

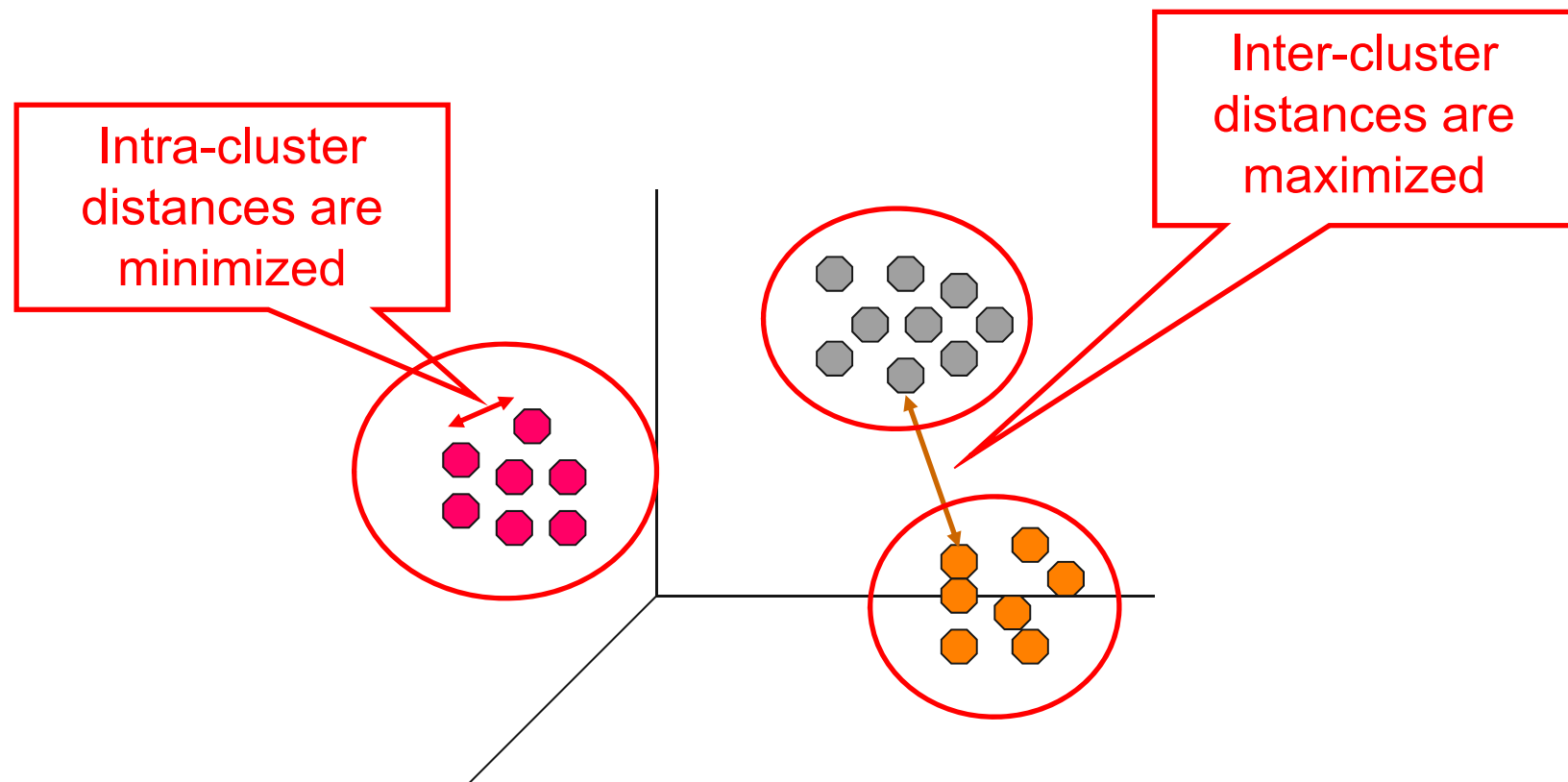
Data Mining I

Cluster Analysis

1. What is Cluster Analysis?
2. K-Means Clustering
3. Density-based Clustering
4. Hierarchical Clustering
5. Proximity Measures

1. What is Cluster Analysis?

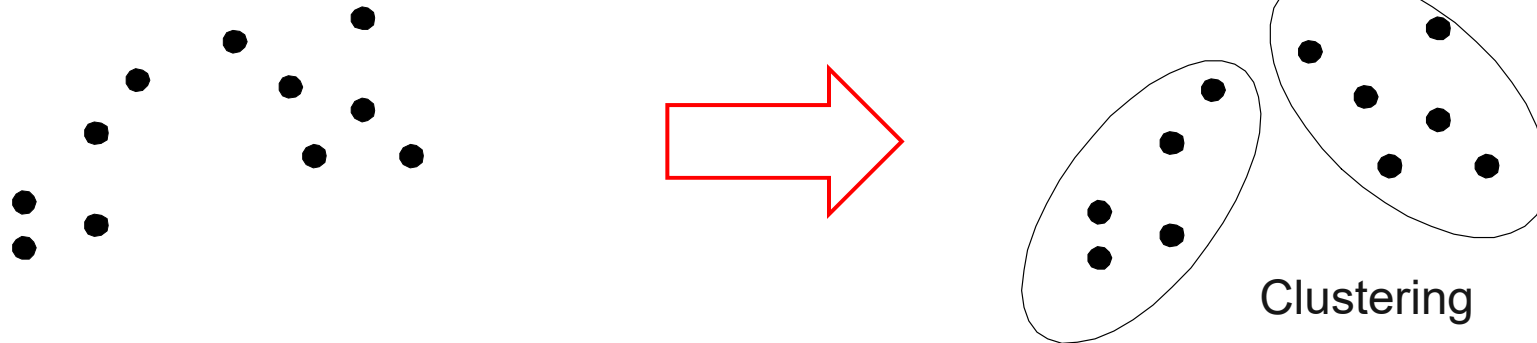
- Finding groups of objects such that
 - the objects in a group will be similar to one another
 - and different from the objects in other groups.
- Goal: Get a better understanding of the data.



Types of Clusterings

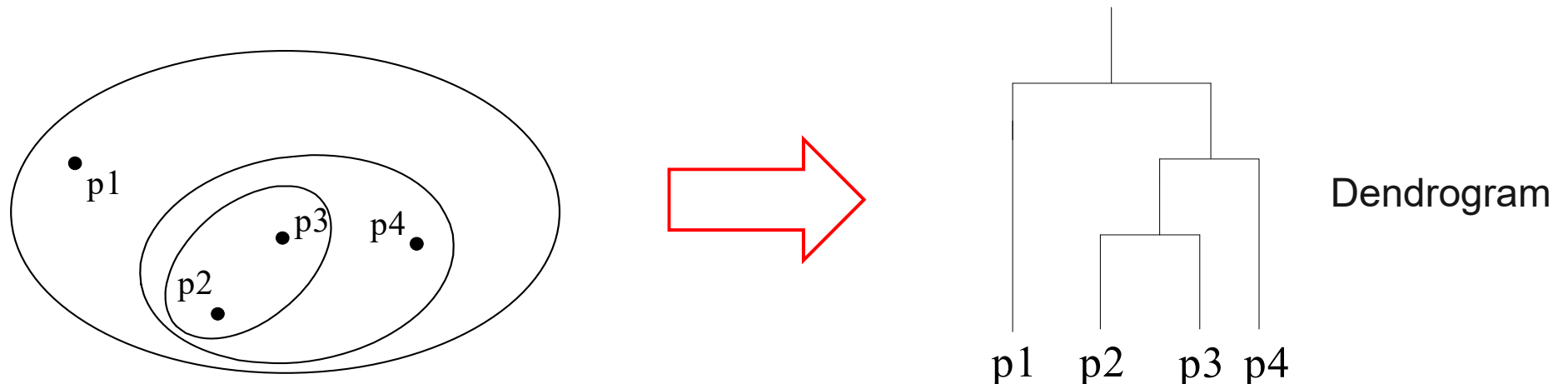
– Partitional Clustering

- A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset



– Hierarchical Clustering

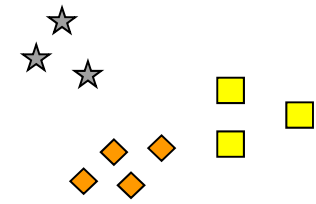
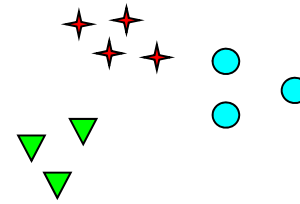
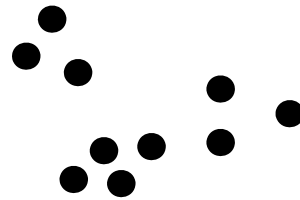
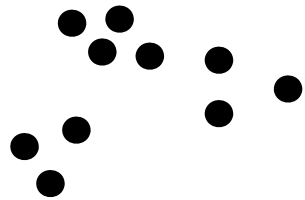
- A set of nested clusters organized as a hierarchical tree



Aspects of Cluster Analysis

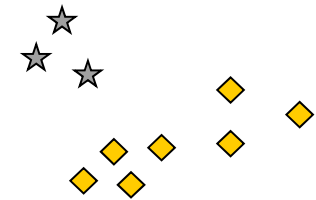
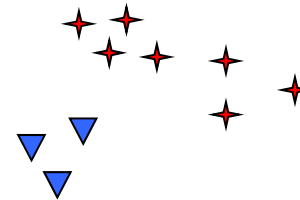
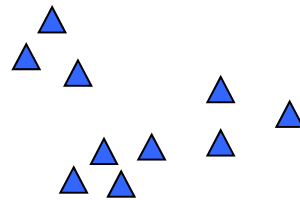
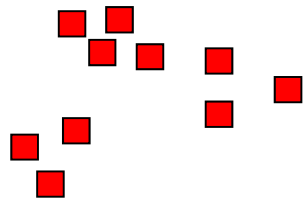
- A clustering algorithm
 - Partitional algorithms
 - Density-based algorithms
 - Hierarchical algorithms
 - ...
- A proximity (similarity, or dissimilarity) measure
 - Euclidean distance
 - Cosine similarity
 - Domain-specific similarity measures
 - ...
- Clustering Quality
 - Intra-clusters distance \Rightarrow minimized.
 - Inter-clusters distance \Rightarrow maximized.

The Notion of a Cluster is Ambiguous



How many clusters?

Six Clusters



Two Clusters

Four Clusters

The usefulness of a clustering depends on the goals of the analysis.

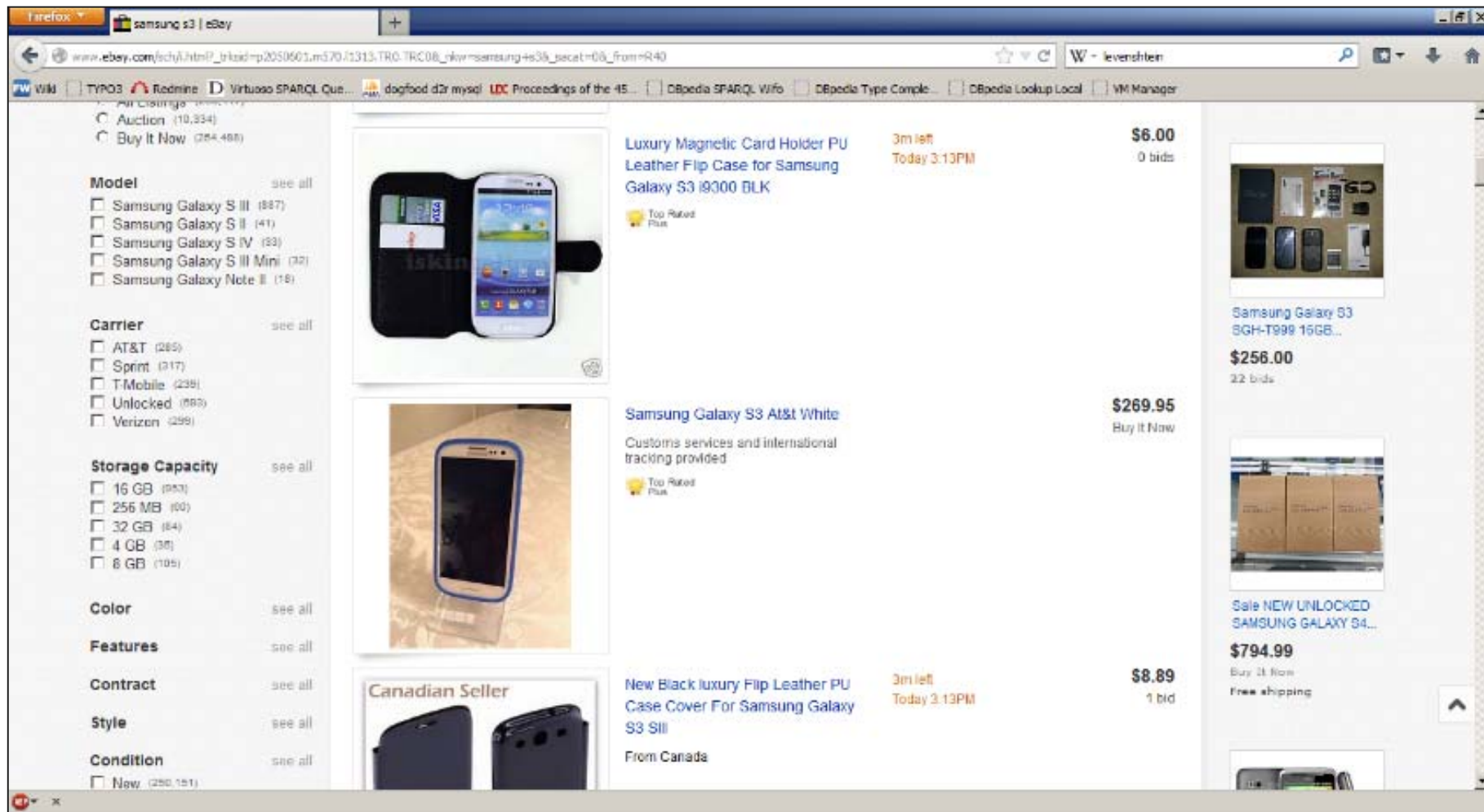
Example Application: Market Research

- Identify groups of similar customers
- Level of granularity depends on the task at hand
- Relevant customer attributes depend on the task at hand



