

Unsupervised Learning: Clustering

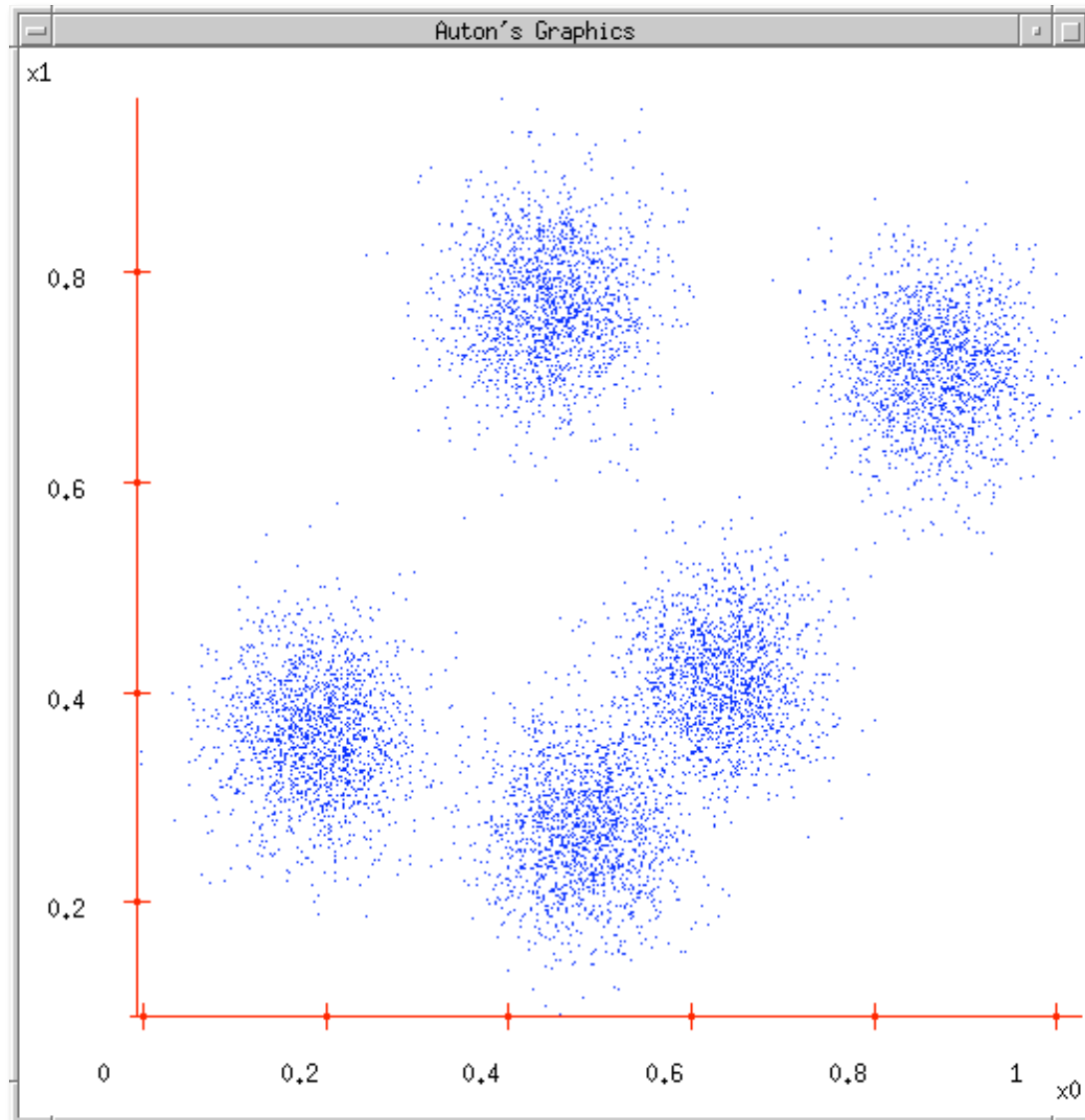
Unsupervised Learning

- Supervised learning used labeled data pairs (x, y) to learn a function $f : X \rightarrow Y$.
- But, what if we don't have labels?
- No labels = **unsupervised learning**
- Only some points are labeled = **semi-supervised learning**
 - Getting labels is expensive, so we only get a few
- **Clustering** is the unsupervised grouping of data points. It can be used for **knowledge discovery**

Clustering algorithms

- Many clustering algorithms
- Clustering typically done using a **distance measure** defined between instances
- Distance defined by instance **feature space**
- **Agglomerative** approach works bottom up:
 - Treat each instance as a cluster
 - Merge two closest clusters
 - Repeat until a stop condition is met
- **Top-down** approach starts cluster with all instances
 - Find a cluster to split into two or more smaller clusters
 - Repeat until stop condition met

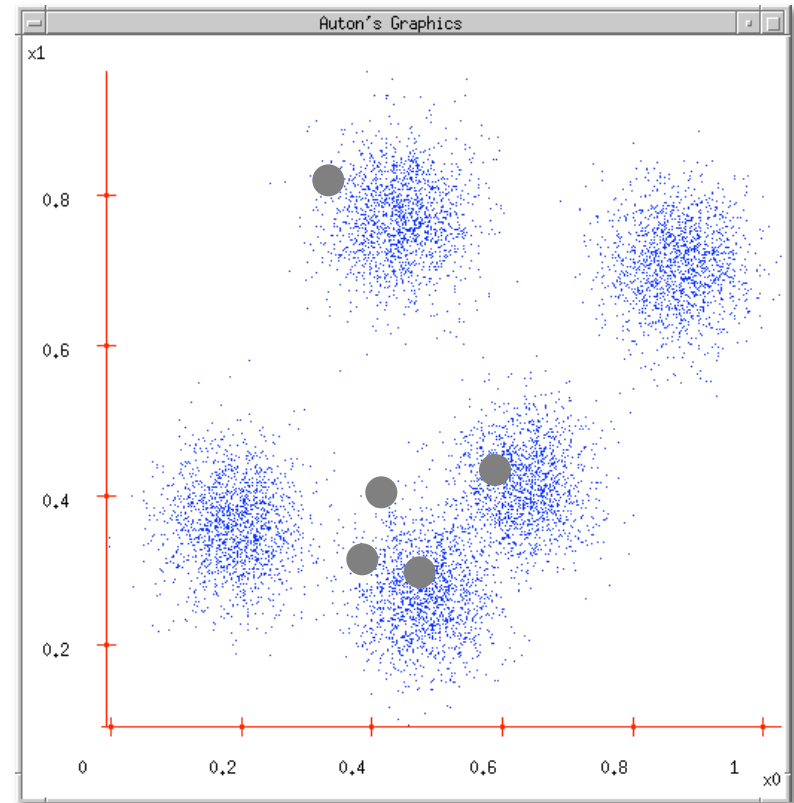
Clustering Data



K-Means Clustering

- Randomly choose k cluster center locations, aka **centroids**
- Loop until convergence
 - assign a point to cluster of the closest centroid
 - re-estimate cluster centroids based on its data assigned
- Convergence: no point is assigned to a different cluster

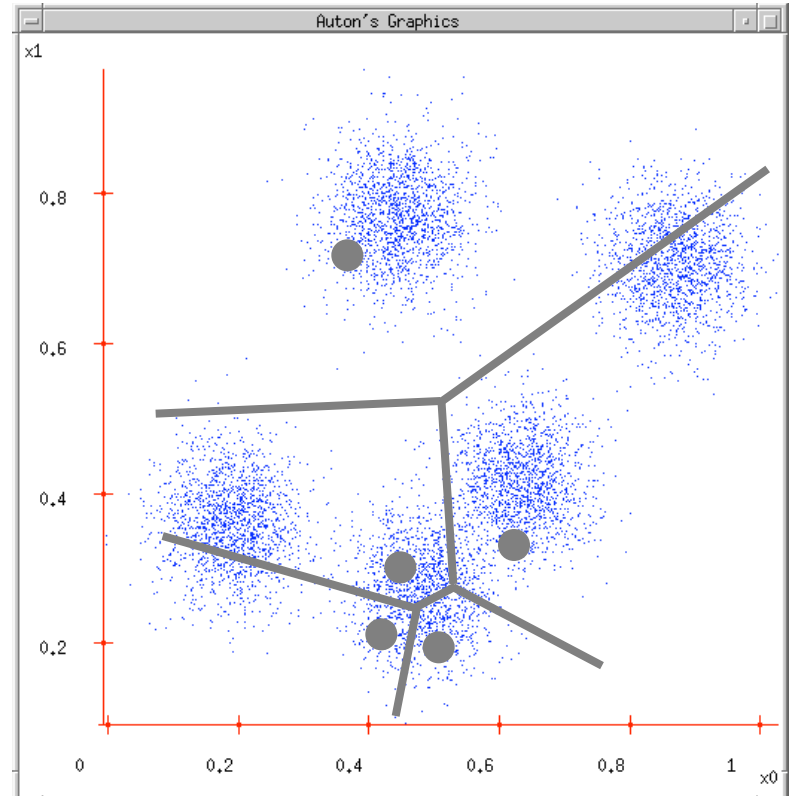
$k = 5$



K-Means Clustering

K-Means (k , data)

