_j^{(1)}$



**Step 4:** If the clusters have not changed (or a given number of iterations has been reached)

$$M_j^{(l+1)} = M_j^{(l)}$$

Terminate. Else, go to step 2.



What's the big O?

$$O(n^2)$$

Since you need to calculate the distance between all  $n$  objects



# K-means Variants

Recall that with K-means, we're minimizing the sum

$$\sum_{j=1}^K \sum_{n \text{ in } S_j} D(X_n - \mu_n)^2$$

where  $\mu_n$  is the mean value of the X's in  $S_j$

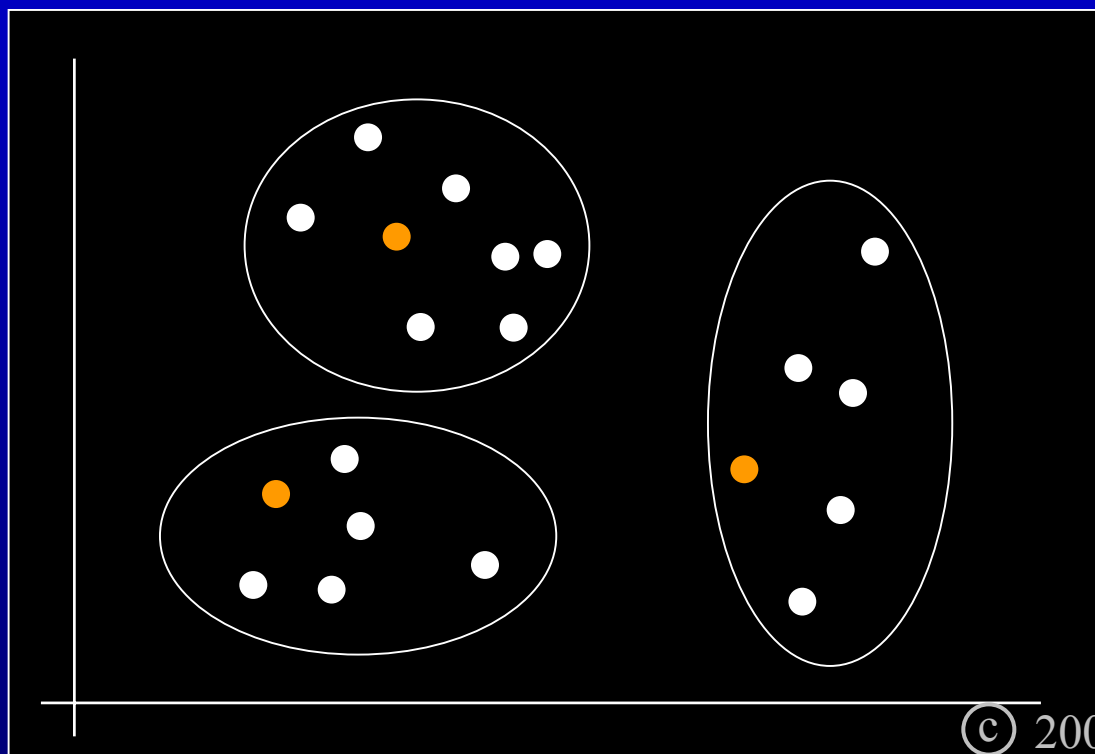
You can replace  $\mu_n$  with other types of values to make variants.

Examples: **K-medoids, K-centroids**

# Covering Algorithms

Covering algorithms are also based on choosing a representative set of objects, but don't require a pre-determined number of clusters.

Instead, they use a distance threshold for objects within a cluster.



# Partitioning lots of objects



With a run time of  $O(n^2)$ , partitioning methods can take a long time for a large number of objects.