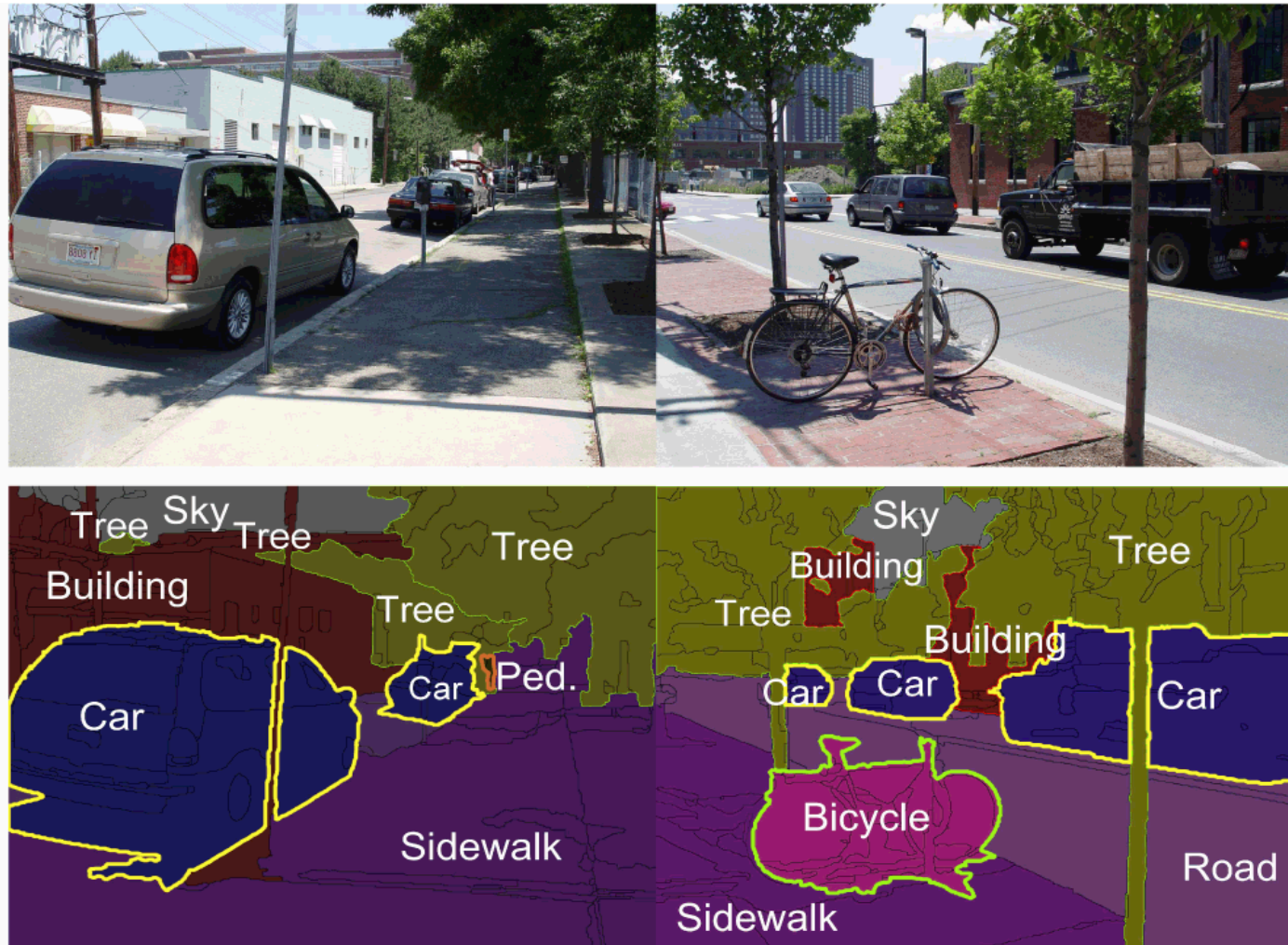
Example Application: E-Commerce

- Identify offers of the same product, e.g. on eBay



Example Application: Image Recognition

- Identify parts of an image that belong to the same object



Cluster Analysis as Unsupervised Learning

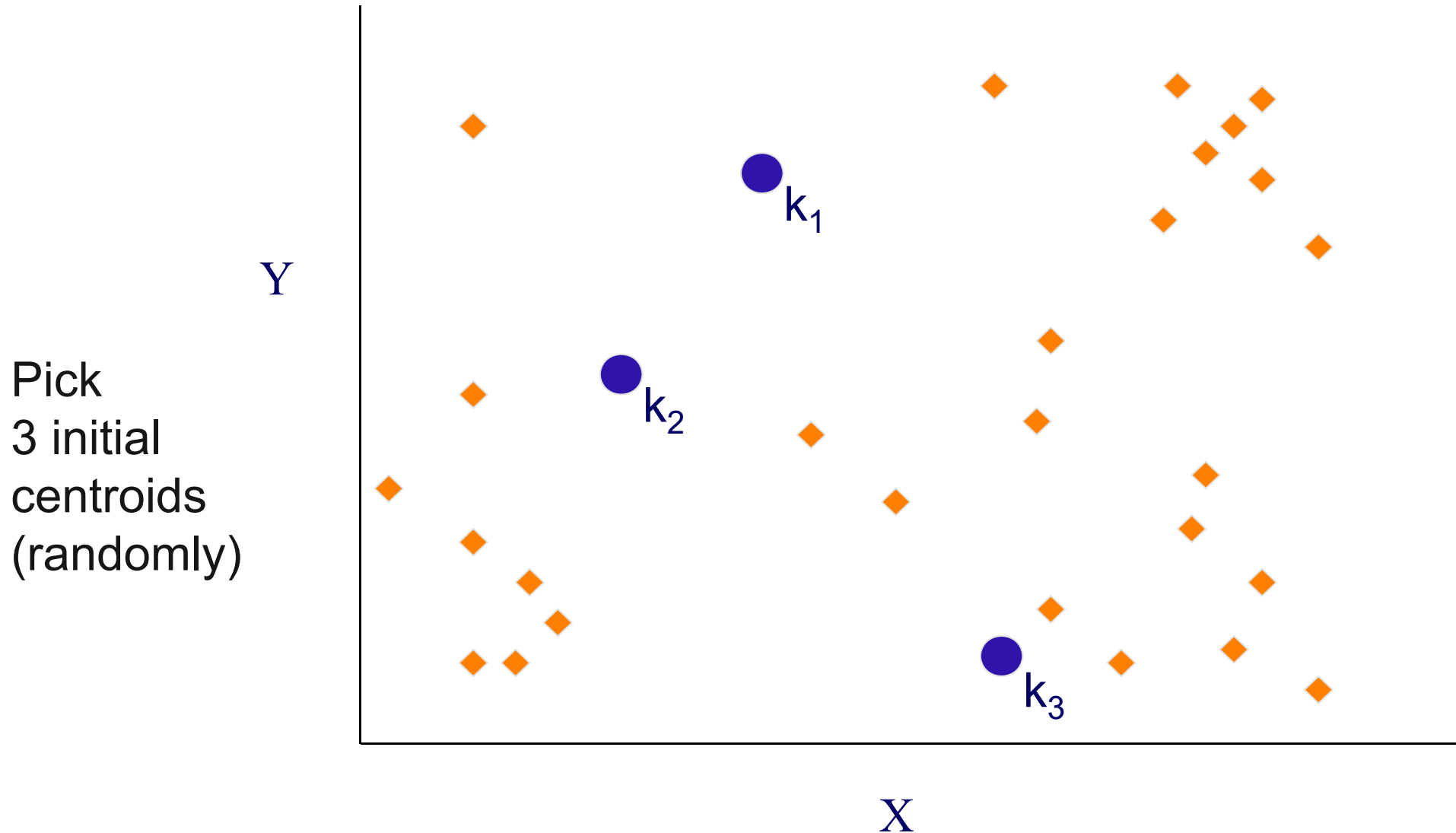
- **Supervised learning:** Discover patterns in the data that relate data attributes with a target (class) attribute.
 - These patterns are then utilized to predict the values of the target attribute in unseen data instances.
 - The set of classes is known before.
 - Training data is often provided by human annotators.
- **Unsupervised learning:** The data has no target attribute.
 - We want to explore the data to find some intrinsic structures in it.
 - The set of classes/clusters is not known before.
 - No training data is used.
- Cluster Analysis is an unsupervised learning task.

2. K-Means Clustering

- Partitional clustering approach.
- Each cluster is associated with a **centroid** (center point).
- Each point is assigned to the cluster with the closest centroid.
- Number of clusters, K , must be specified manually.
- The basic algorithm is very simple:

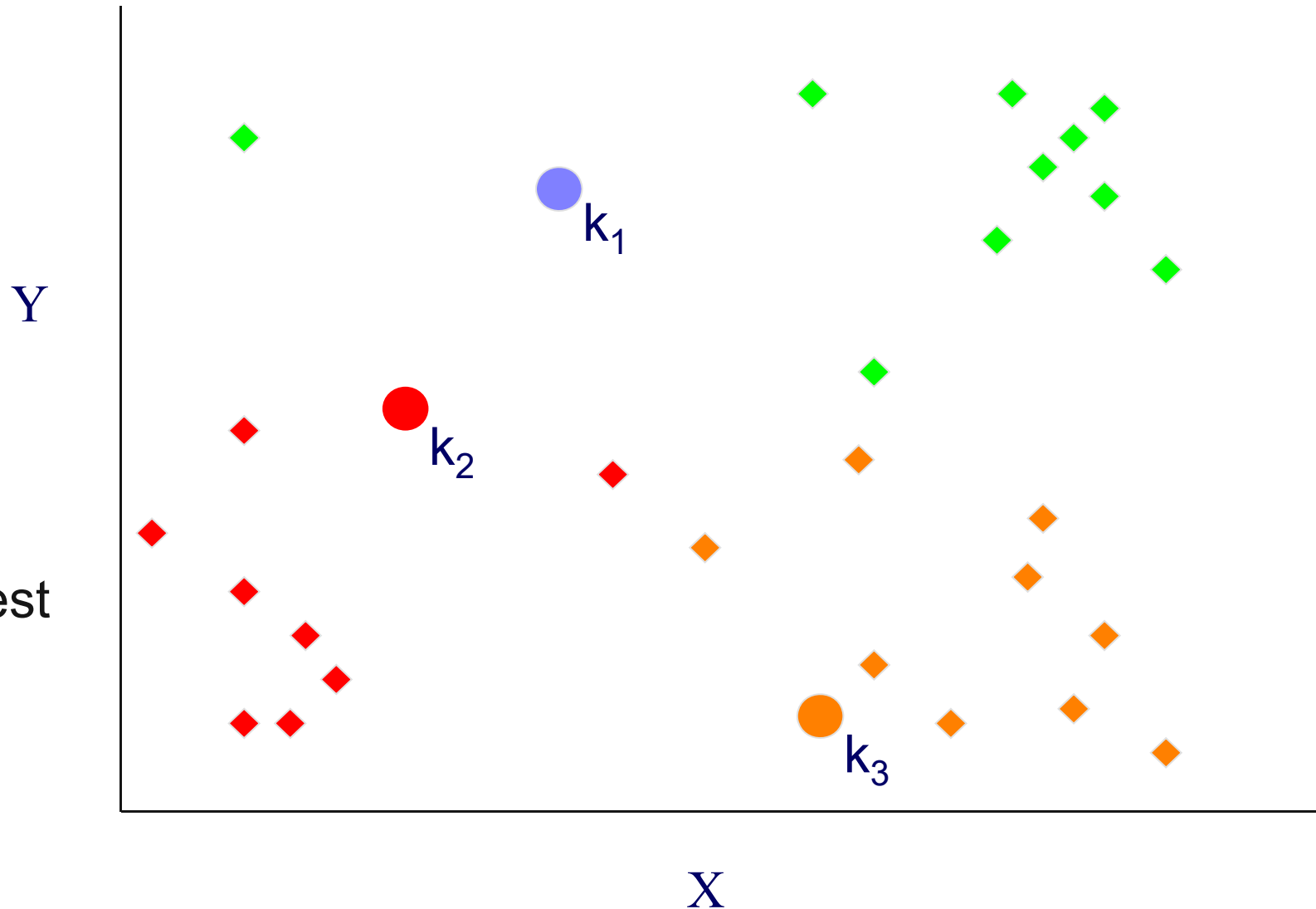
-
- 1: Select K points as the initial centroids.
 - 2: **repeat**
 - 3: Form K clusters by assigning all points to the closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** The centroids don't change
-