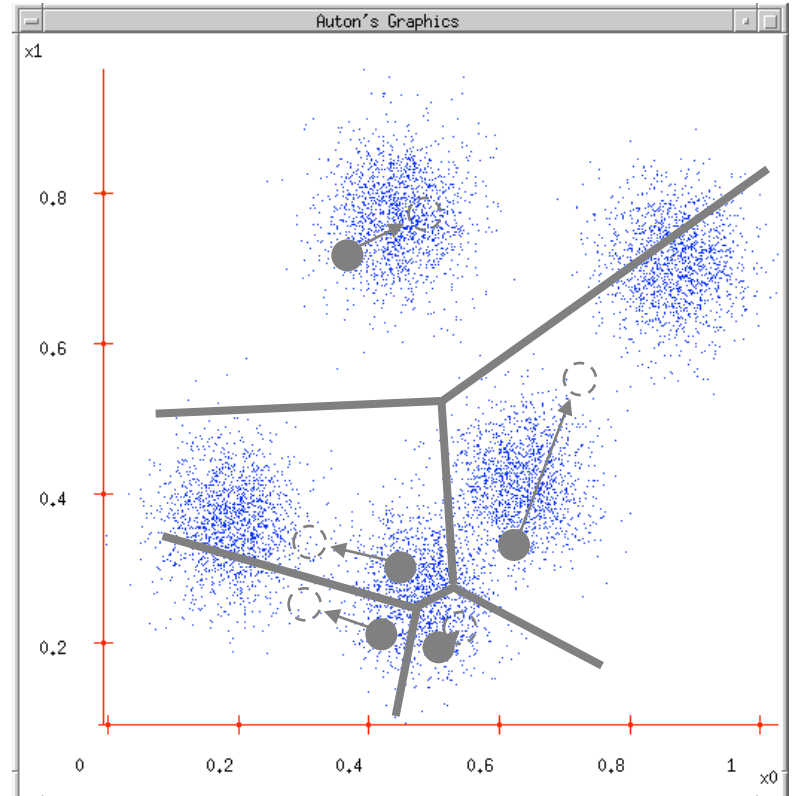
- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid.
 - Re-estimate the cluster centroids based on the data assigned to each
- Convergence: no point is assigned to a different cluster



K-Means Clustering

K-Means (k , data)

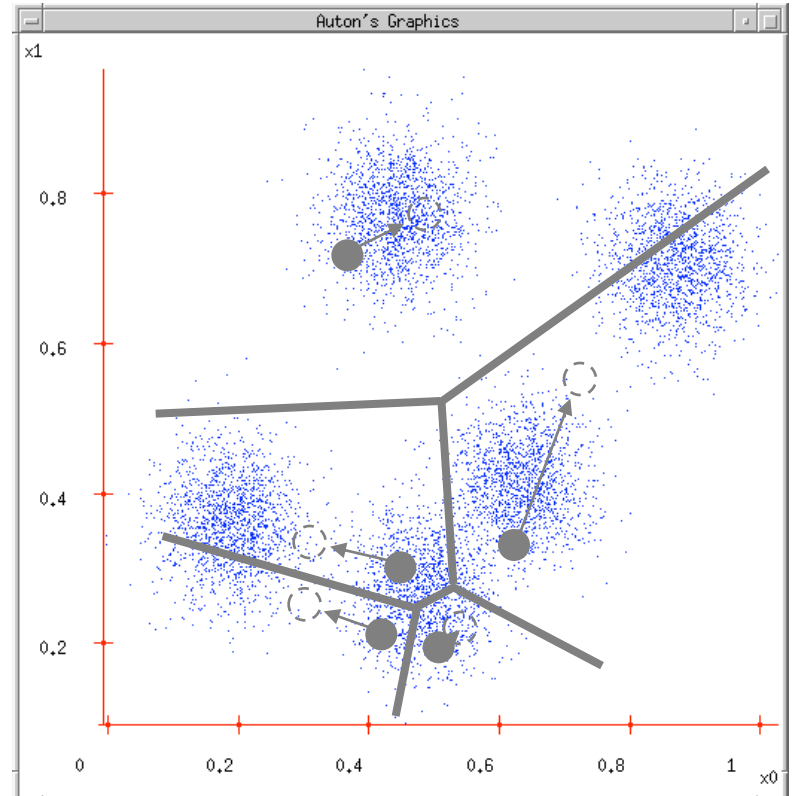
- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each
- Convergence: no point is assigned to a different cluster



K-Means Clustering

K-Means (k , data)

- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each
- Convergence: no point is assigned to a different cluster



Clusterer

Choose SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-

Cluster mode

- Use training set
- Supplied test set
- Percentage split %
- Classes to clusters evaluation
- Store clusters for visualization

Ignore attributes

Start

Stop

Result list (right-click for options)

11:17:51 - SimpleKMeans

Clusterer output

```

within-cluster sum of squared errors: 7.6174566256574

Initial starting points (random):

Cluster 0: 6.1,2.9,4.7,1.4,Iris-versicolor
Cluster 1: 6.2,2.9,4.3,1.3,Iris-versicolor
Cluster 2: 6.9,3.1,5.1,2.3,Iris-virginica

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute          Full Data          Cluster#
                   (150.0)            0           1           2
                   (50.0)            (50.0)      (50.0)      (50.0)
=====
sepalength         5.8433             5.936       5.006       6.588
sepalwidth         3.054              2.77        3.418       2.974
petallength        3.7587             4.26        1.464       5.552
petalwidth         1.1987             1.326       0.244       2.026
class              Iris-setosa Iris-versicolor  Iris-setosa  Iris-virginica

Time taken to build model (full training data) : 0 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      50 ( 33%)
1      50 ( 33%)
2      50 ( 33%)
    
```

Status

OK

Log



x 0

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose `SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance -R first-`

Cluster mode

Use training set

Supplied test set

Percentage split %

Classes to clusters evaluation

Store clusters for visualization

Result list (right-click for options)

- 11:17:51 - SimpleKMeans
- 11:21:09 - SimpleKMeans

Clusterer output

```

sepalength    5.8433    5.8885    5.006    6.8462
sepalwidth    3.054     2.7377    3.418    3.0821
petallength   3.7587    4.3967    1.464    5.7026
petalwidth    1.1987    1.418     0.244    2.0795

Time taken to build model (full training data) : 0 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      61 ( 41%)
1      50 ( 33%)
2      39 ( 26%)

Class attribute: class
Classes to Clusters:

 0 1 2 <-- assigned to cluster
0 50 0 | Iris-setosa
47 0 3 | Iris-versicolor
14 0 36 | Iris-virginica

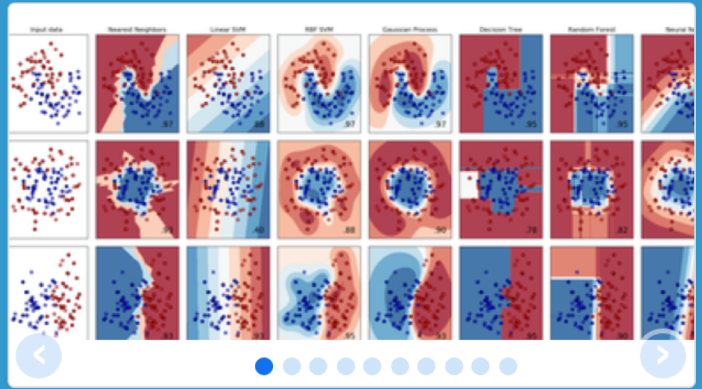
Cluster 0 <-- Iris-versicolor
Cluster 1 <-- Iris-setosa
Cluster 2 <-- Iris-virginica

Incorrectly clustered instances :      17.0      11.3333 %

```

Status

OK x 0



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ... — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization. — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics. — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction. — Examples



Previous sklearn.clust ...
Next sklearn.clust ...
Up API Reference

scikit-learn v0.19.1
Other versions

Please cite us if you use the software.

sklearn.cluster.KMeans
Examples using
sklearn.cluster.KMeans

sklearn.cluster.KMeans

```
class sklearn.cluster.KMeans (n_clusters=8, init='k-means++', n_init=10, max_iter=300, tol=0.0001, precompute_distances='auto', verbose=0, random_state=None, copy_x=True, n_jobs=1, algorithm='auto')
```

[\[source\]](#)

K-Means clustering

Read more in the [User Guide](#).

Parameters: **n_clusters** : int, optional, default: 8

The number of clusters to form as well as the number of centroids to generate.

init : {'k-means++', 'random' or an ndarray}

Method for initialization, defaults to 'k-means++':

'k-means++' : selects initial cluster centers for k-mean clustering in a smart way to speed up convergence. See section Notes in `k_init` for more details.

'random': choose k observations (rows) at random from data for the initial centroids.