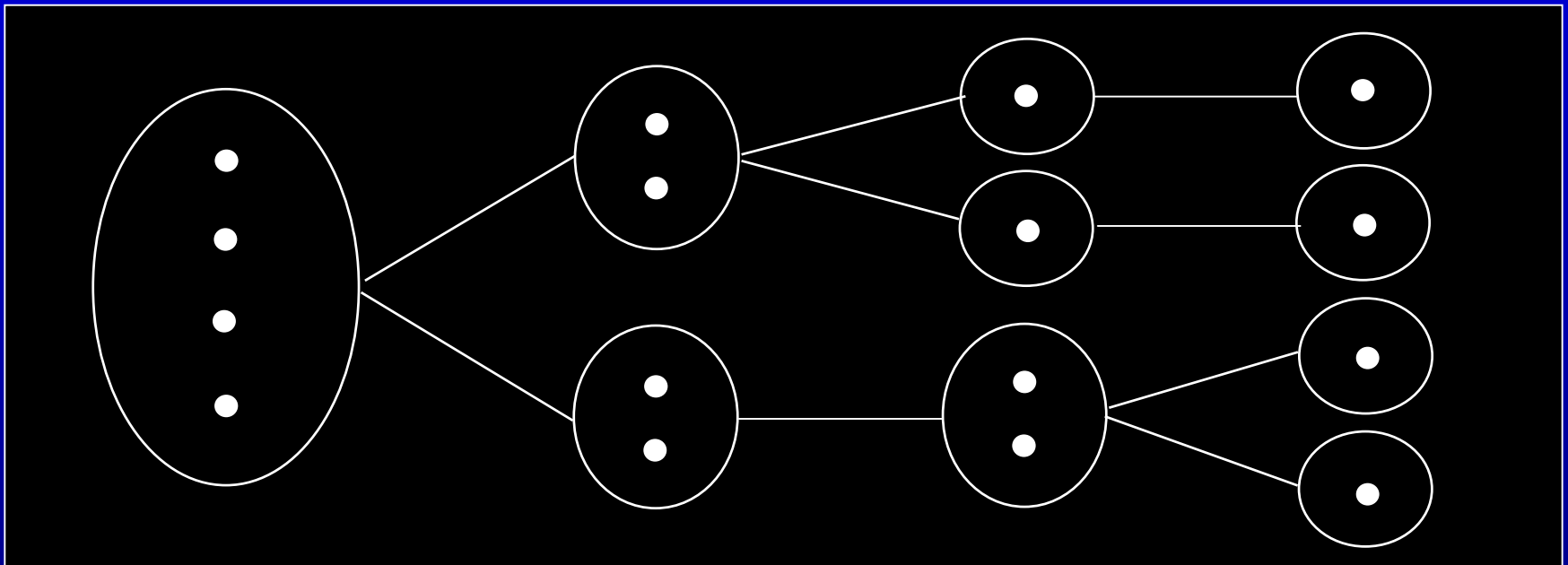
One common trick is to cluster only a small sample of the total set of objects, and then go back and assign the remaining objects to the representatives used in the sample clustering.

Since this depends on the original sample, you should do this several times with different samples.