K-Means Example, Step 1



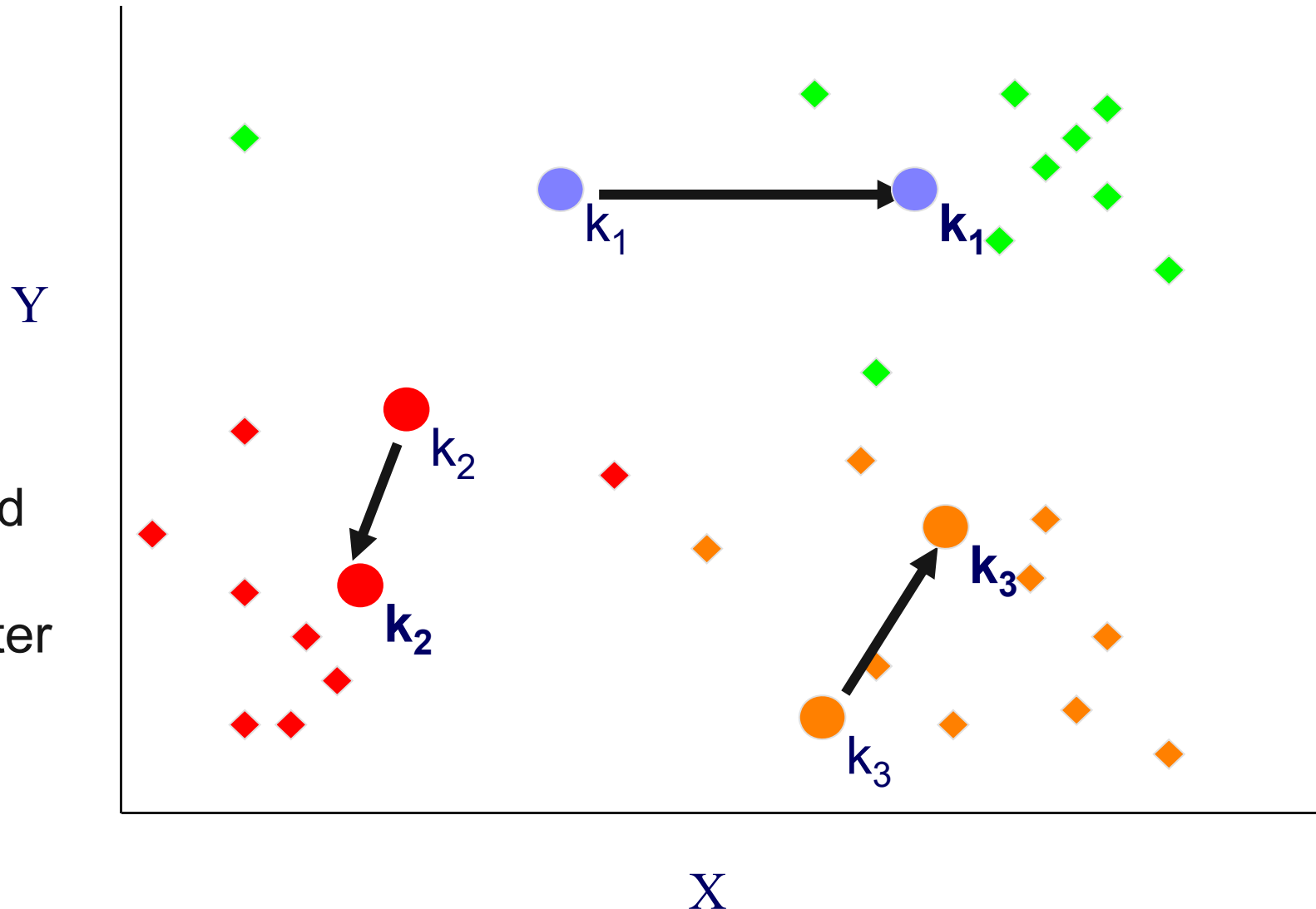
K-Means Example, Step 2

Assign
each point
to the closest
centroid



K-Means Example, Step 3

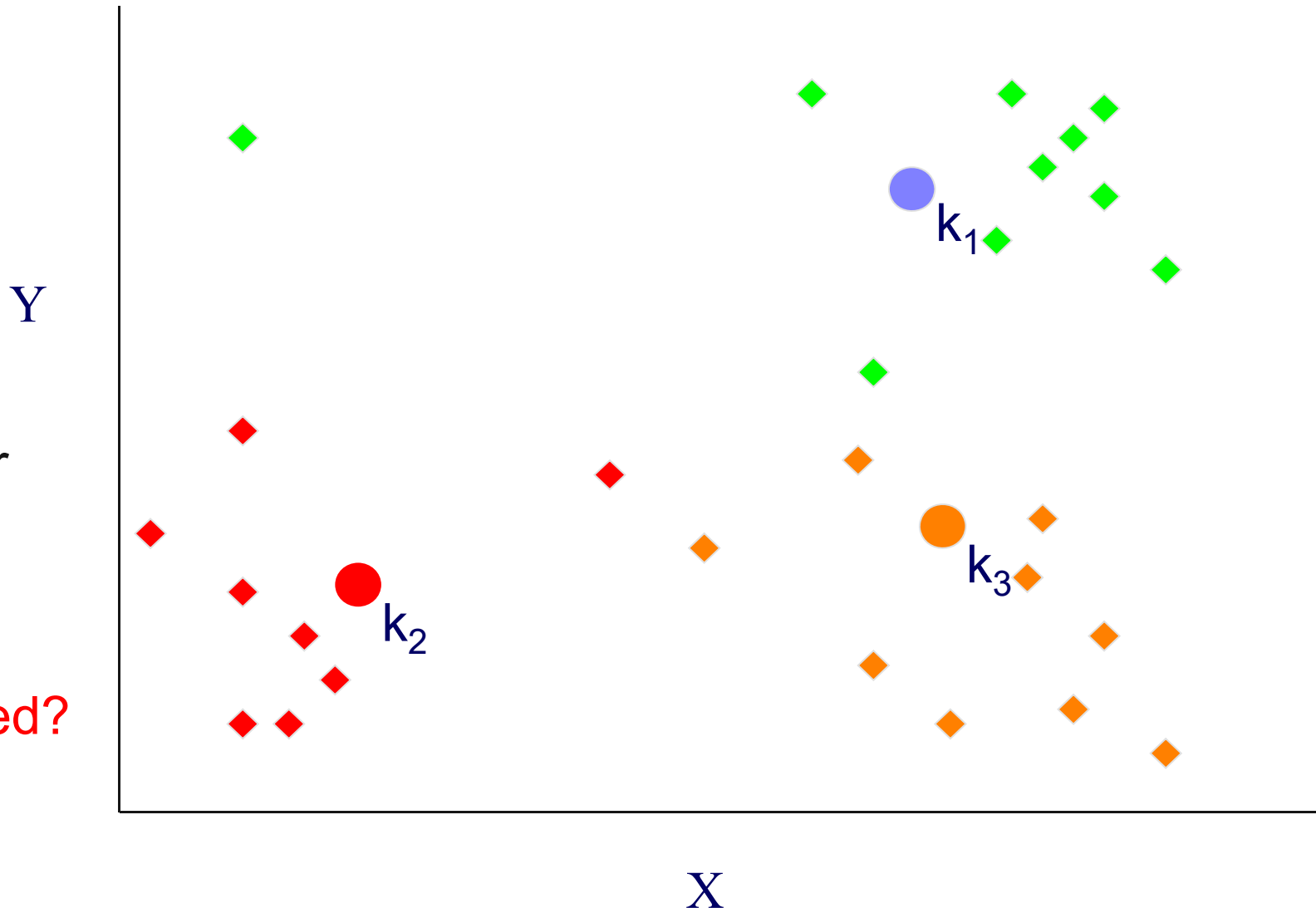
Move
each centroid
to **the mean**
of each cluster



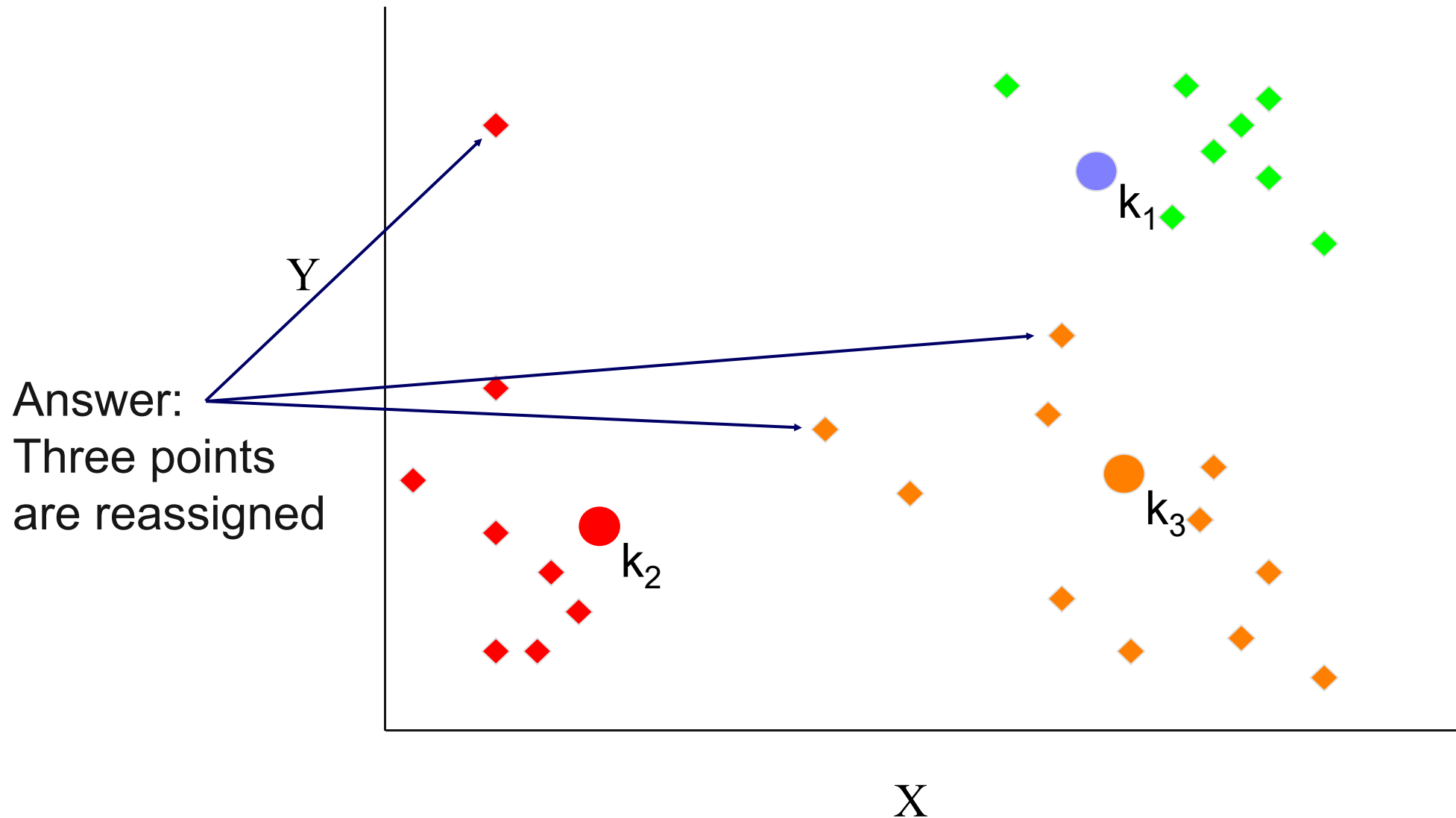
K-Means Example, Step 4

Reassign
points
closest to a
different new
cluster center

Question:
Which points
are reassigned?