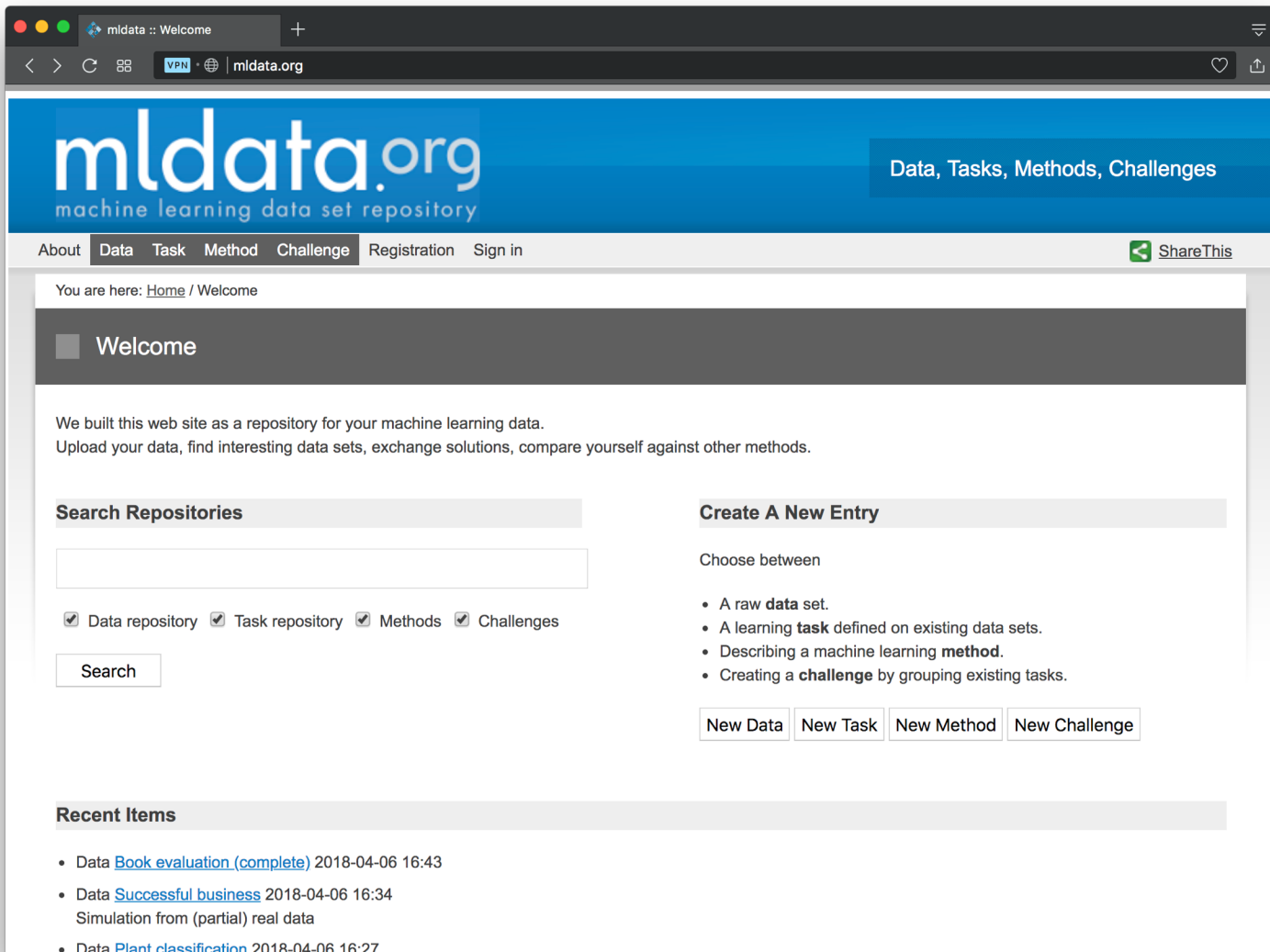
If an ndarray is passed, it should be of shape (n_clusters, n_features) and gives the initial centers.

n_init : int, default: 10

Number of time the k-means algorithm will be run with different centroid seeds. The final results will be the best output of n_init consecutive runs in terms of inertia.

max_iter : int, default: 300

- sklearn has a few test datasets, including IRIS
- Can load data directly from mldata.org & from CSV files
- mldata.org has ~900 datasets for machine learning



The screenshot shows the mldata.org website in a browser window. The browser's address bar displays 'mldata.org'. The website's header features the logo 'mldata.org machine learning data set repository' and a navigation menu with links for 'About', 'Data', 'Task', 'Method', 'Challenge', 'Registration', and 'Sign in'. A 'Share This' button is also present. Below the header, a breadcrumb trail indicates 'You are here: Home / Welcome'. The main content area is titled 'Welcome' and contains the following text: 'We built this web site as a repository for your machine learning data. Upload your data, find interesting data sets, exchange solutions, compare yourself against other methods.'

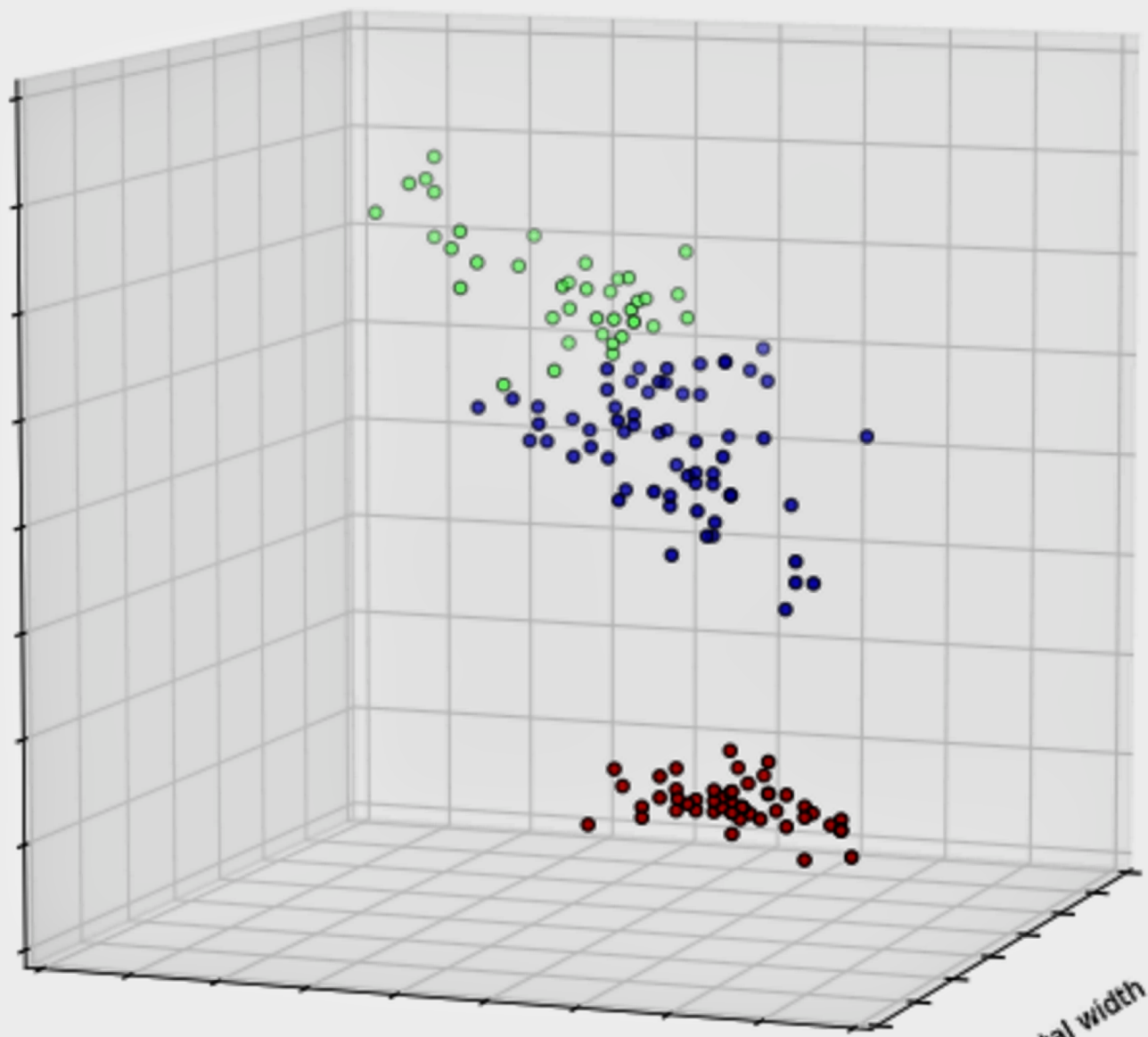
The page is divided into two main sections:

- Search Repositories:** Includes a search input field, a 'Search' button, and four checked checkboxes for 'Data repository', 'Task repository', 'Methods', and 'Challenges'.
- Create A New Entry:** Features the text 'Choose between' followed by a bulleted list:
 - A raw **data** set.
 - A learning **task** defined on existing data sets.
 - Describing a machine learning **method**.
 - Creating a **challenge** by grouping existing tasks.Below the list are four buttons: 'New Data', 'New Task', 'New Method', and 'New Challenge'.