# Hierarchical Algorithms

Hierarchical clustering algorithms come in two basic flavors

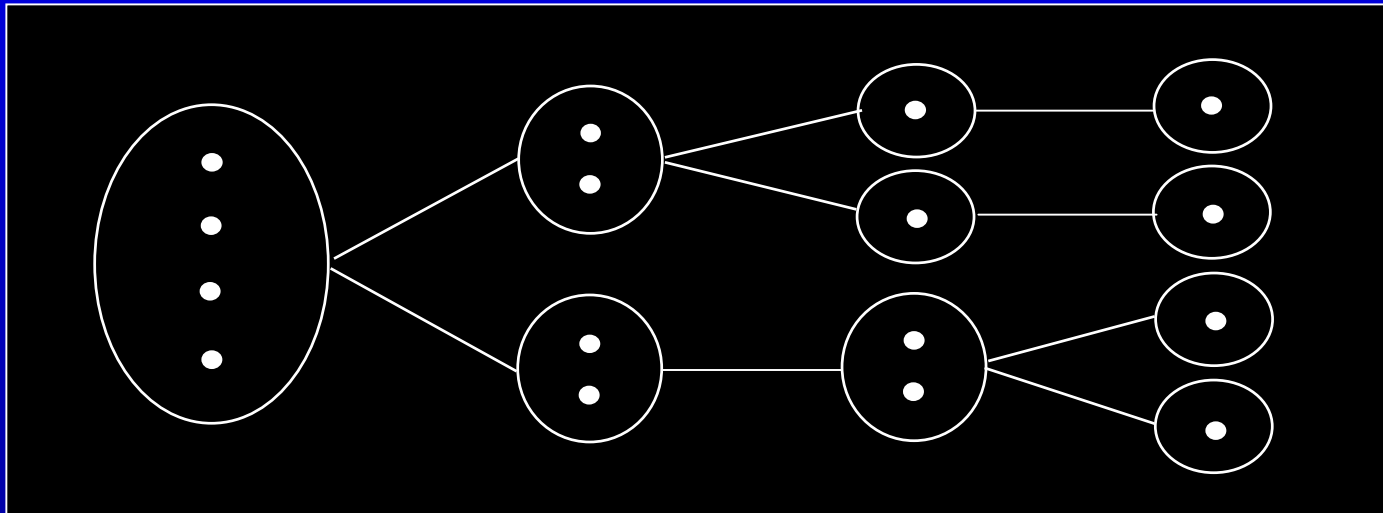
Divisive   Agglomerative



Note that the output here is not a single clustering, but a set of clusterings at different resolutions.

# Hierarchical Algorithms

The key step in hierarchical algorithms is the decision of which sets to join (or split).



One of the drawbacks is that once you've made this decision, you can't go back and split (or re-join) clusters.

There are lots of ways to make these decisions.

# Agglomerative Algorithms

Agglomerative algorithms work as follows:

Step1: Put each object in its own cluster.

Step2: Choose (via a **linkage criterion**), two clusters and merge them.

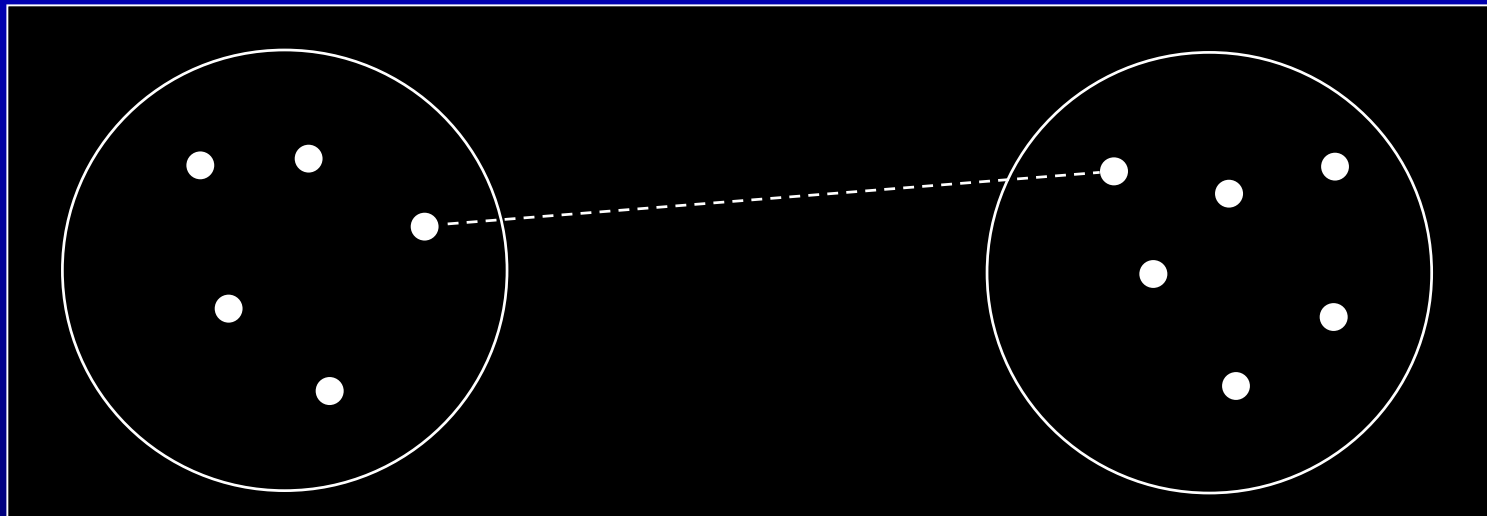
Step3: If there is only one cluster left, stop.  
Else, go to step 2.