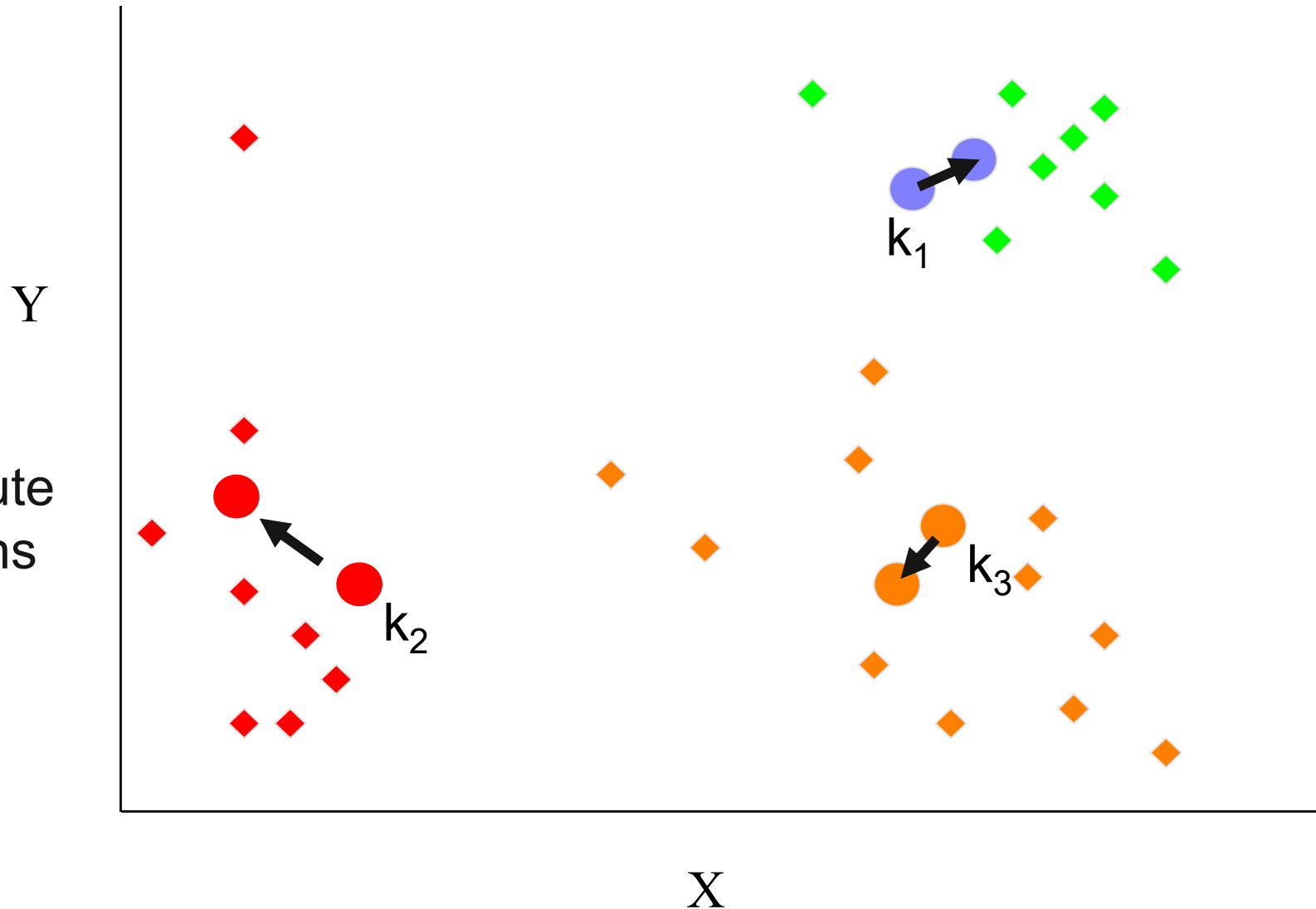


K-Means Example, Step 4

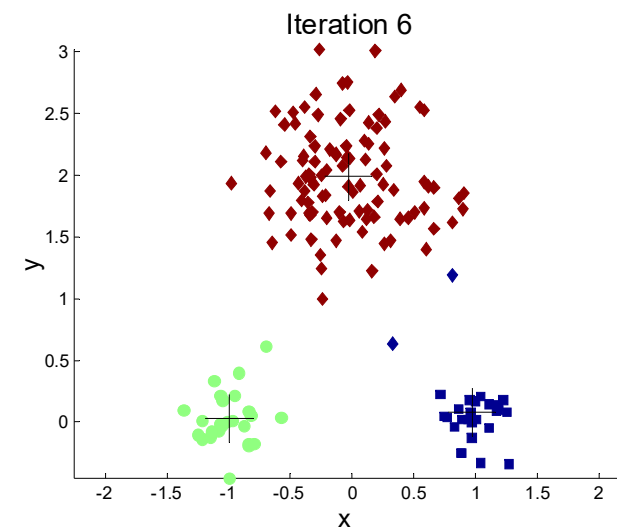
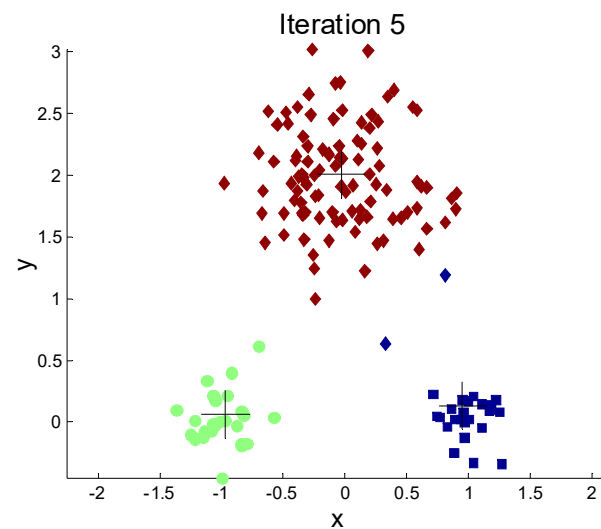
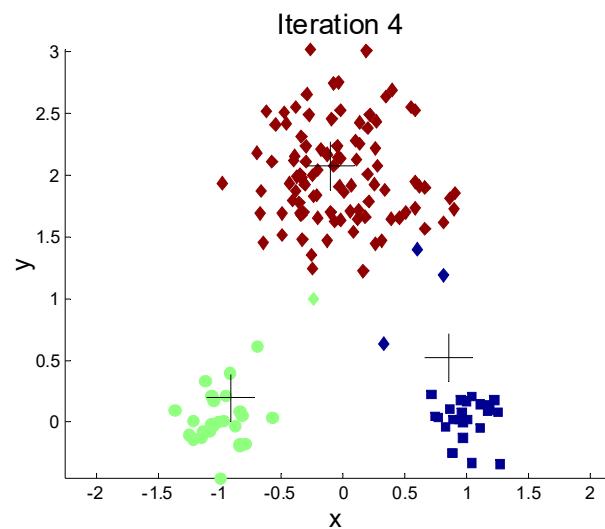
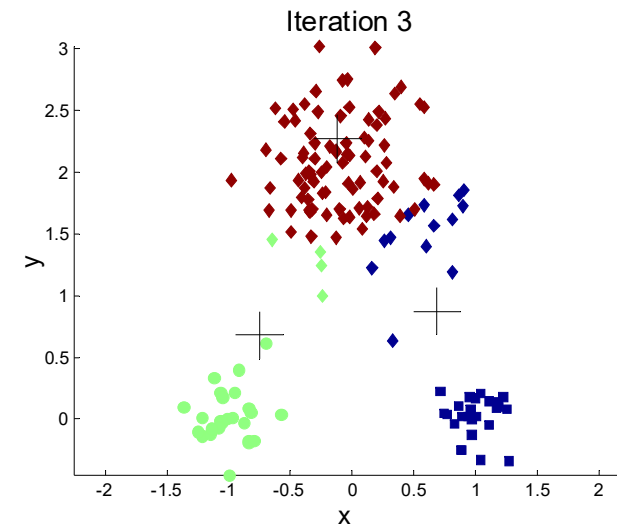
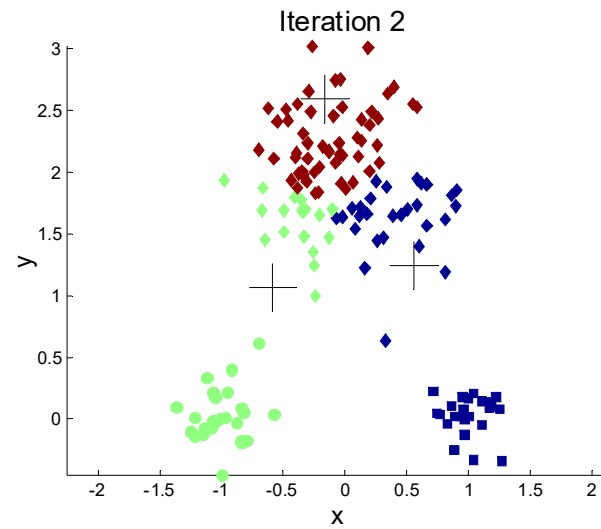
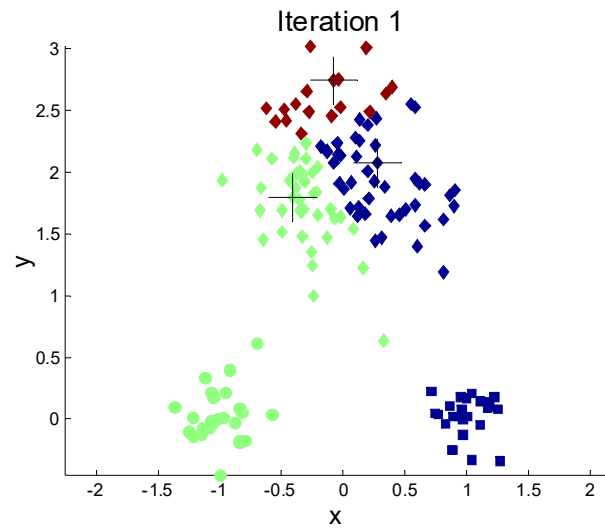


K-Means Example, Step 5

1. Re-compute cluster means
2. Move centroids to new cluster means



K-Means Clustering – Second Example



Convergence Criteria

Standard Convergence Criterion

1. no (or minimum) change of centroids

Alternative Convergence Criteria

1. no (or minimum) re-assignments of data points to different clusters
2. stop after X iterations
3. minimum decrease in the sum of squared error (SSE)
 - see next slide

Evaluating K-Means Clusterings

- Most common cohesion measure: **Sum of Squared Error (SSE)**

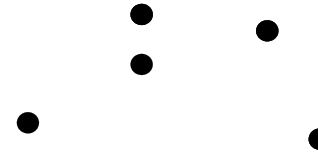
- For each point, the error is the distance to the nearest centroid.
- To get SSE, we square these errors and sum them.

$$SSE = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \text{dist}(\mathbf{x}, \mathbf{m}_j)^2$$

- C_j is the j -th cluster
 - \mathbf{m}_j is the centroid of cluster C_j (the mean vector of all the data points in C_j),
 - $\text{dist}(\mathbf{x}, \mathbf{m}_j)$ is the distance between data point \mathbf{x} and centroid \mathbf{m}_j .
- Given several clusterings, we should prefer the one with the smallest SSE.