At the bottom, a 'Recent Items' section lists three entries:

- Data [Book evaluation \(complete\)](#) 2018-04-06 16:43
- Data [Successful business](#) 2018-04-06 16:34
Simulation from (partial) real data
- Data [Plant classification](#) 2018-04-06 16:27

Petal length



Sepal length

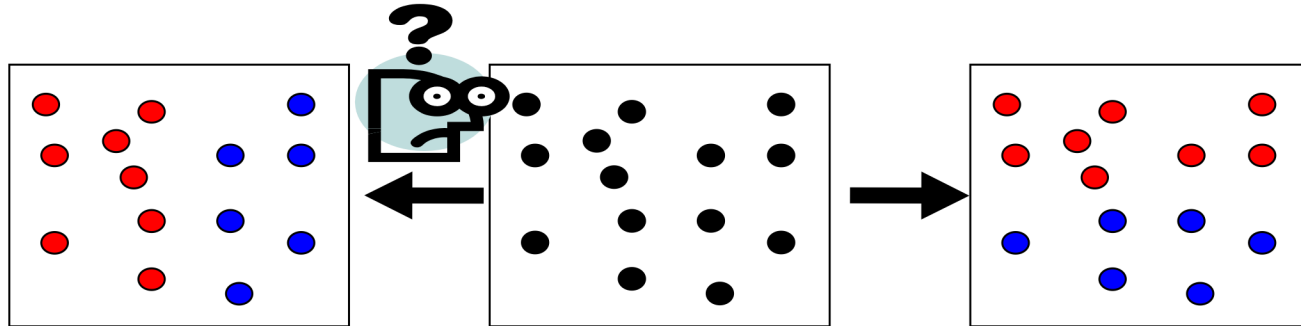
Petal width

Problems with K-Means

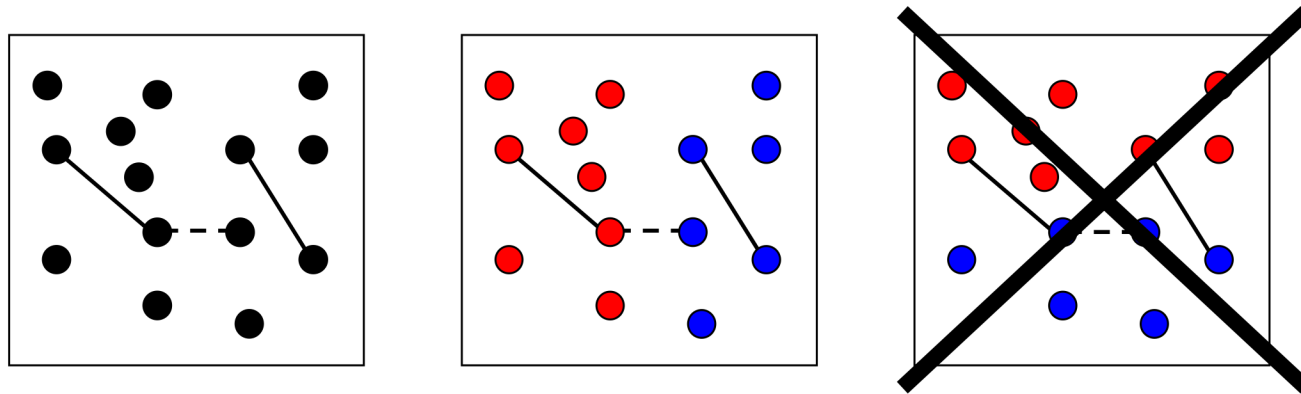
- Only works for numeric data (typically reals)
- **Very** sensitive to the initial points
 - Do many runs of k-Means, each with different initial centroids
 - Seed centroids using better method than random (e.g., Farthest-first sampling)
- Must manually choose k
 - Learn optimal k for clustering
 - Note: requires a performance measure

Problems with K-Means

- How do you tell it which clustering you want?



- Constrained clustering technique



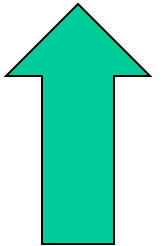
—— Same-cluster constraint
(must-link)

- - - Different-cluster constraint
(cannot-link)

Hierarchical clustering

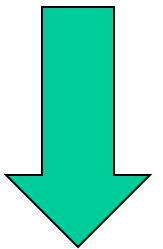
- **Agglomerative**

- **bottom up** approach: elements start as individual clusters & clusters are merged as one moves up the hierarchy



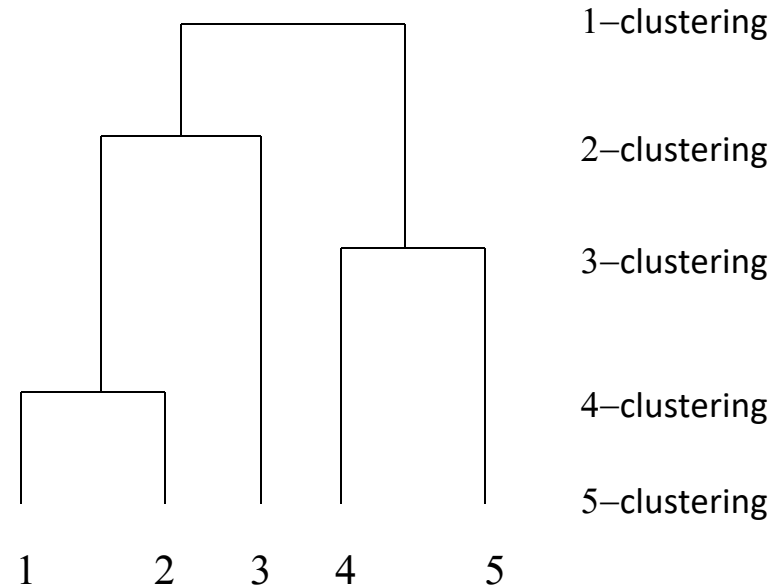
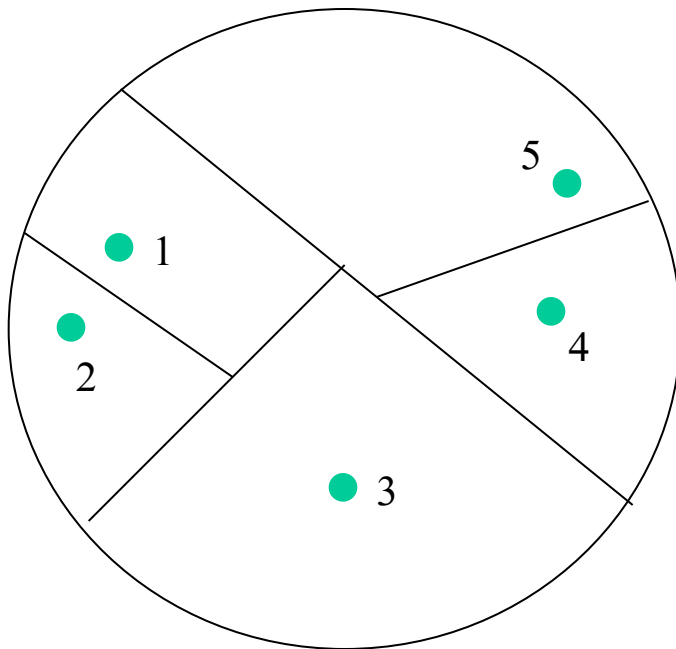
- **Divisive**

- **top down** approach: elements start as a single cluster & clusters are split as one moves down the hierarchy