# Single Linkage

One of the oldest and simplest methods.

aka “Nearest-Neighbor”

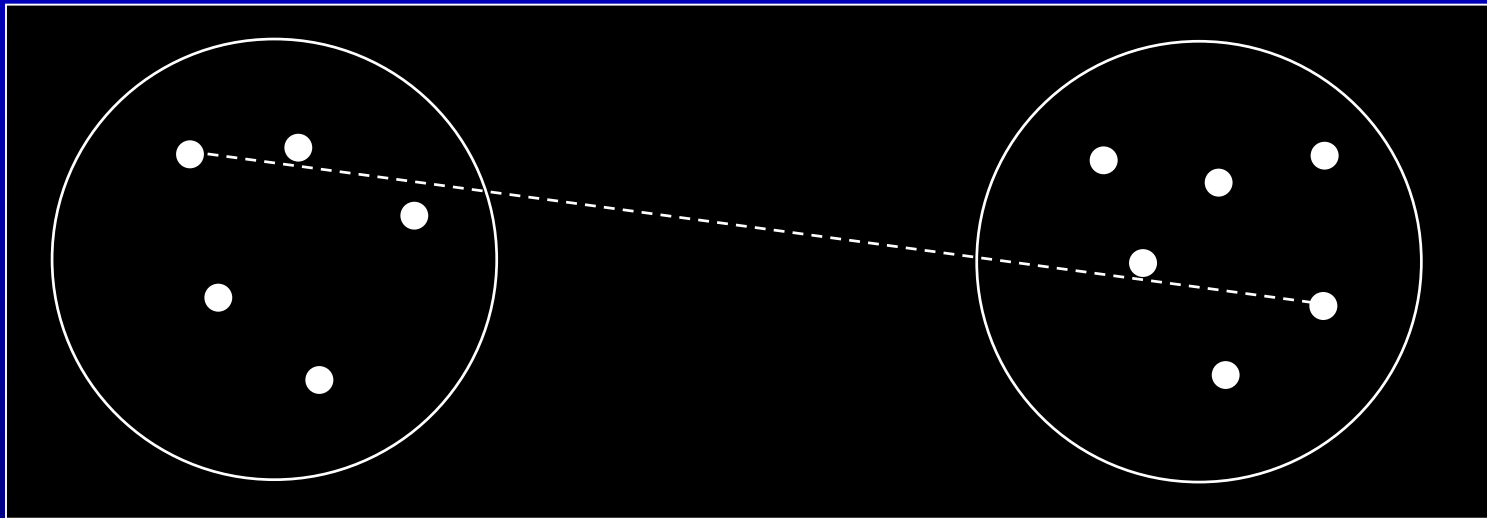
Choose the clusters with the shortest distance between the closest object in one cluster and the closest object in the other cluster.



# Complete Linkage

aka “Furthest-Neighbor”

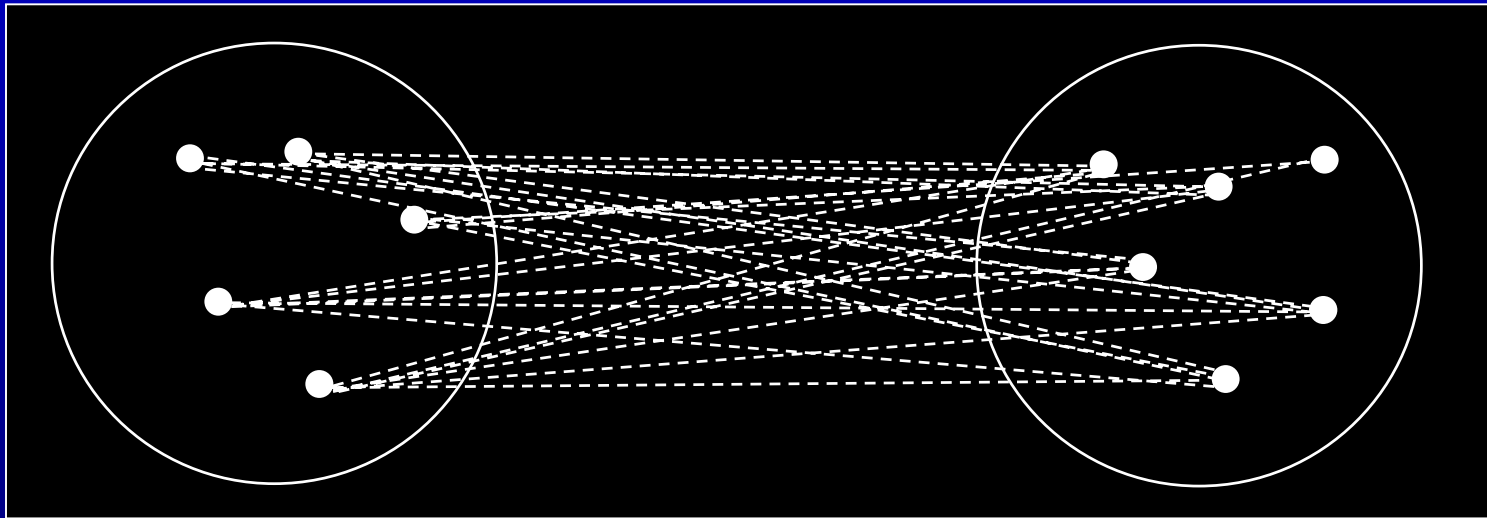
Choose the clusters with the shortest distance between the furthest object in one cluster and the furthest object in the other cluster.