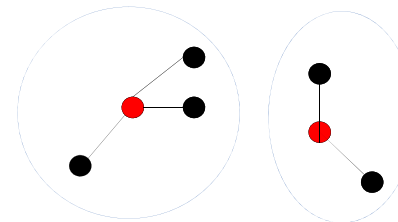
Illustration: Sum of Squared Error

- Clustering problem given:



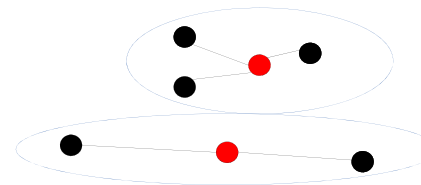
- Good solution:

- small distances to centroids



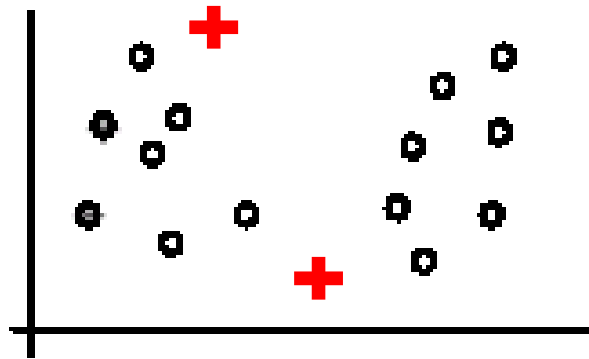
- Not so good solution:

- larger distances to centroids

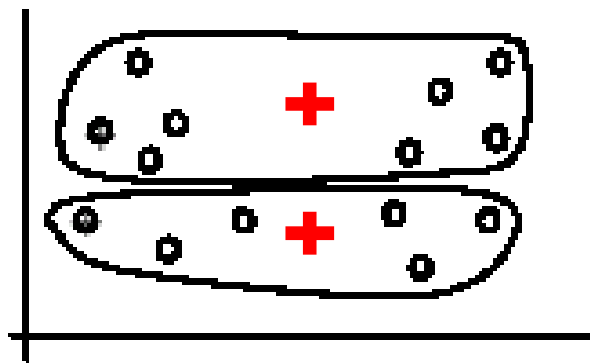


Weaknesses of K-Means: Initial Seeds

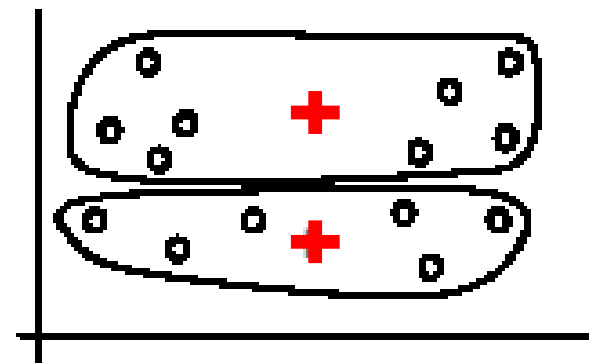
Clustering results may vary significantly depending on initial choice of seeds (**number** and **position** of seeds).



(A). Random selection of seeds (centroids)



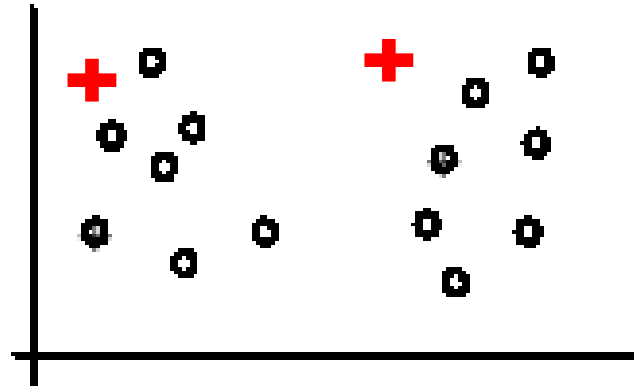
(B). Iteration 1



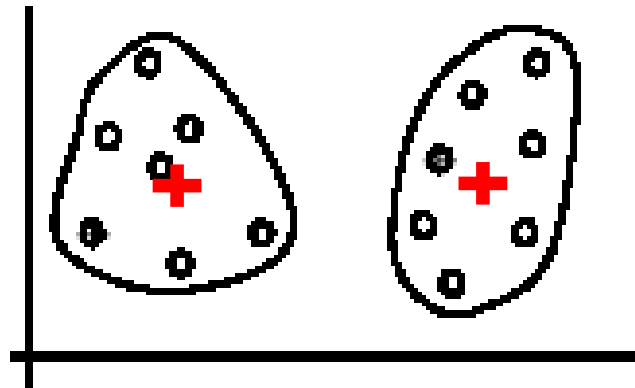
(C). Iteration 2

Weaknesses of K-Means: Initial Seeds

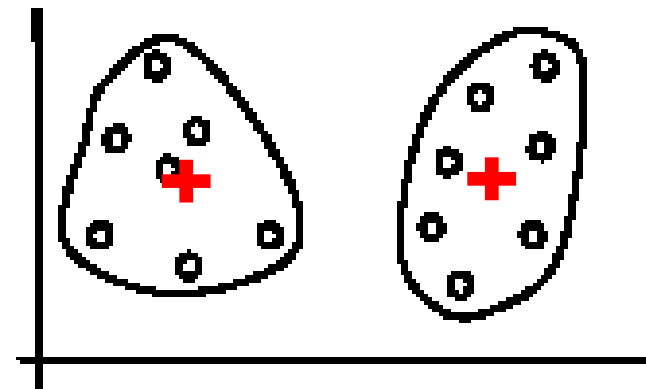
If we use **different seeds**, we get good results.



(A). Random selection of k seeds (centroids)

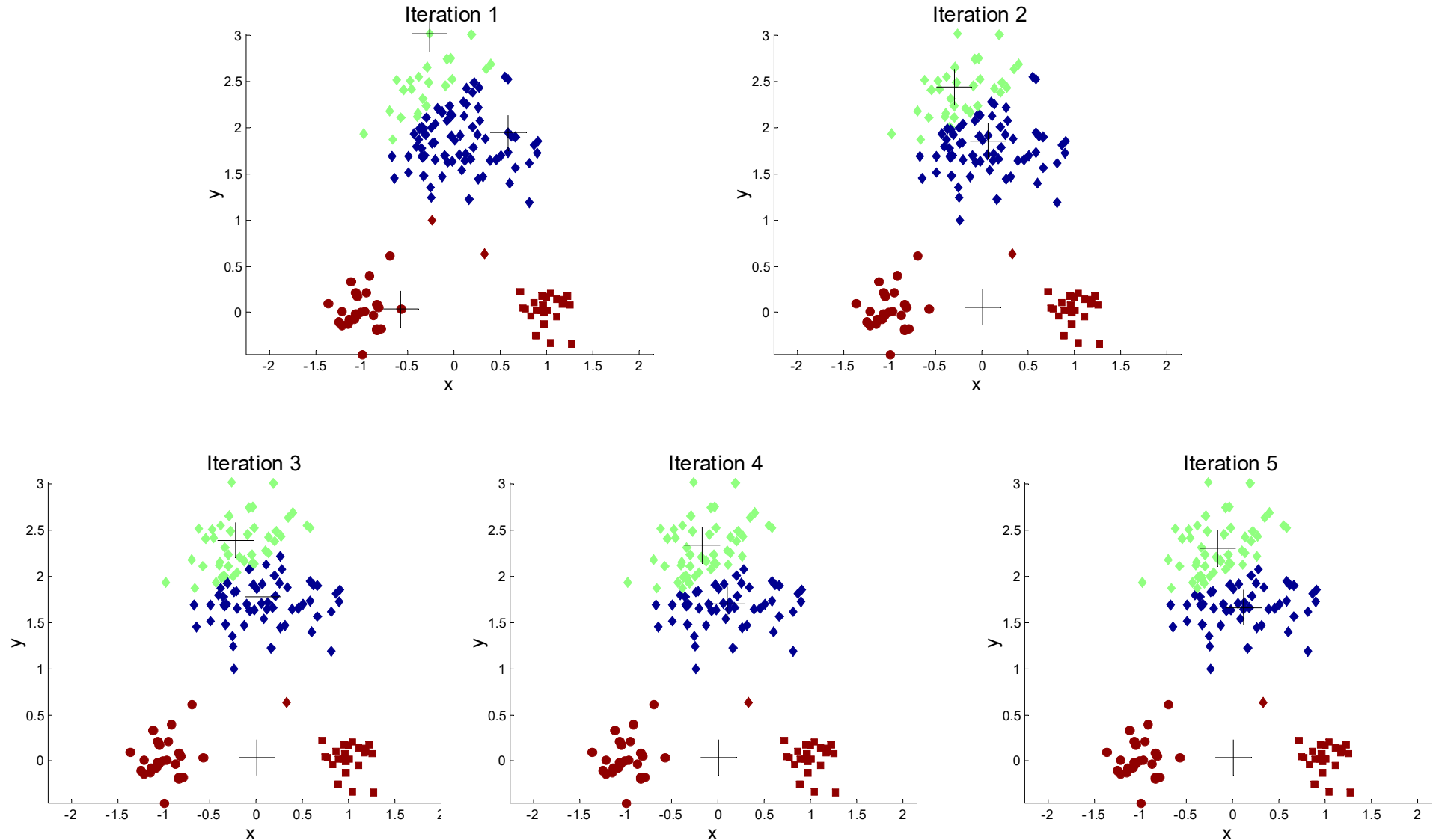


(B). Iteration 1



(C). Iteration 2

Bad Initial Seeds – Second Example



Weaknesses of K-Means: Initial Seeds

Approaches to increase the chance of finding good clusters:

1. Restart a number of times with different random seeds

- chose the resulting clustering with the smallest sum of squared error (SSE)

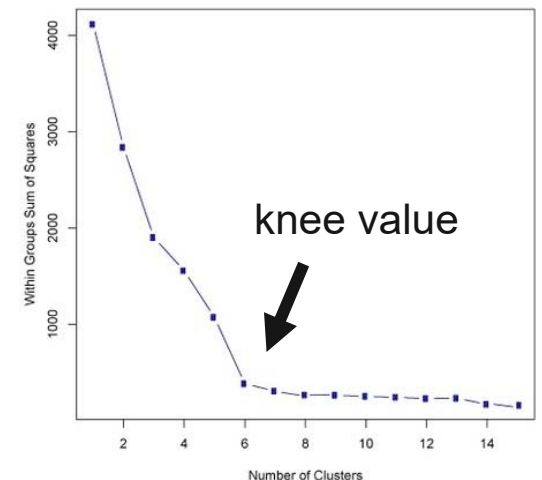
2. Run k-means with different values of k

- The SSE for different values of k cannot directly be compared.
- Think: What happens for $k \rightarrow$ number of examples?
- Workarounds

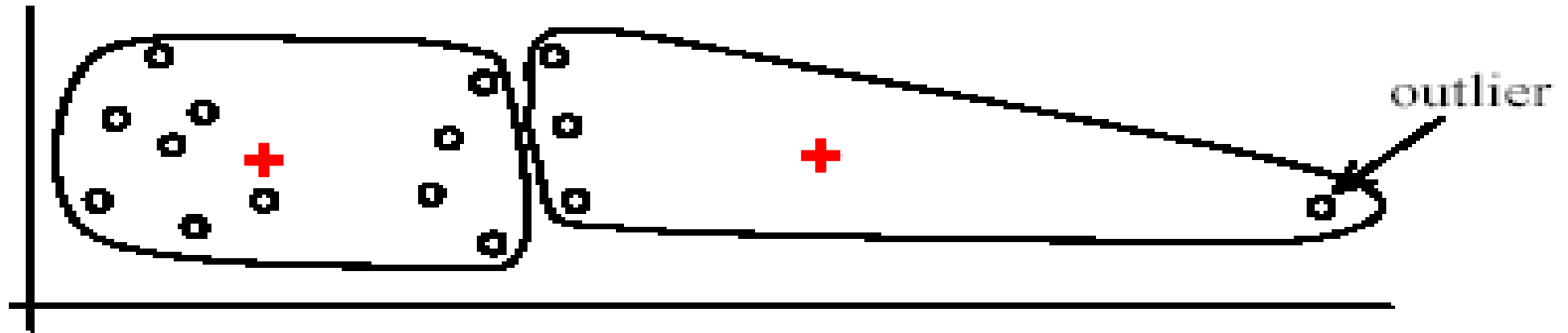
1. Choose k where SSE improvement decreases (knee value of k)

2. Employ X-Means

- Variation of K-Means algorithm that automatically determines k
- starts with small k, then splits large clusters and checks if result improves



Weaknesses of K-Means: Problems with Outliers



(A): Undesirable clusters



(B): Ideal clusters

Weaknesses of K-Means: Problems with Outliers

Approaches to deal with outliers:

1. K-Medoids

- K-Medoids is a K-Means variation that uses the **median** of each cluster instead of the mean.
- Medoids are the **most central existing data points** in each cluster.
- K-Medoids is more robust against outliers as the median is less affected by extreme values:
 - Mean and Median of 1, 3, 5, 7, 9 is **5**
 - Mean of 1, 3, 5, 7, 1009 is **205**
 - Median of 1, 3, 5, 7, 1009 is **5**

2. DBSCAN

- Density-based clustering method that **removes outliers**.
 - see next section

K-Means Clustering Summary

Advantages

- Simple, understandable
- Efficient time complexity: $O(t \cdot k \cdot n)$
where
 - n is the number of data points
 - k is the number of clusters
 - t is the number of iterations