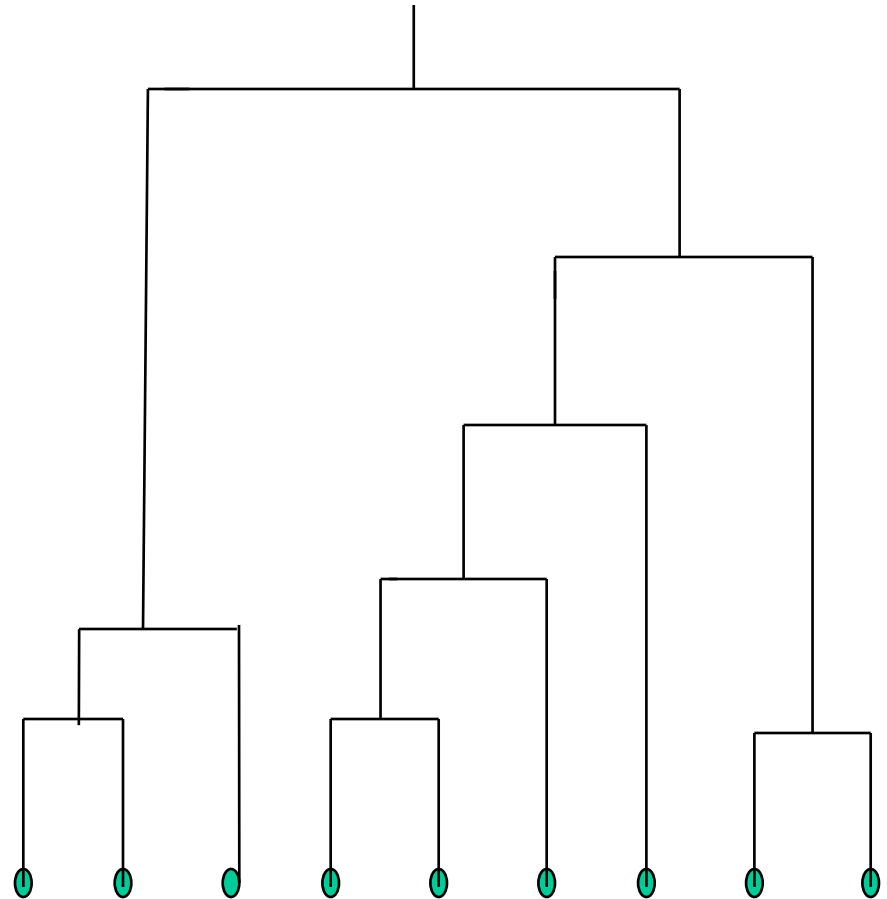
Hierarchical Clustering

Recursive partitioning/merging of a data set



Dendrogram

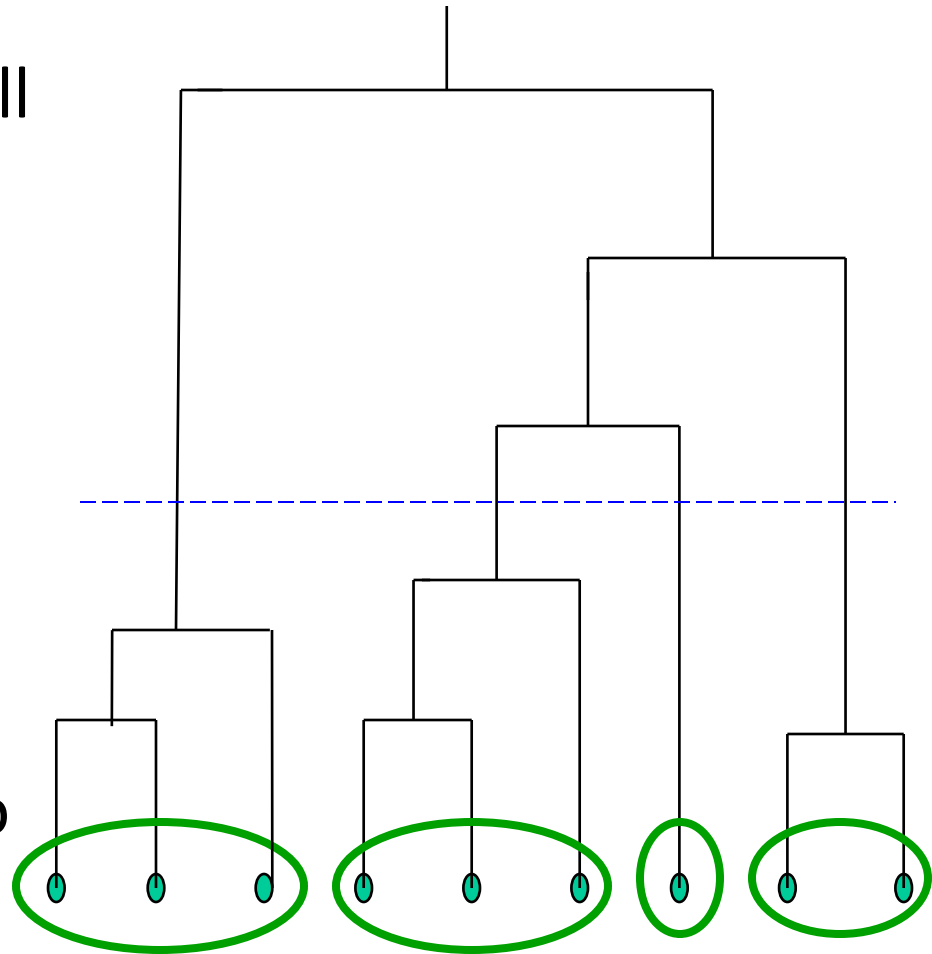
- Tree structure representing all data partitionings
- Constructed as clustering proceeds



Nine items

Dendrogram

- Tree structure representing all data partitionings
- Constructed as clustering proceeds
- Get a K-clustering by looking at **connected** components at any given level
- Frequently binary dendograms, but n-ary ones easy to obtain with minor algorithm changes

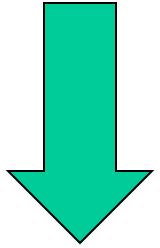


Four clusters

Hierarchical clustering advantages

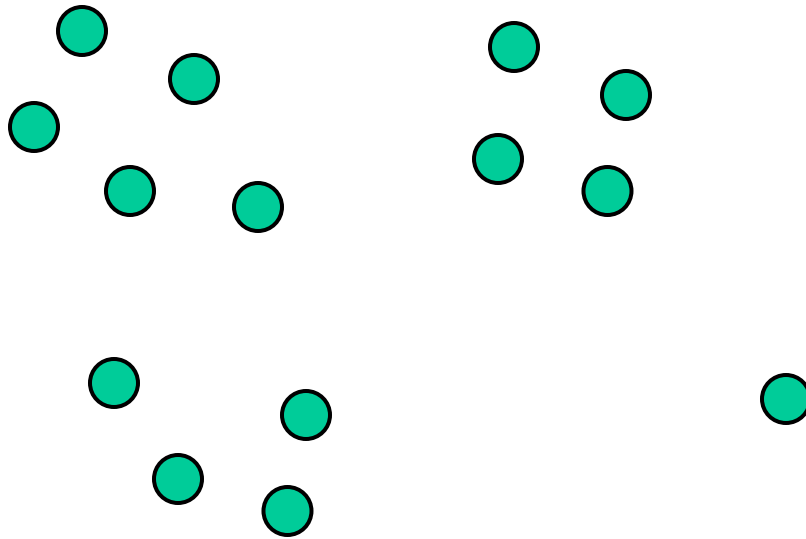
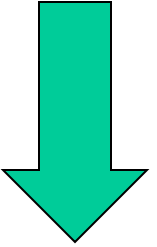
- Need not specify number of clusters
- Good for data visualization
 - See how data points interact at many levels
 - Can view data at multiple granularity levels
 - Understand how all points interact
- Specifies all of the K clusterings/partitions

Divisive hierarchical clustering

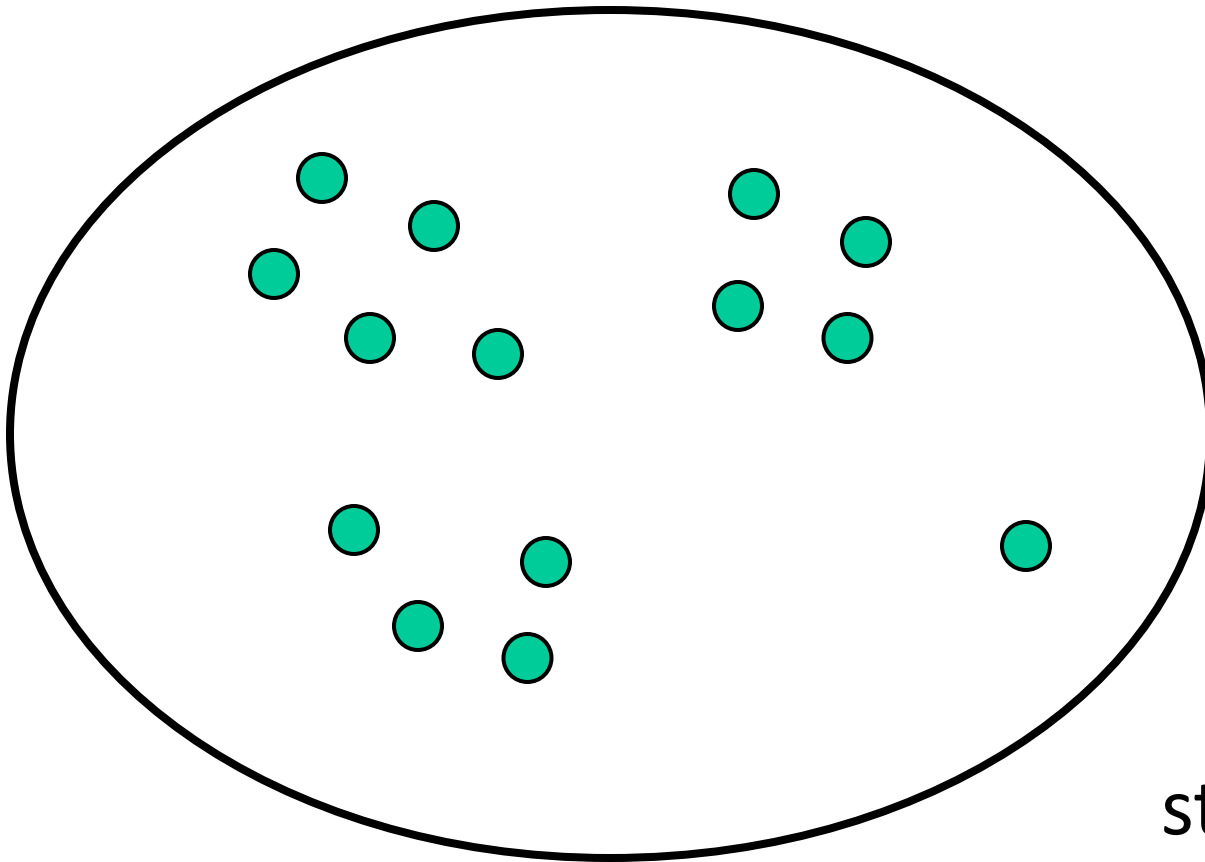
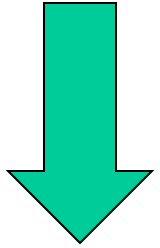


- Top-down
- Finding best partitioning of data generally exponential in time
- Common approach:
 - Let \mathbf{C} be a set of clusters
 - Initialize \mathbf{C} to be a one-clustering of data
 - While there exists a cluster c in \mathbf{C}
 - remove c from \mathbf{C}
 - partition c into 2 clusters (c_1 and c_2) using a flat clustering algorithm (e.g., k-means)
 - Add to c_1 and c_2 \mathbf{C}
- Bisecting k-means