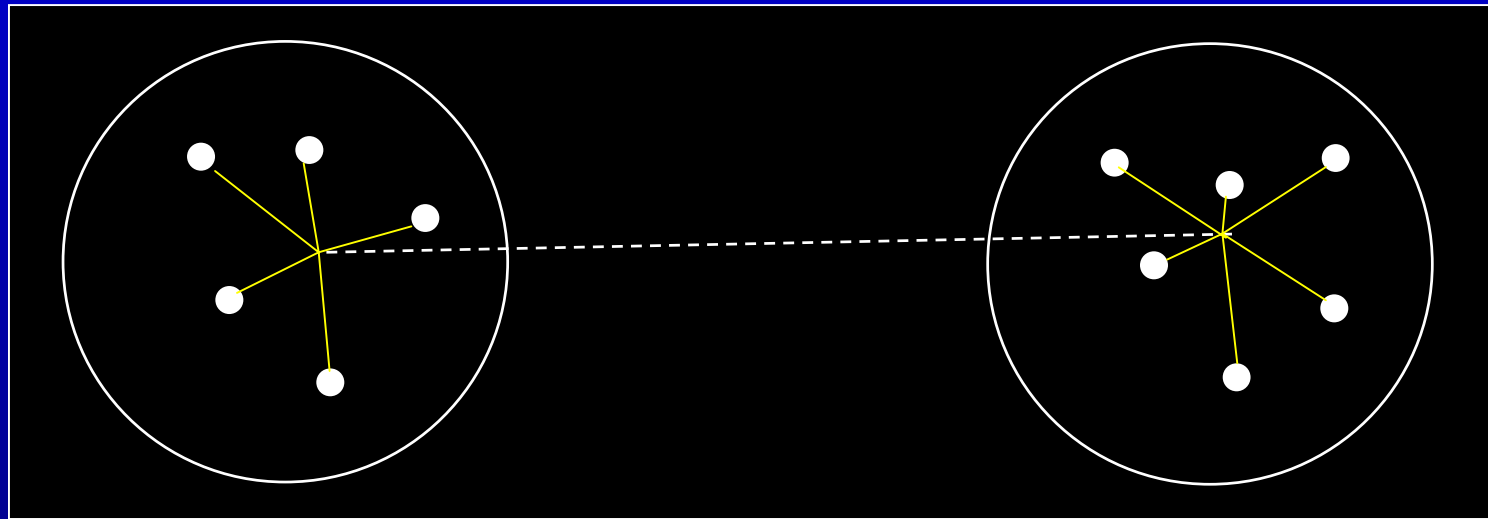
# Average Linkage

Choose the clusters with the shortest average distance between all objects in one cluster and all objects in the other cluster.