Disadvantages

- Need to determine number of clusters
- All items are forced into a cluster
- Sensitive to outliers
- Sensitive to initial seeds

K-Means Clustering in RapidMiner

The screenshot displays the RapidMiner software interface. On the left, a workflow is visible with two main components: 'Retrieve Data' and 'Clustering'. The 'Retrieve Data' component is a purple box with a green checkmark and an 'out' port. The 'Clustering' component is a green box with a green checkmark, an 'exa' port, and two 'clu' ports. A purple line connects the 'out' port of 'Retrieve Data' to the 'exa' port of 'Clustering'. The 'Clustering' component has three output ports labeled 'res', 'res', and 'res'. On the right, the 'Parameters' panel for 'Clustering (k-Means)' is open. It contains several settings:

- ☒ add cluster attribute
- ☐ add as label
- ☐ remove unlabeled
- k: 4
- max runs: 10
- ☒ determine good start values
- measure types: MixedMeasures
- mixed measure: MixedEuclideanDis...
- max optimization steps: 100

K-Means Clustering Results

Result History | ExampleSet (Clustering) | Cluster Model (Clustering)

Open in: Turbo Prep | Auto Model

Filter (150 / 150 examples):

Row No.	id	label	cluster	a1	a2	a3	
1	id_1	Iris-setosa	cluster_1	6.100	5.500	1.400	0.200
2	id_2	Iris-setosa	cluster_1	4.900	3	1.400	0.200
3	id_3	Iris-se					
4	id_4	Iris-se					
5	id_5	Iris-se					
6	id_6	Iris-se					
7	id_7	Iris-se					
8	id_8	Iris-se					
9	id_9	Iris-se					
10	id_10	Iris-se					
11	id_11	Iris-se					
12	id_12	Iris-se					
13	id_13	Iris-se					
14	id_14	Iris-se					

New cluster attribute

Result History | ExampleSet (Clustering) | Cluster Model (Clustering)

Description

Attribute	cluster_0	cluster_1	cluster_2
a1	6.913	5.006	6.252
a2	3.100	3.418	2.855
a3			
a4			

Folder View

Graph

Centroid Table

Result History | ExampleSet (Clustering) | Cluster Model (Clustering)

Cluster Model

Description

Cluster 0: 32 items
Cluster 1: 50 items
Cluster 2: 40 items
Cluster 3: 28 items
Total number of items: 150

Folder View

Graph

Centroid Table

K-Medoids Clustering in RapidMiner

The screenshot displays the RapidMiner software interface. On the left, a workflow is visible with two main components: 'Retrieve Data' and 'Clustering'. The 'Retrieve Data' component has an 'out' port connected to the 'Clustering' component's 'exa' port. The 'Clustering' component has two 'clu' ports, each connected to a 'res' port. The 'Clustering' component is highlighted with an orange border. On the right, the 'Parameters' panel for 'Clustering (k-Medoids)' is open. It contains several settings:

- ☒ add cluster attribute
- ☐ add as label
- ☐ remove unlabeled
- k: 2
- max runs: 10
- max optimization steps: 100
- ☐ use local random seed
- measure types: MixedMeasures
- mixed measure: MixedEuclideanDis...

X-Means Clustering in RapidMiner

The screenshot displays the X-Means clustering tool in RapidMiner. On the left, the main interface shows a workflow with an 'X-Means' node. The node has an input port labeled 'exa' and two output ports labeled 'clu'. The 'clu' output is connected to a 'res' port. The 'Parameters' dialog box is open on the right, showing the following settings:

- X-Means**
- ☒ add cluster attribute
- ☐ add as label
- ☐ remove unlabeled
- k min: 3
- k max: 60
- ☒ determine good start values
- measure types: NumericalMeasures
- numerical measure: EuclideanDistance
- clustering algorithm: KMeans

Boundaries
for testing
k values

3. Density-based Clustering



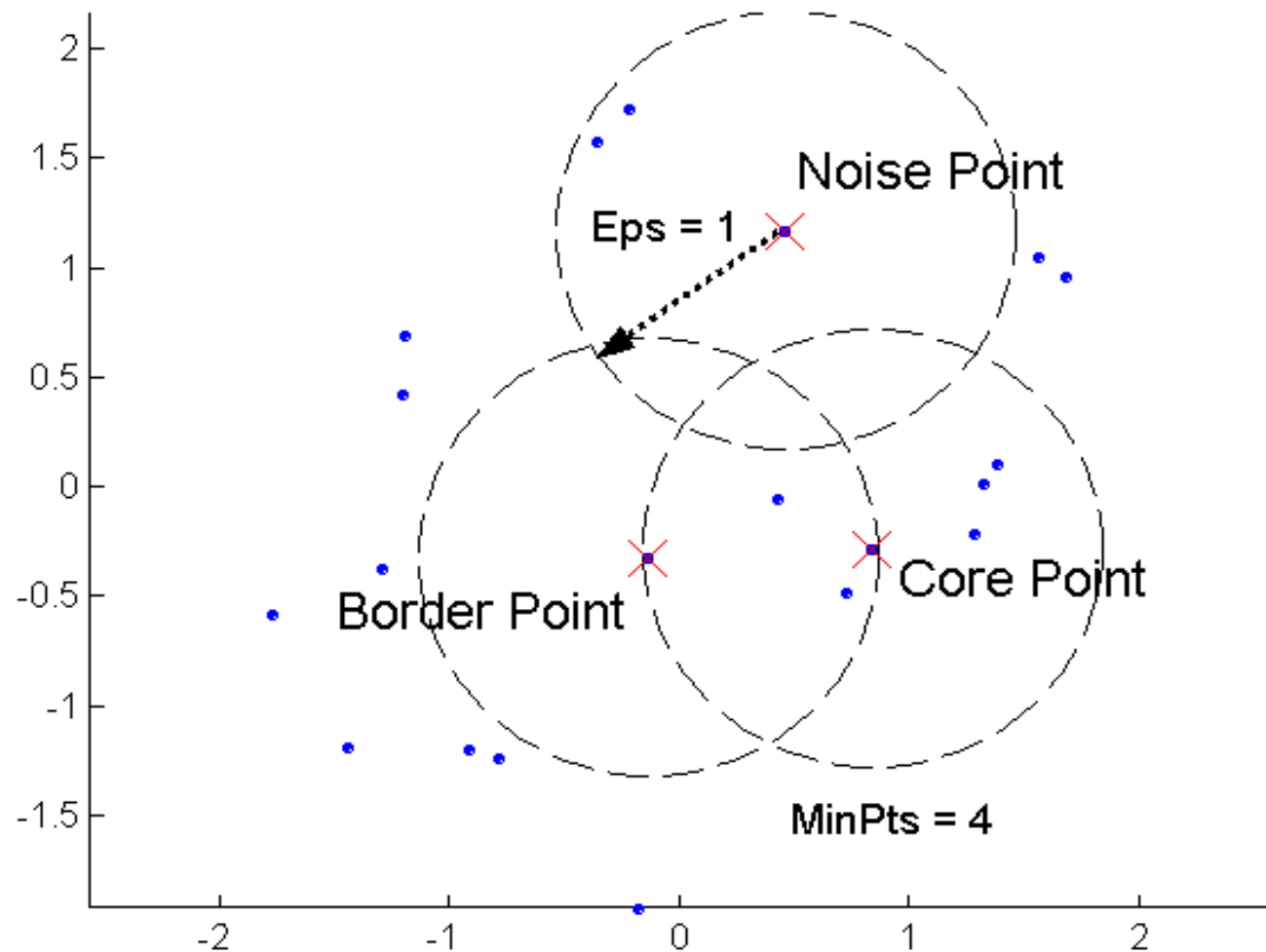
Challenging use case for K-Means:

- Problem 1: Non-globular shapes
- Problem 2: Outliers / Noise points

DBSCAN

- DBSCAN is a density-based algorithm
 - **Density** = number of points within a specified radius Epsilon (Eps)
- Divides data points in three classes:
 1. A point is a **core point** if it has at least a specified number of neighboring points (MinPts) within the specified radius Eps
 - The point itself is counted as well
 - These points form the interior of a dense region (cluster)
 2. A **border point** has fewer than MinPts within Eps, but is in the neighborhood of a core point
 3. A **noise point** is any point that is not a core point or a border point

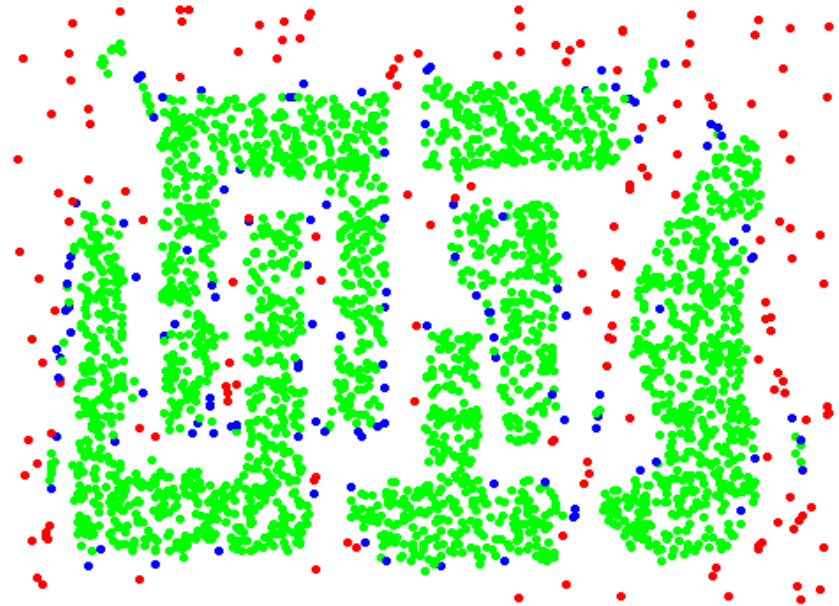
Examples of Core, Border, and Noise Points 1



Examples of Core, Border, and Noise Points 2



Original Points



Point types: **core**,
border and **noise**

The DBSCAN Algorithm

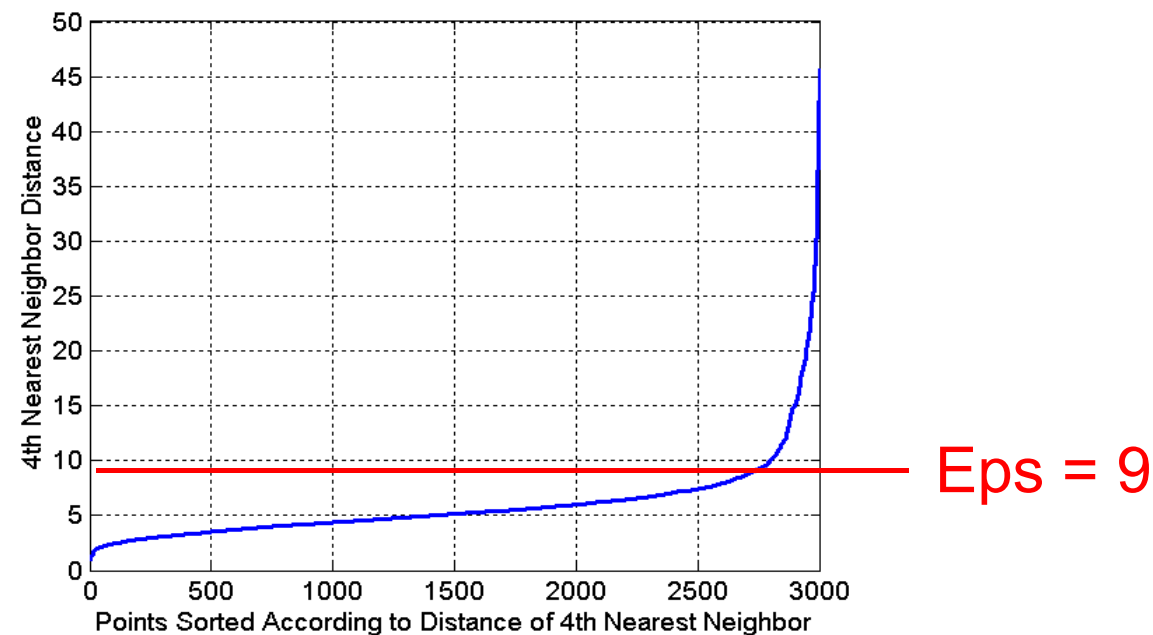
Eliminates noise points and returns clustering of the remaining points:

1. Label all points as core, border, or noise points
2. Eliminate all noise points
3. Put an edge between all core points that are within Eps of each other
4. Make each group of connected core points into a separate cluster
5. Assign each border point to one of the clusters of its associated core points
 - as a border point can be at the border of multiple clusters
 - use voting if core points belong to different clusters.
 - if equal vote, than assign border point randomly

How to Determine suitable Eps and MinPts values?