Divisive clustering

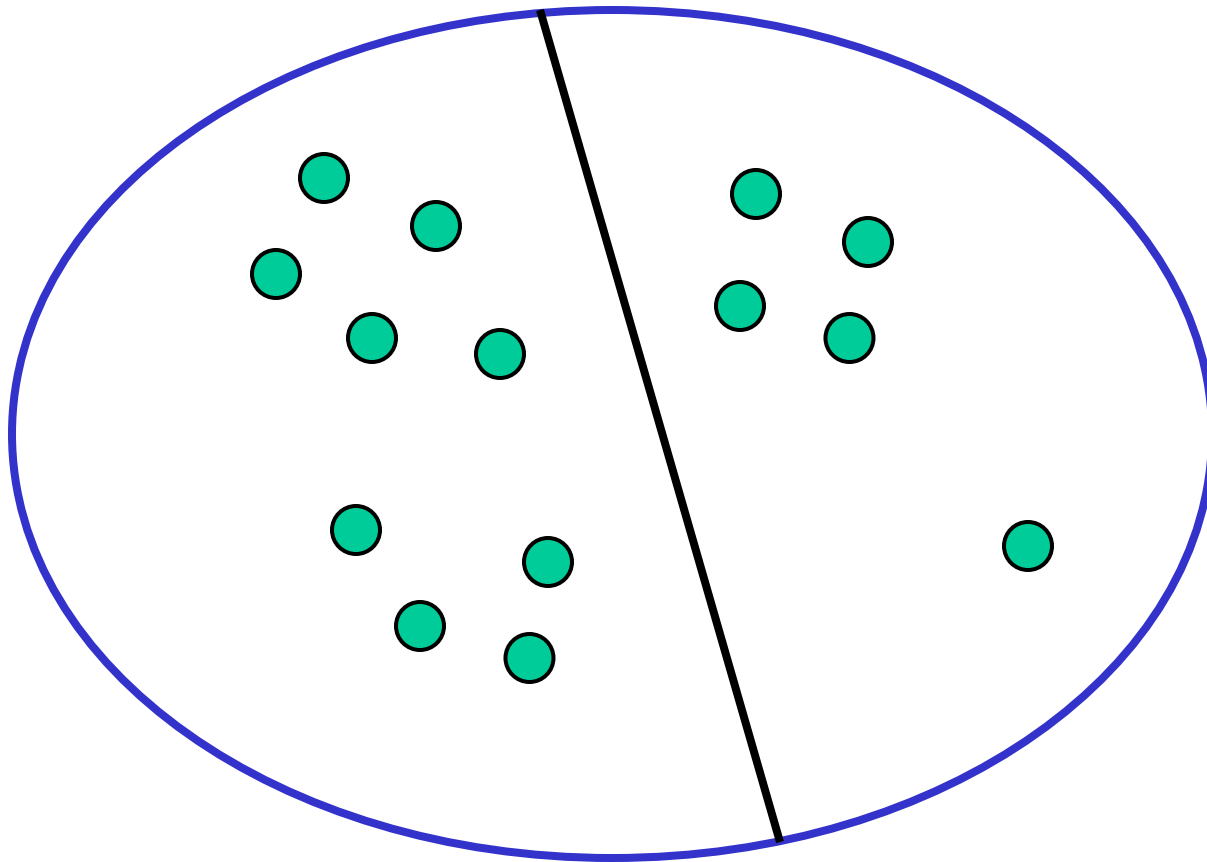
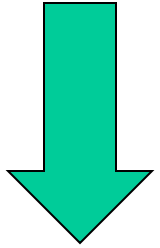


Divisive clustering



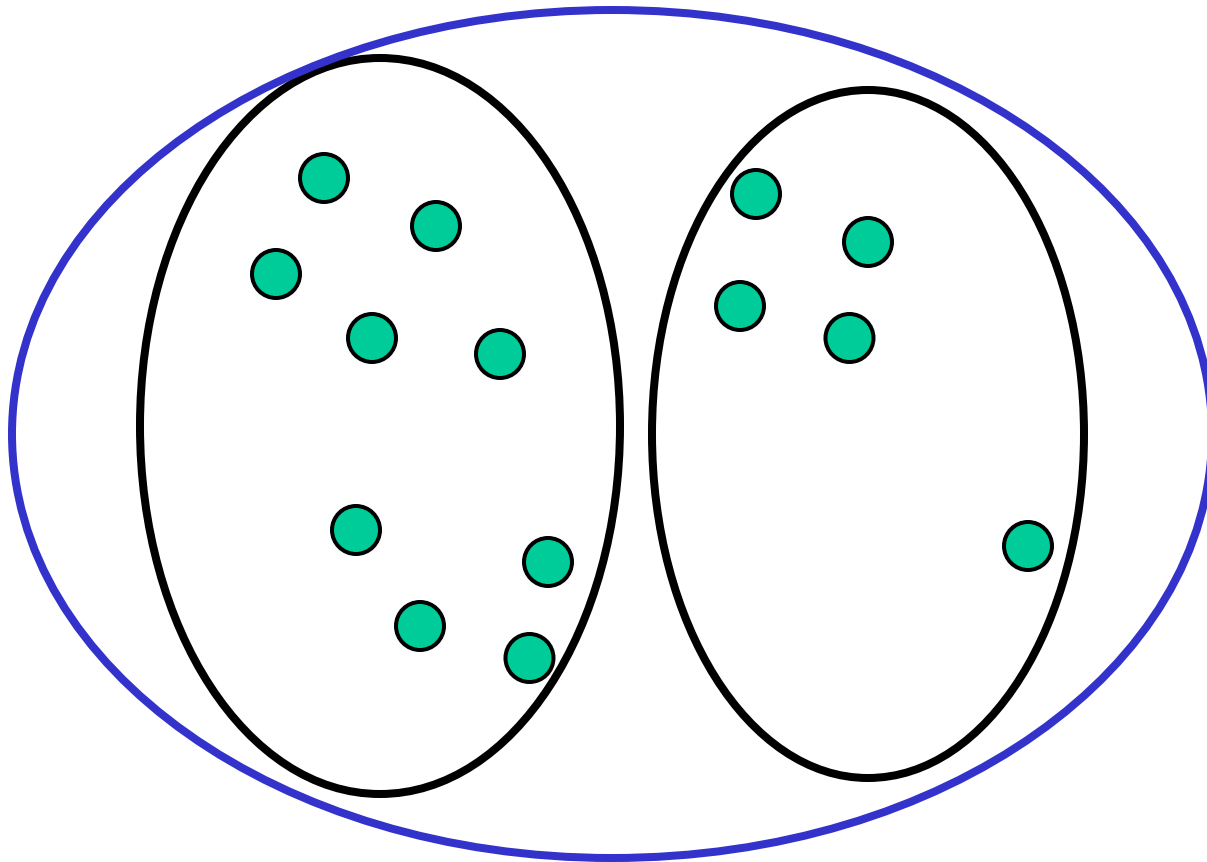
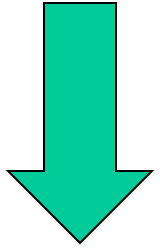
start with one
cluster

Divisive clustering

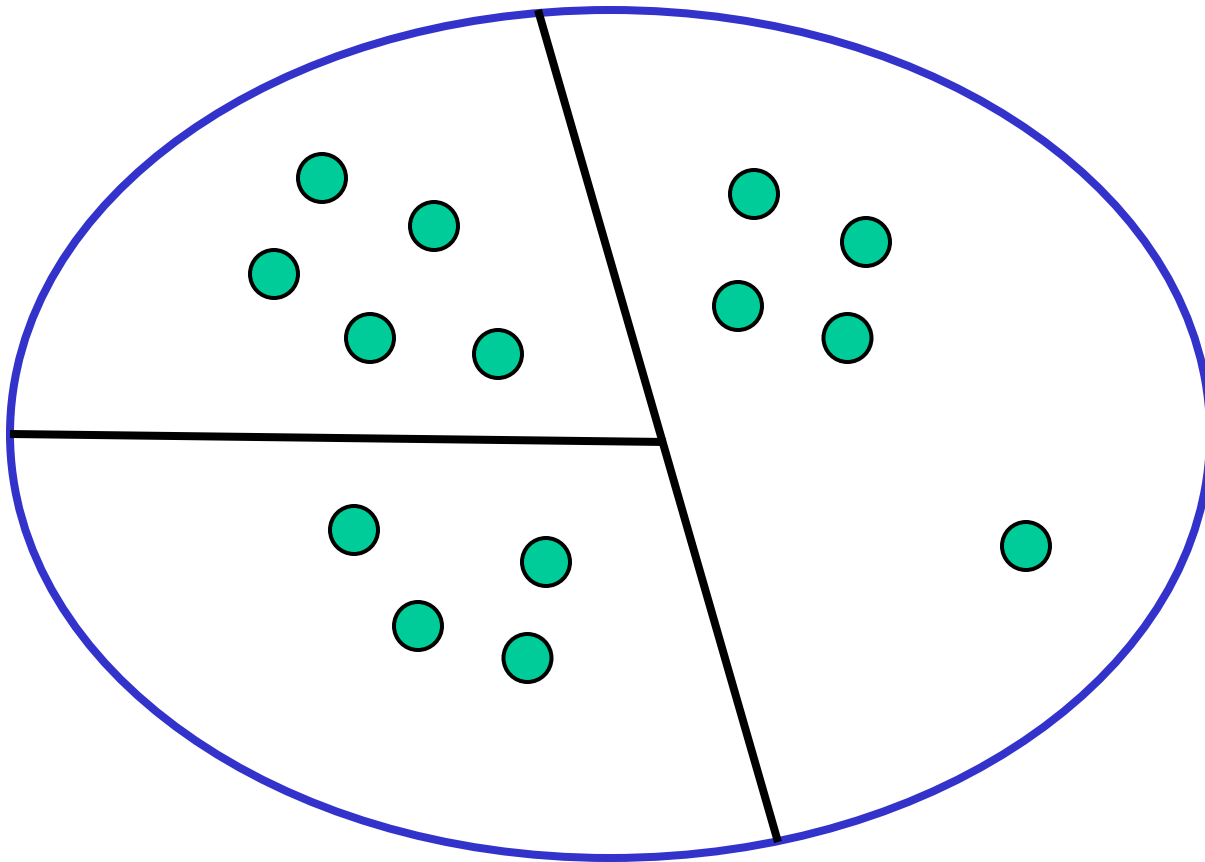
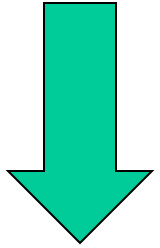