# Centroid Linkage

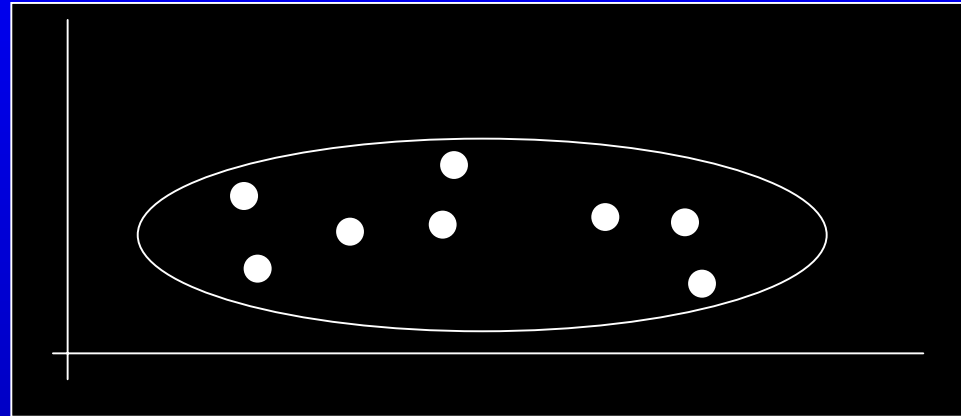
Choose the clusters with the shortest distance between the centroid of one cluster and the centroid of the other cluster.



# Some Comparisons

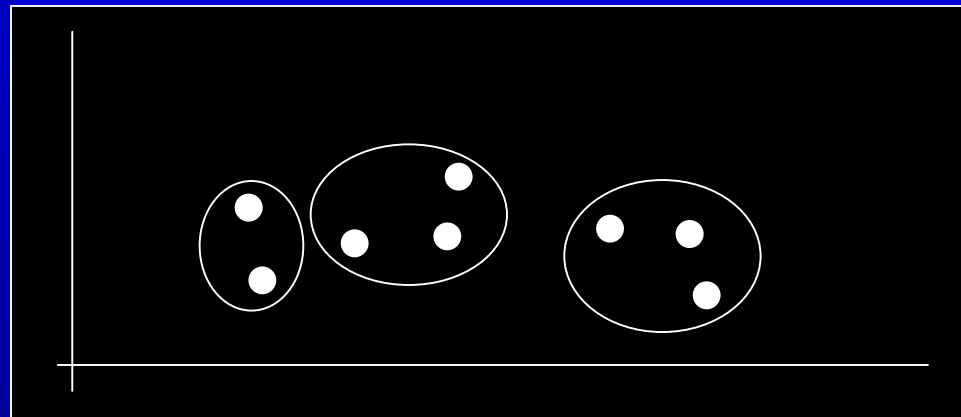
Single linkage

Makes long,  
drawn-out clusters



Complete linkage

Makes compact  
clusters



Centroid linkage

Somewhere in between

# Complexity

Almost all implementations of hierarchical clustering are agglomerative because of the theoretical run times.

In an agglomerative algorithm, the first step considers all of the possible fusions of two (single element) clusters,

$$\frac{n(n-1)}{2} = O(n^2)$$

# Complexity

A divisive algorithm must first consider all divisions of the entire data set into two non-empty sets,

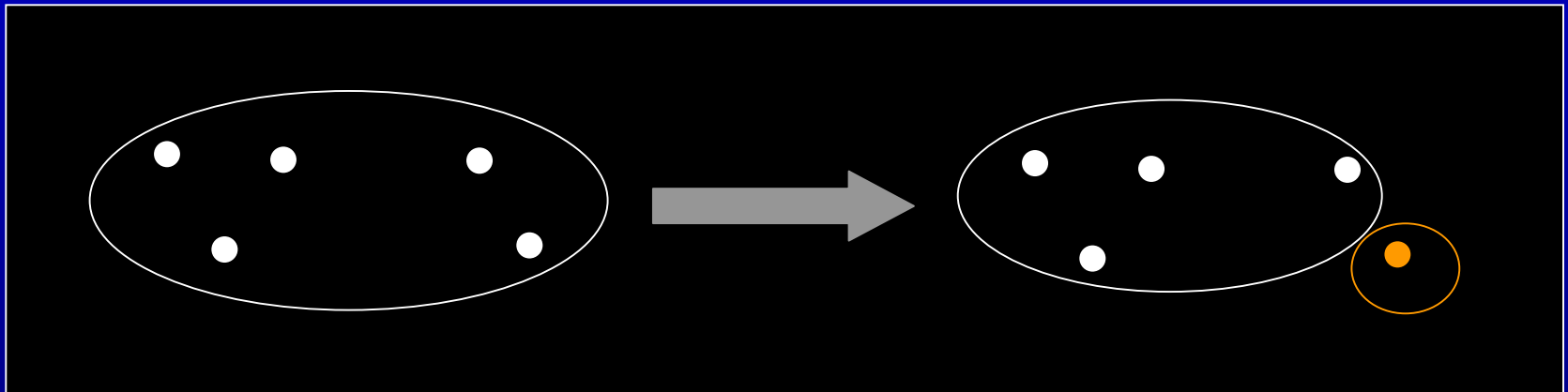
$$2^{n-1} - 1 = O(2^n)$$

There do exist some divisive algorithms that use tricks to get around the expensive first steps.