For points in a cluster, their k^{th} nearest neighbor (single point) should be at roughly the same distance. Noise points have their k^{th} nearest neighbor at farther distance.

1. Start with setting MinPts = 4 (rule of thumb)
2. Plot sorted distance of every point to its k^{th} nearest neighbor:

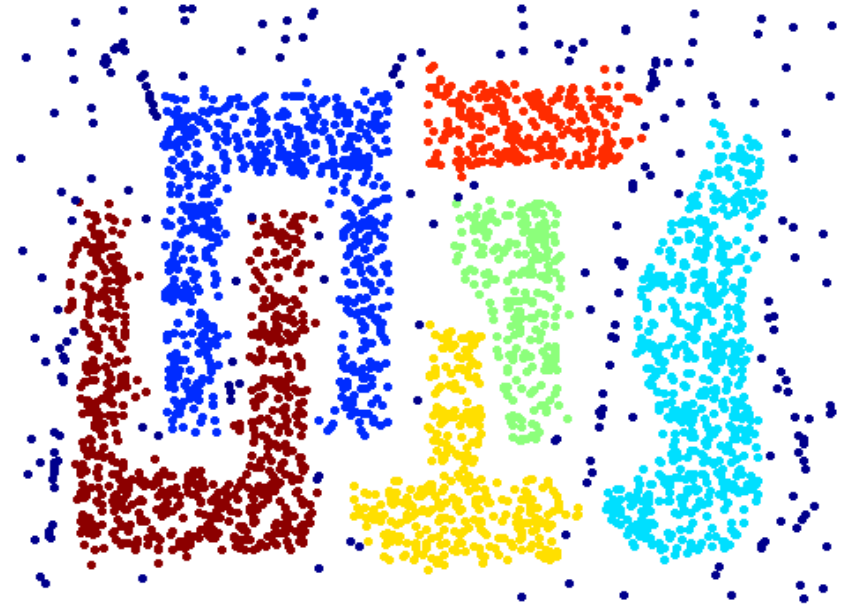


3. Set Eps to the sharp increase of the distances (start of noise points)
4. Decrease k if small clusters are labeled as noise (subjective decision)
5. Increase k if outliers are included into the clusters (subjective decision)

When DBSCAN Works Well



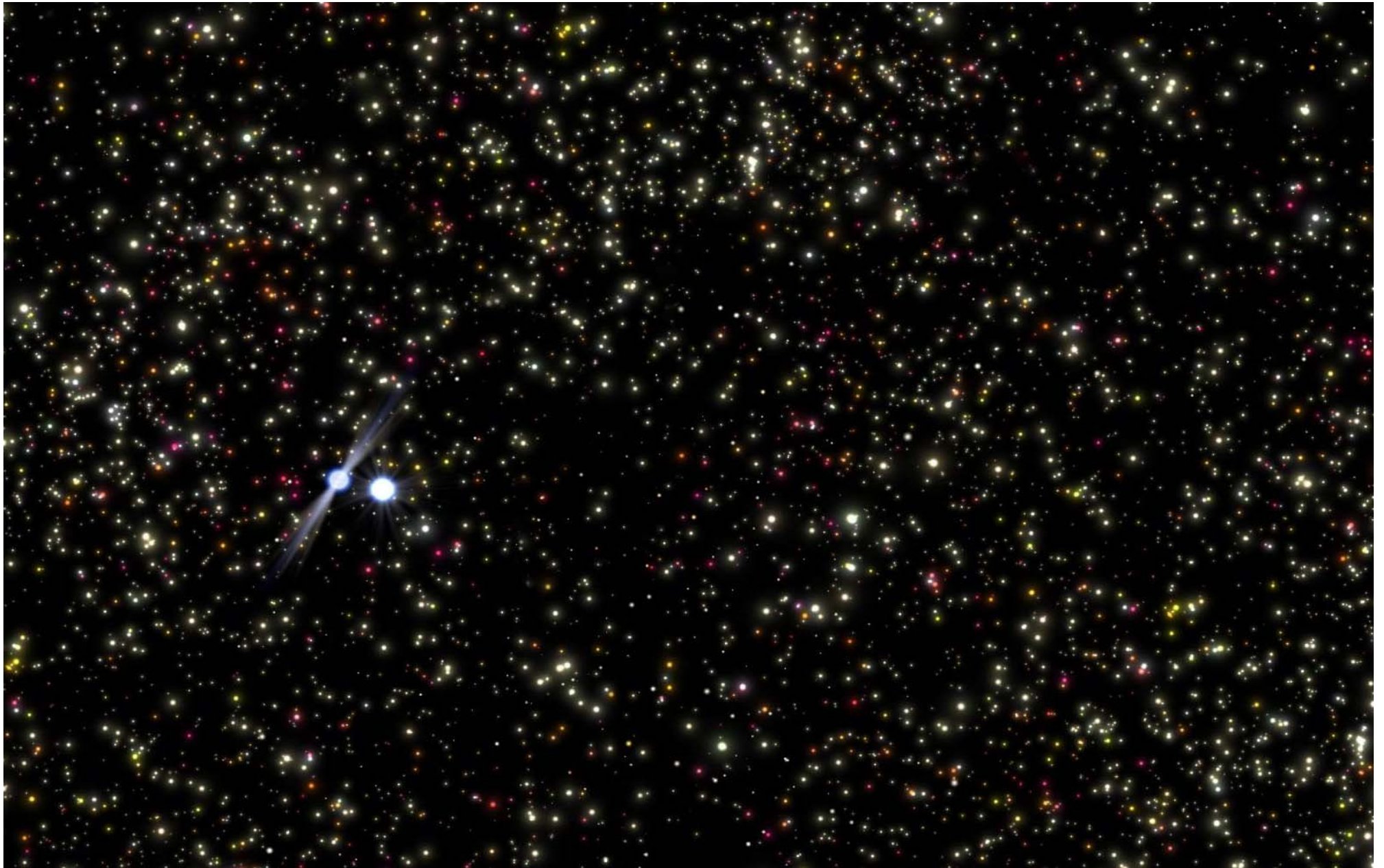
Original Points



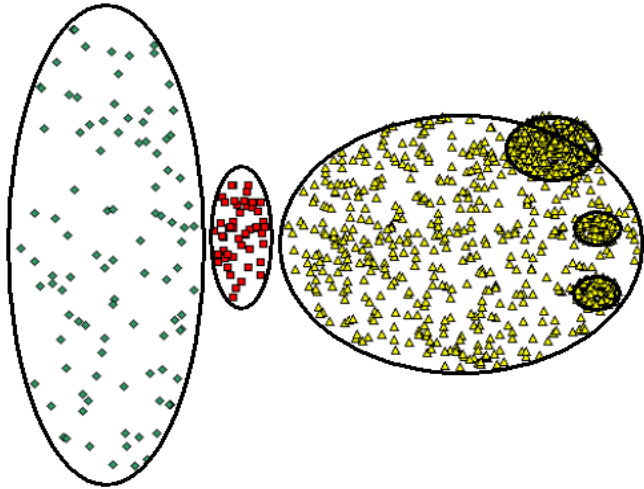
Clusters

- Resistant to noise
- Can handle clusters of different shapes and sizes

Application: Sky Images

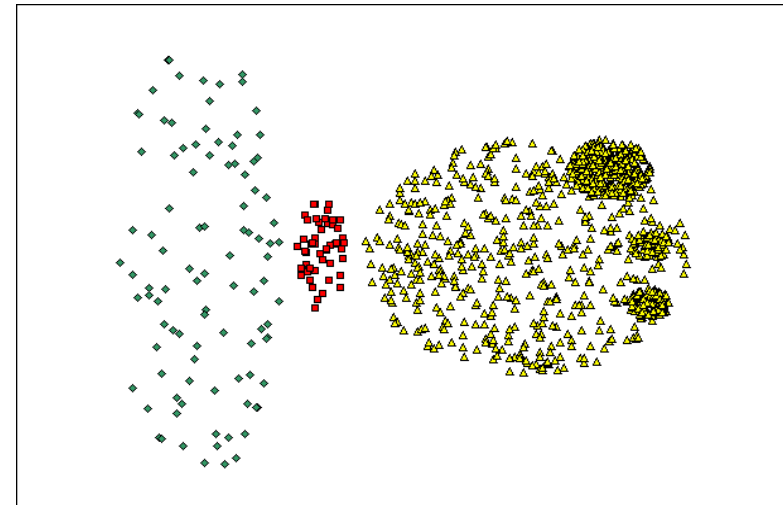


When DBSCAN Does NOT Work Well

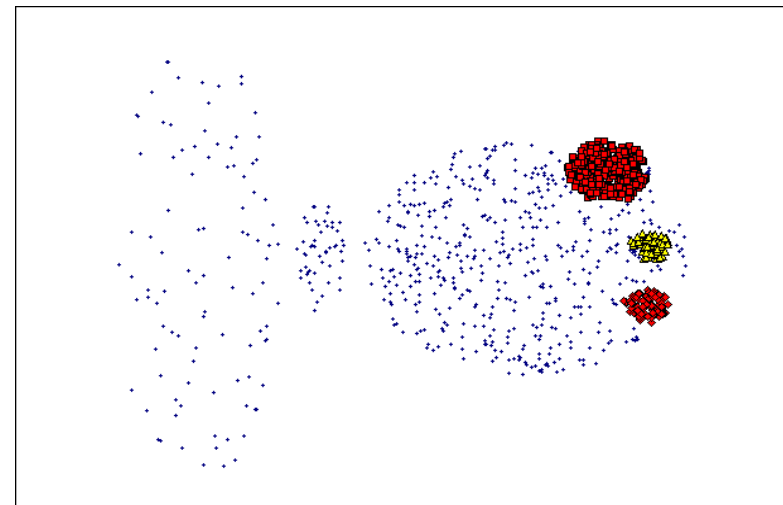


Original Points

DBSCAN has problems with datasets of varying densities.



(MinPts=4, Eps=9.92)



(MinPts=4, Eps=9.75)

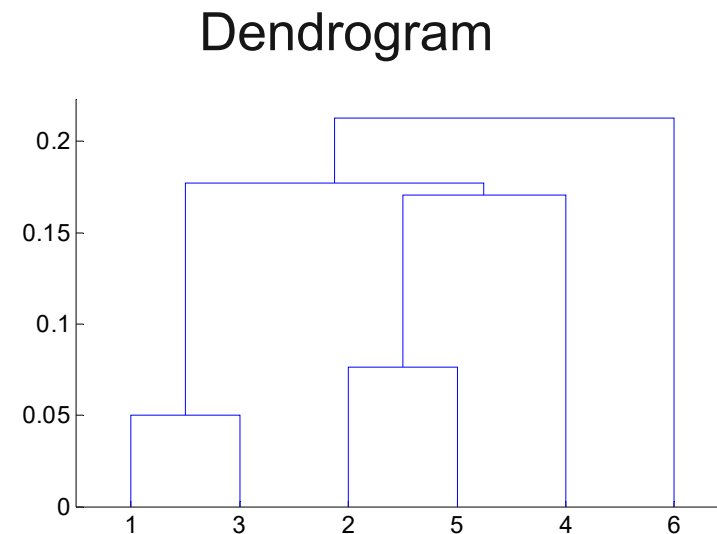
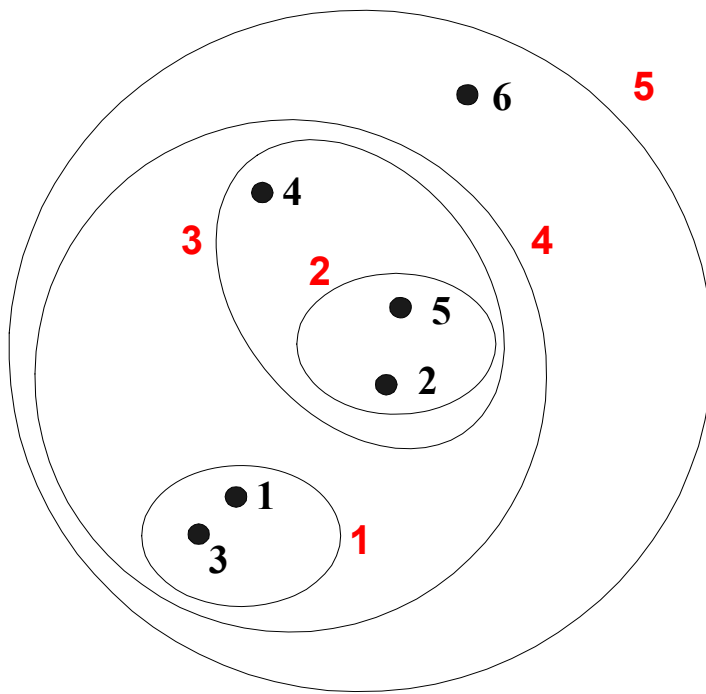
DBSCAN in RapidMiner