split using flat
clustering,
e.g., Kmeans

Divisive clustering

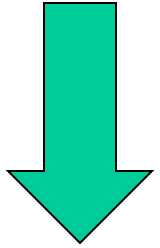


Divisive clustering

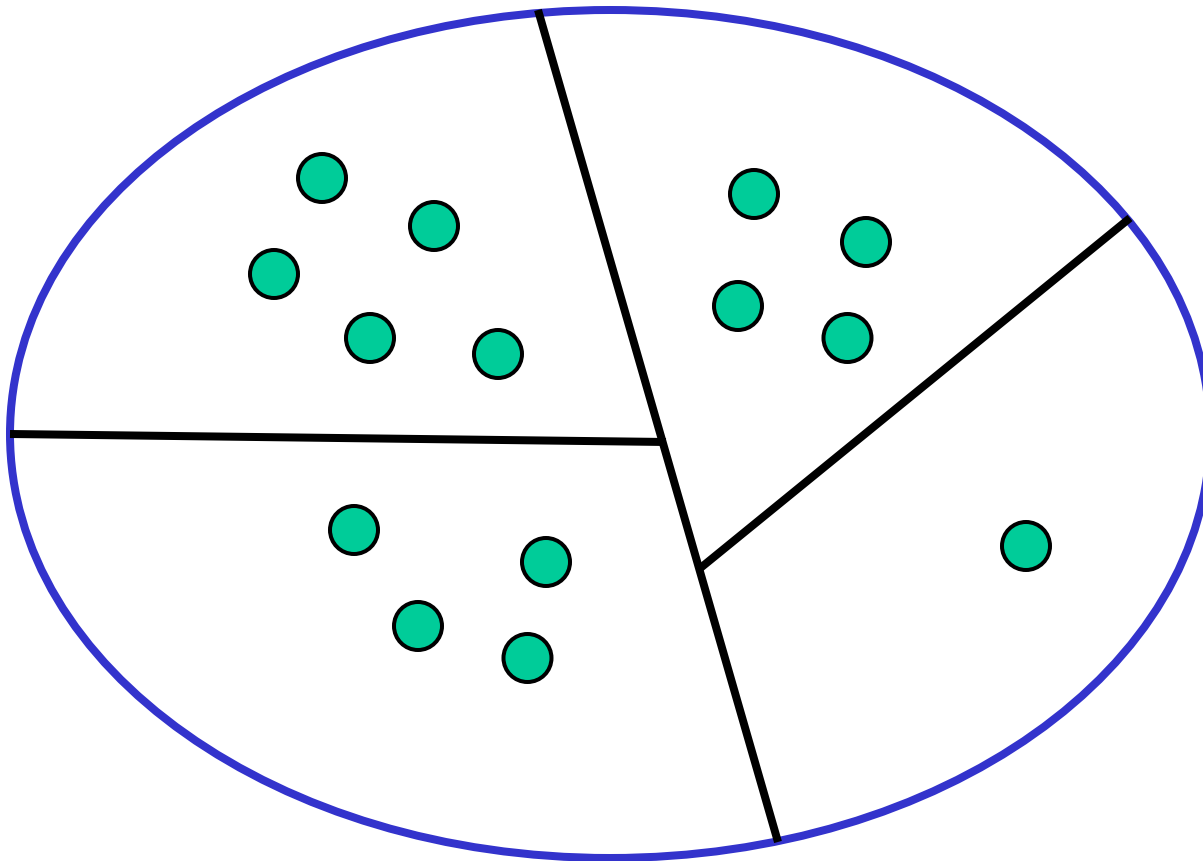


split using flat
clustering,
e.g., Kmeans

Divisive clustering

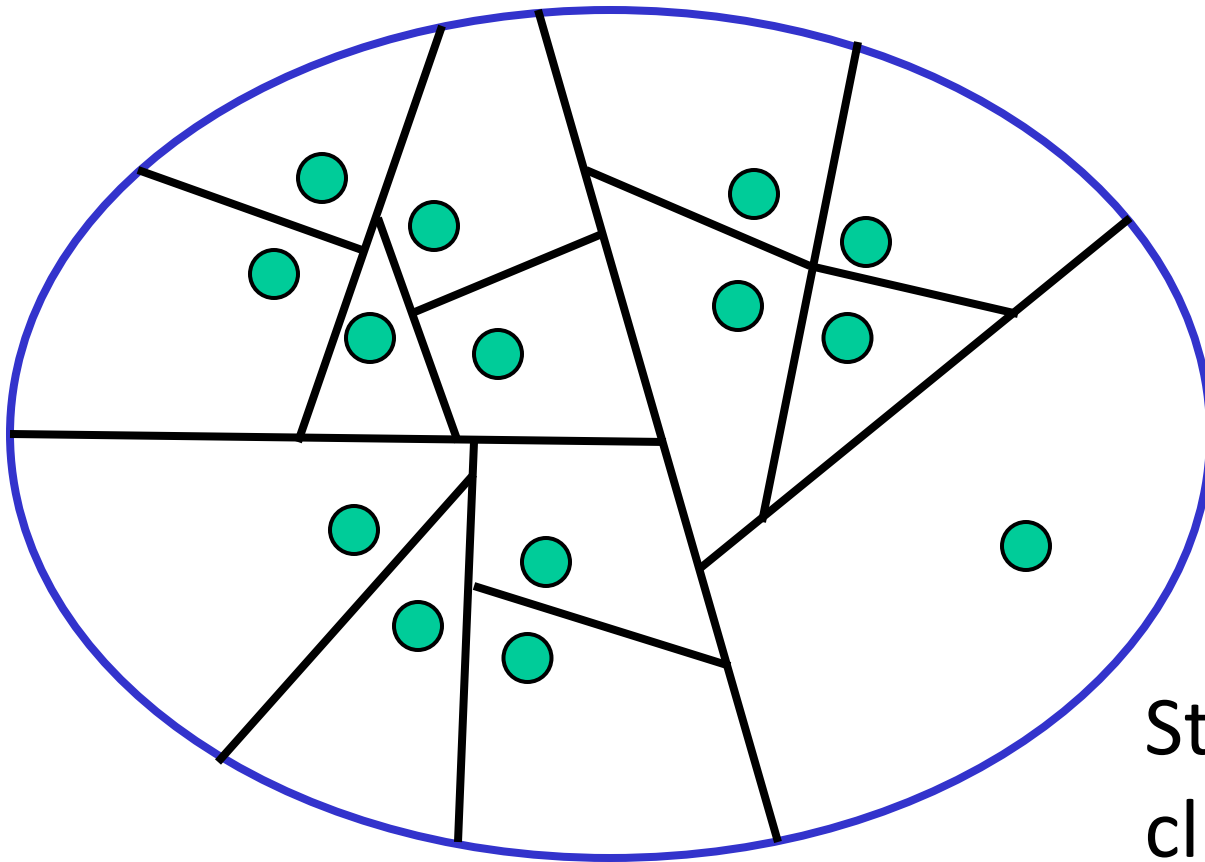
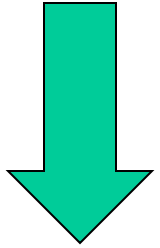


split using flat clustering



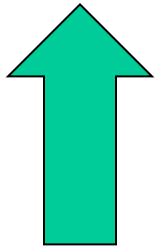
split using flat clustering,
e.g., Kmeans

Divisive clustering



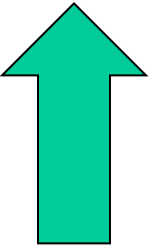
Stop when
clusters reach
some constraint

Hierarchical Agglomerative Clustering

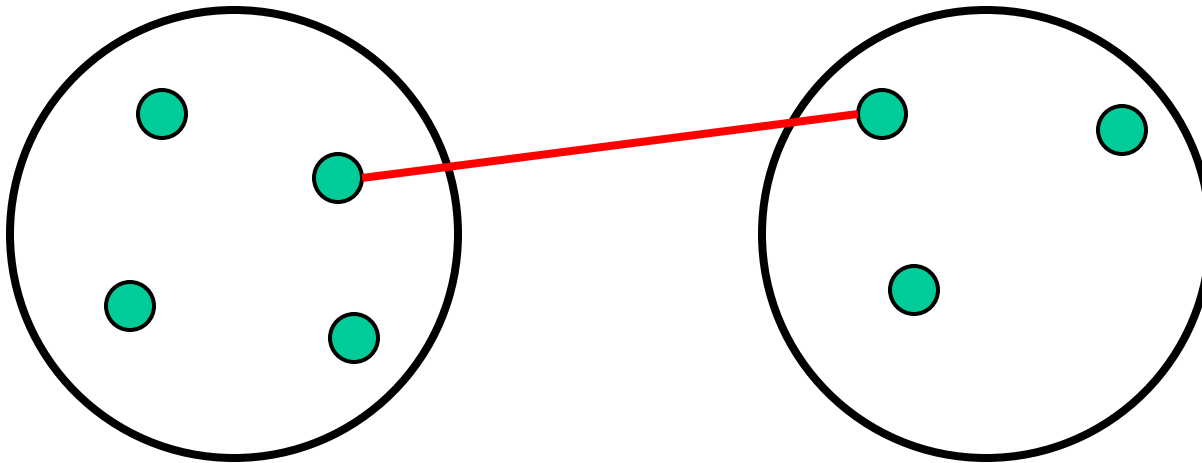


- Let \mathbf{C} be a set of clusters
- Initialize \mathbf{C} to all points/docs as separate clusters
- While \mathbf{C} contains more than one cluster
 - find c_1 and c_2 in \mathbf{C} that are **closest together**
 - remove c_1 and c_2 from \mathbf{C}
 - merge c_1 and c_2 and add resulting cluster to \mathbf{C}
- Merging history forms a binary tree or hierarchy
- **Q: How to measure distance between clusters?**