# The Macnaughton-Smith Method

**Step 1:** Find the object with the highest average distance from all of the others.

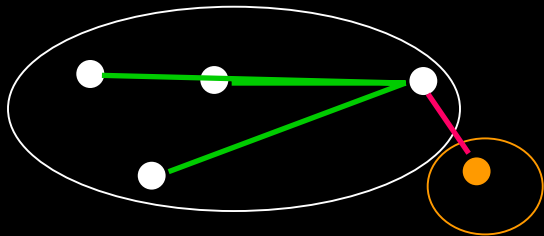
This object forms a **splinter set**.



# The Macnaughton-Smith Method

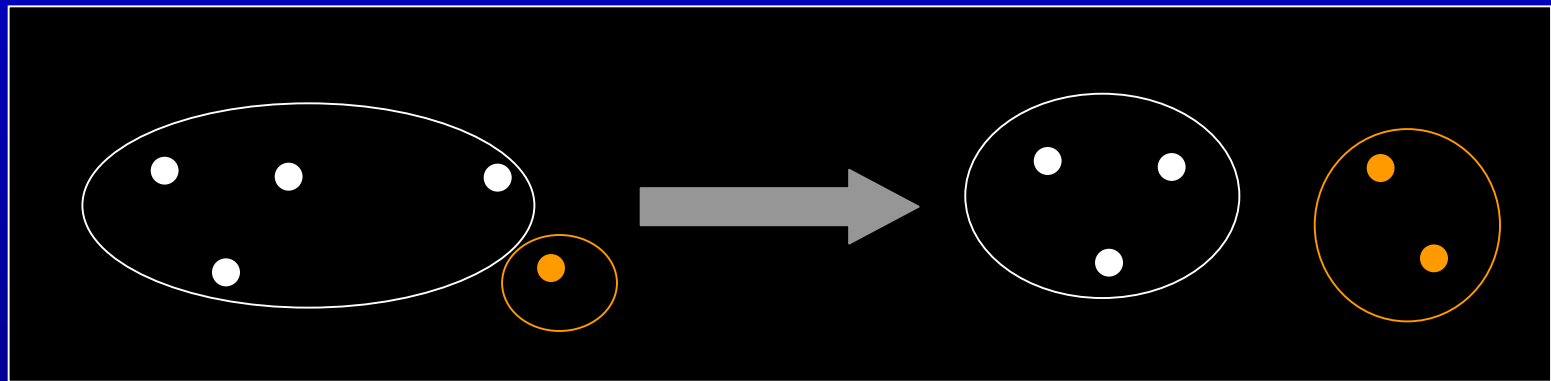
**Step 2:** For each object left in the original set, calculate the difference between the average distance from objects in the set and the average distance from objects in the splinter set.

$$\text{Diff}(X) = \frac{1}{|\text{orig set}|} \sum_{A \text{ in orig set}} D(X, A) - \frac{1}{|\text{spliter set}|} \sum_{B \text{ in orig set}} D(X, B)$$



# The Macnaughton-Smith Method

**Step 3:** If all the differences are negative, stop.  
Else, move the object with the most positive difference to the splinter set and go to Step 2.