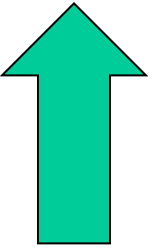
Distance between clusters



Single-link: Similarity of the *most* similar (single-link)



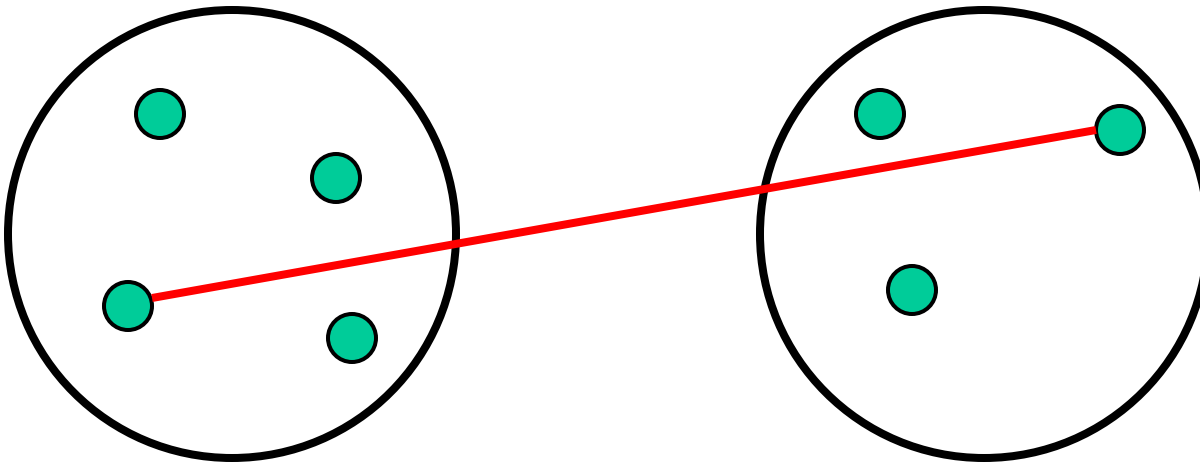
$$\max_{l \in L, r \in R} \text{sim}(l, r)$$



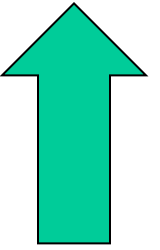
Distance between clusters

Complete-link: Similarity of the “furthest” points, the *least* similar

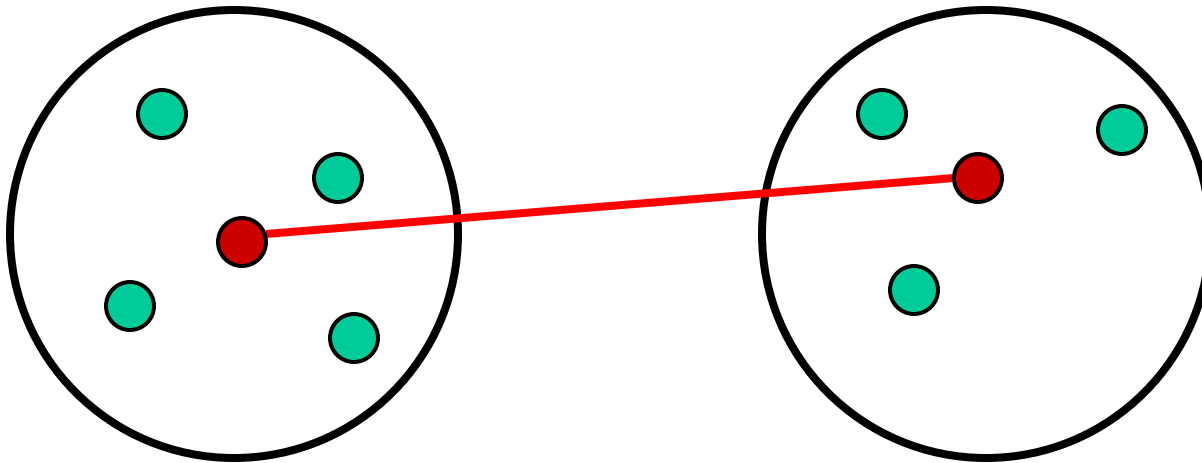
$$\min_{l \in L, r \in R} \text{sim}(l, r)$$



Distance between clusters

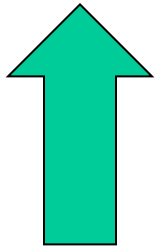


Centroid: Clusters whose centroids (centers of gravity) are the most similar

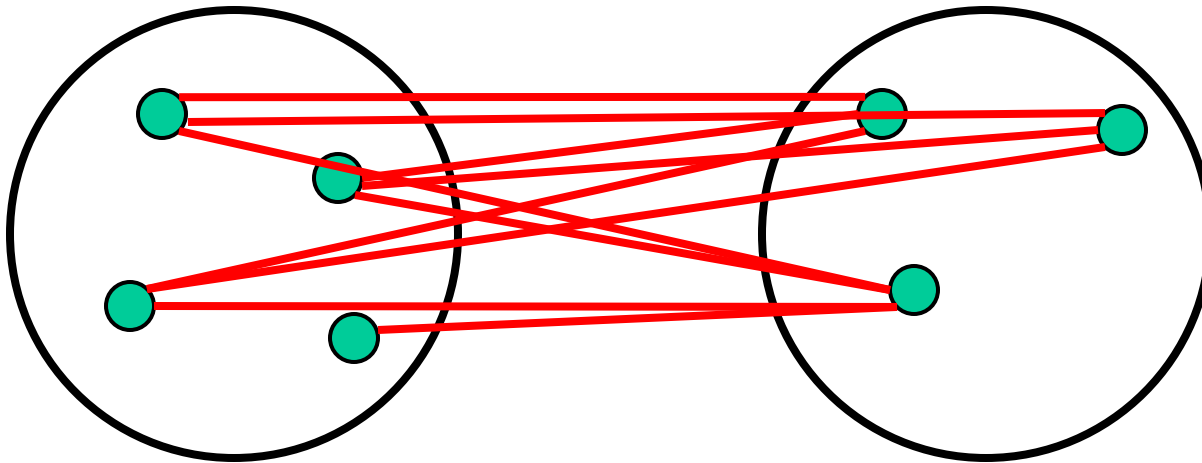


$$\|\mu(L) - \mu(R)\|^2$$

Distance between clusters



Average-link: Average similarity between all pairs of elements



$$\frac{1}{|L| \cdot |R|} \sum_{x \in L, y \in R} \|x - y\|^2$$

Knowing when to stop

- A general issue for hierarchical clustering is knowing when to stop merging/splitting a cluster
- We may have a problem specific desired range of clusters (e.g., 3-6)
- There are some general metrics for assessing the quality of a cluster
- There are also domain specific heuristics for cluster quality

Preprocess Classify Cluster Associate Select attributes Visualize

Clusterer

Choose HierarchicalClusterer -N 3 -L NEIGHBOR_JOINING -P -A "weka.core.EuclideanDistance -R first-last"

Cluster mode

 Use training set Supplied test set

Set...

 Percentage split

% 66

 Classes to clusters evaluation

(Nom) class

 Store clusters for visualization

Ignore attributes

Start

Stop

Result list (right-click for options)

15:06:44 - HierarchicalClusterer
 15:11:24 - HierarchicalClusterer
 15:11:52 - HierarchicalClusterer
 15:12:50 - HierarchicalClusterer

Clusterer output

```
Cluster 1
((((1.4:0.07344,1.5:0.07344):0.08446,((1.5:0.09914,1.4:0.09914):0.0122,(1.3:0.08407,((1.4:
```

```
Cluster 2
(((2.5:0.10622,(2.3:0.0975,(2.4:0.06047,2.3:0.06047):0.03703):0.00872):0.29975,(((2.1:0.10
```

Time taken to build model (full training data) : 0.04 seconds

=== Model and evaluation on training set ===

Clustered Instances

```
0      50 ( 33%)
1      75 ( 50%)
2      25 ( 17%)
```

Class attribute: class

Classes to Clusters:

```
0 1 2 <-- assigned to cluster
50 0 0 | Iris-setosa
0 50 0 | Iris-versicolor
0 25 25 | Iris-virginica
```

```
Cluster 0 <-- Iris-setosa
Cluster 1 <-- Iris-versicolor
Cluster 2 <-- Iris-virginica
```

Incorrectly clustered instances : 25.0 16.6667 %

Using **WARD** cluster distance measure improves results