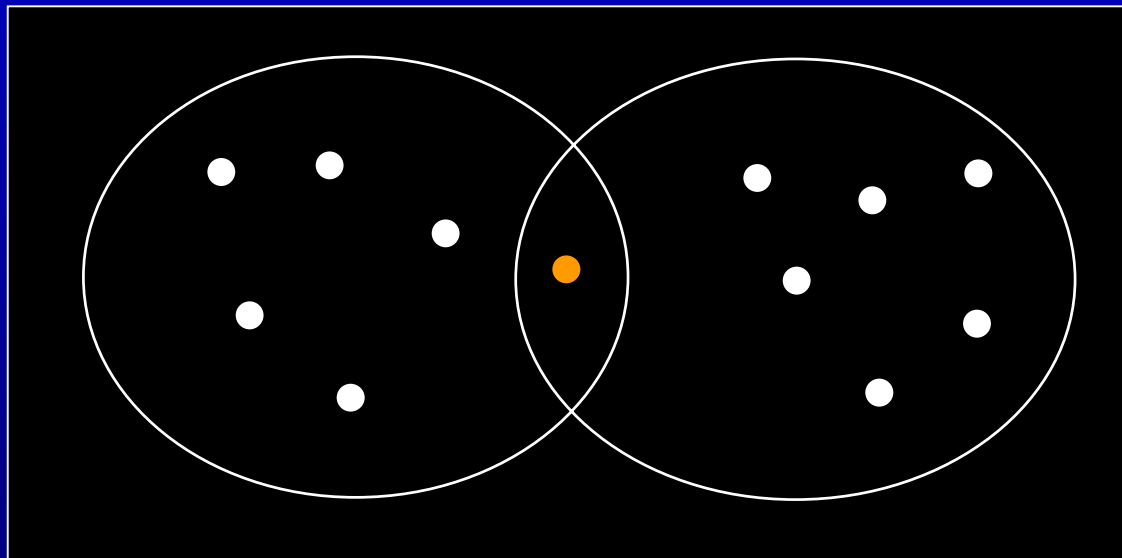
This algorithm takes  $O(n^2)$  instead of  $O(2^n)$  to split a group



# “Fuzzy” Variants

There are “fuzzy” variants to both partitioning and hierarchical methods.

These variants allow objects to be assigned to more than one cluster at a time.



# A Sample Algorithm

*Input:* A set of data points,  $\{x_1, x_2, \dots, x_n\}$ , and a threshold  $t$ .

$i = 1, k = 1$

assign  $x_1$  to cluster  $C_1$

if not all points assigned:

$i = i + 1$

$\text{dist} = \min(\text{distance between } x_i \text{ and all } x \text{ assigned})$

$C_m = \text{cluster with corresponding to dist}$

    if  $\text{dist} < t$ : assign  $x_i$  to  $C_m$

    else:

$k = k + 1$

        assign  $x_i$  to  $C_k$

Is this a partitioning algorithm, or a hierarchical algorithm?

What parts of the algorithm effect the final clustering?

What's the runtime of this algorithm?

# Choosing an Algorithm

Often, there is no clear choice of clustering algorithm to use.

Considerations include:

- Total number of objects
- Number of likely clusters
- Shape of clusters

You may want to try more than one algorithm.

# Evaluating Clusterings

Clusterings can often be quite arbitrary.

No algorithm works best on all sets of objects.

How do you know when you've clustered correctly?

Unless you know the answer ahead of time, in general you don't.

But that doesn't mean that people haven't tried to evaluate clusterings mathematically.

# Sum of Distances Method

This often-used quantitative method looks at the sum of distances between objects within a cluster (or to a representative object).

$$Q = \sum_{n=1}^K \sum_{i,j \text{ in } S_n} D(X_i - X_j) + C$$

The lower the sum, the tighter (and often “better”) the clustering.

# Back to Binary Relations

One qualitative way to look at the results of a clustering is to look at the equivalence relation defined by the clustering itself.

Is it essentially the same as the relationship between objects that the clustering algorithm was using?

Does it make sense in the context of your particular problem?

# Cautions

Clustering is often used as a “fishing” tool-- cluster first and ask questions later.

Think about the question you’re trying to answer and choose a relationship and algorithm accordingly.

Boundaries between clusters are almost never “nice”.

Don’t take them as gospel, especially if they change when different algorithms are used on the same data.

# Clustering in Bioinformatics

Objects	Variables	Algorithms
genes	sequence, organism, promoters	agglomerative
proteins	sequence, structure, function	partitioning, hierarchical
microarray data	fluorescence, time, conditions	agglomerative
small compounds	structure, bioactivity	partitioning, hierarchical
patients	treatment, outcome, location	partitioning, hierarchical



# Summary

What have we learned today?

- Some relational algebra
- Some useful relations between objects
- Data handling
- Partitioning methods (K-means, etc.)
- Hierarchical methods (single linkage, etc.)
- How to choose an algorithm
- How to evaluate a clustering