The screenshot displays the RapidMiner software interface. On the left, a workflow canvas shows a 'Clustering' process node with an input port labeled 'exa' and two output ports labeled 'clu'. A data stream is connected to the input, and two output streams are shown. The right-hand pane is titled 'Parameters' and contains the configuration for the 'Clustering (DBSCAN)' process. The parameters are as follows:

- epsilon**: 1.0
- min points**: 5
- add cluster attribute**: ☒
- add as label**: ☐
- remove unlabeled**: ☐
- measure types**: MixedMeasures
- mixed measure**: MixedEuclideanDistance

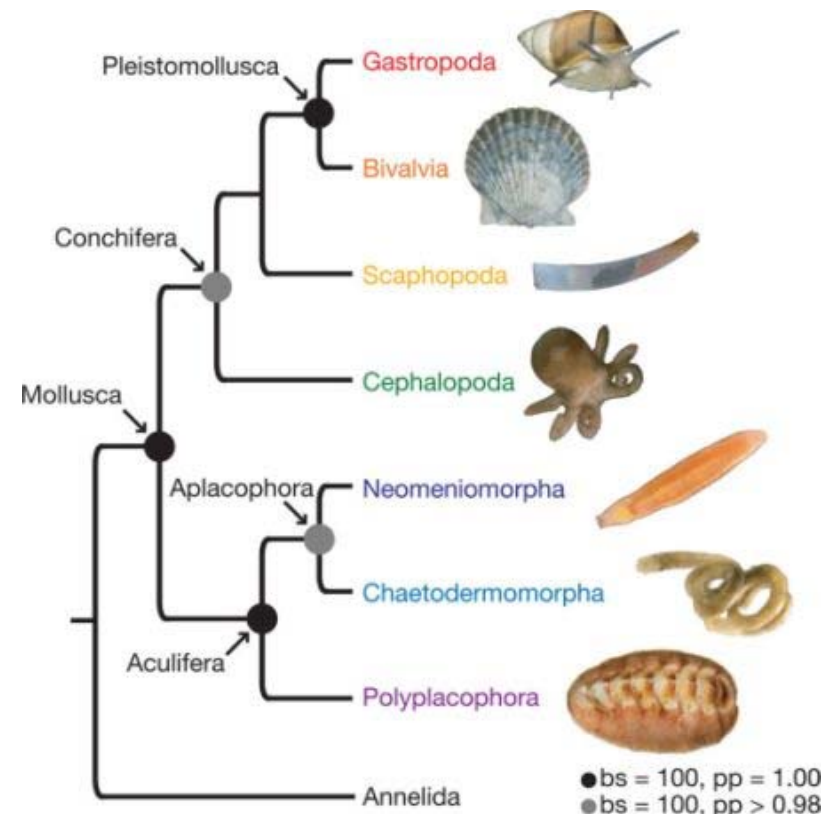
4. Hierarchical Clustering

- Produces a set of **nested clusters** organized as a hierarchical tree.
- Can be visualized as a Dendrogram
 - A tree like diagram that records the sequences of merges or splits.
 - The y-axis displays the former distance between merged clusters.



Strengths of Hierarchical Clustering

- We do not have to assume any particular number of clusters
 - any desired number of clusters can be obtained by ‘cutting’ the dendrogram at the proper level.
- May be used to discover meaningful taxonomies
 - taxonomies of biological species
 - taxonomies of different customer groups



Two Main Types of Hierarchical Clustering

– Agglomerative

- Start with the points as individual clusters
- At each step, merge the closest pair of clusters until only one cluster (or k clusters) is left

– Divisive

- Start with one, all-inclusive cluster
- At each step, split a cluster until each cluster contains a point (or there are k clusters)

– Agglomerative Clustering is more widely used.

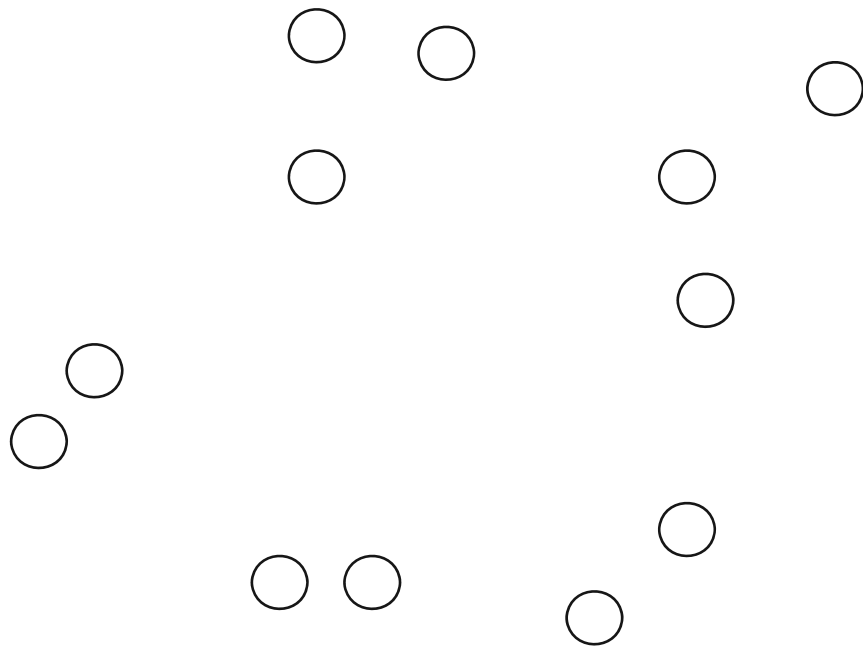
Agglomerative Clustering Algorithm

The basic algorithm is straightforward:

1. Compute the proximity matrix
 2. Let each data point be a cluster
 3. Repeat
 1. Merge the two closest clusters
 2. Update the proximity matrixUntil only a single cluster remains
- The key operation is the computation of the proximity of two clusters.
 - The different approaches to defining the distance between clusters distinguish the different algorithms.

Starting Situation

Start with clusters of individual points and a proximity matrix.



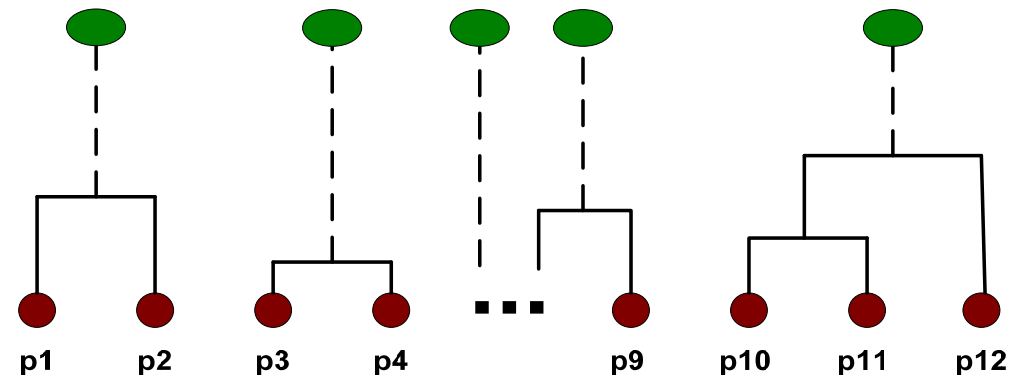
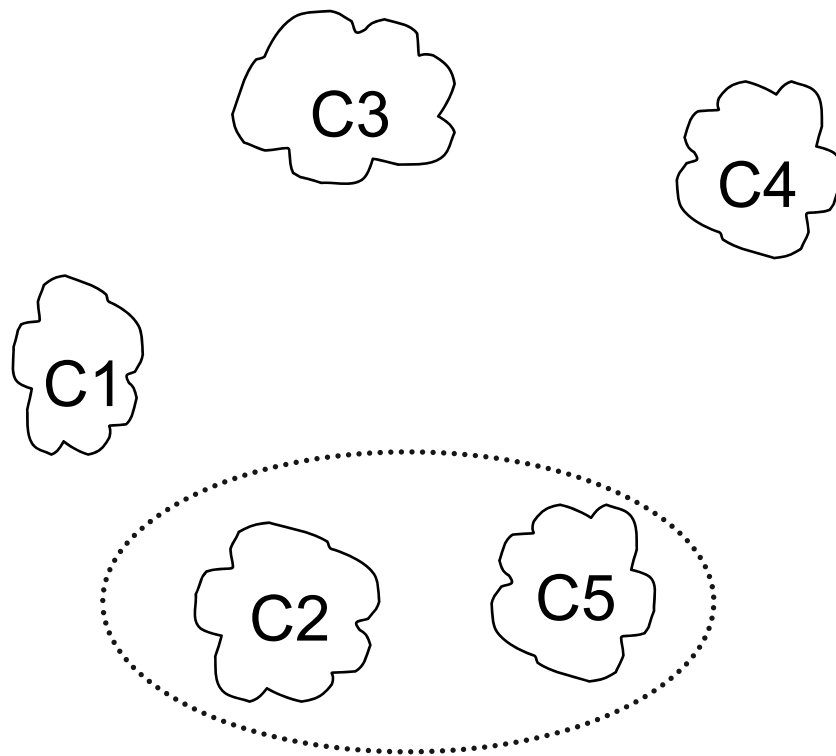
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Proximity Matrix

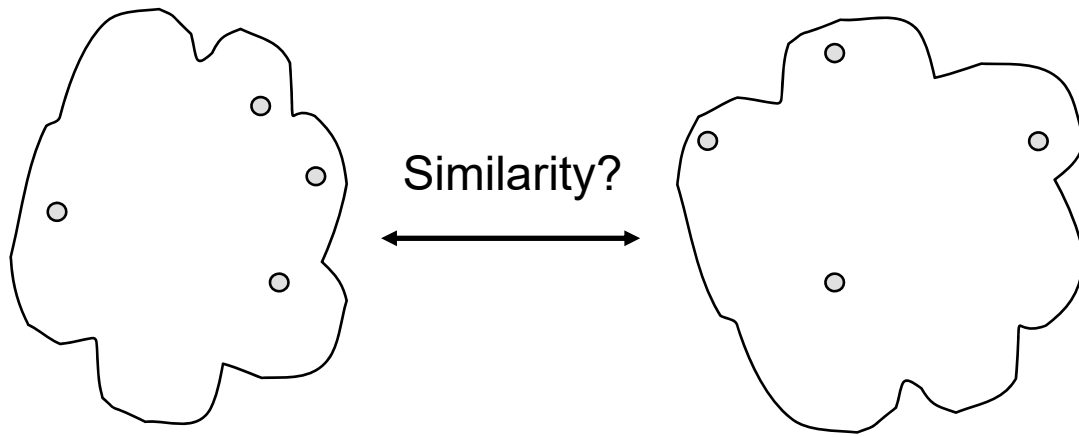


Intermediate Situation

- After some merging steps, we have larger clusters.
- We want to keep on merging the two closest clusters (C2 and C5?)



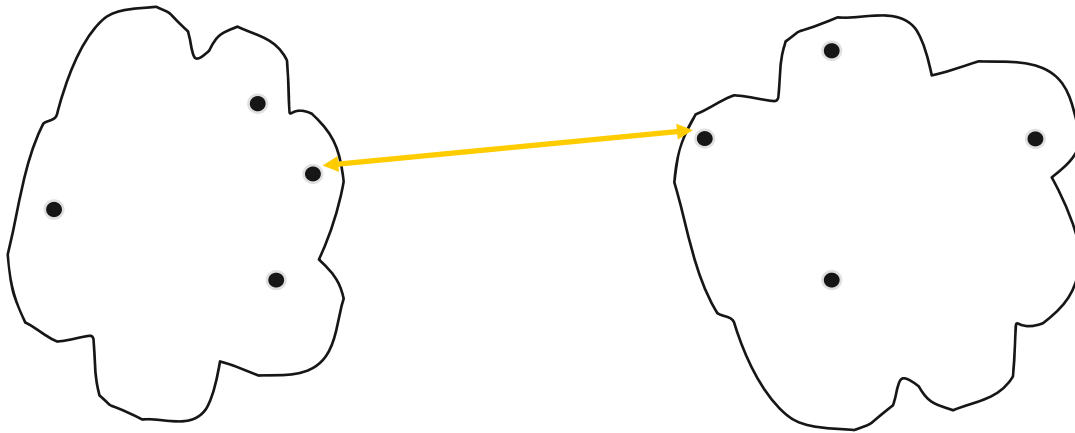
How to Define Inter-Cluster Similarity?



Different Approaches possible:

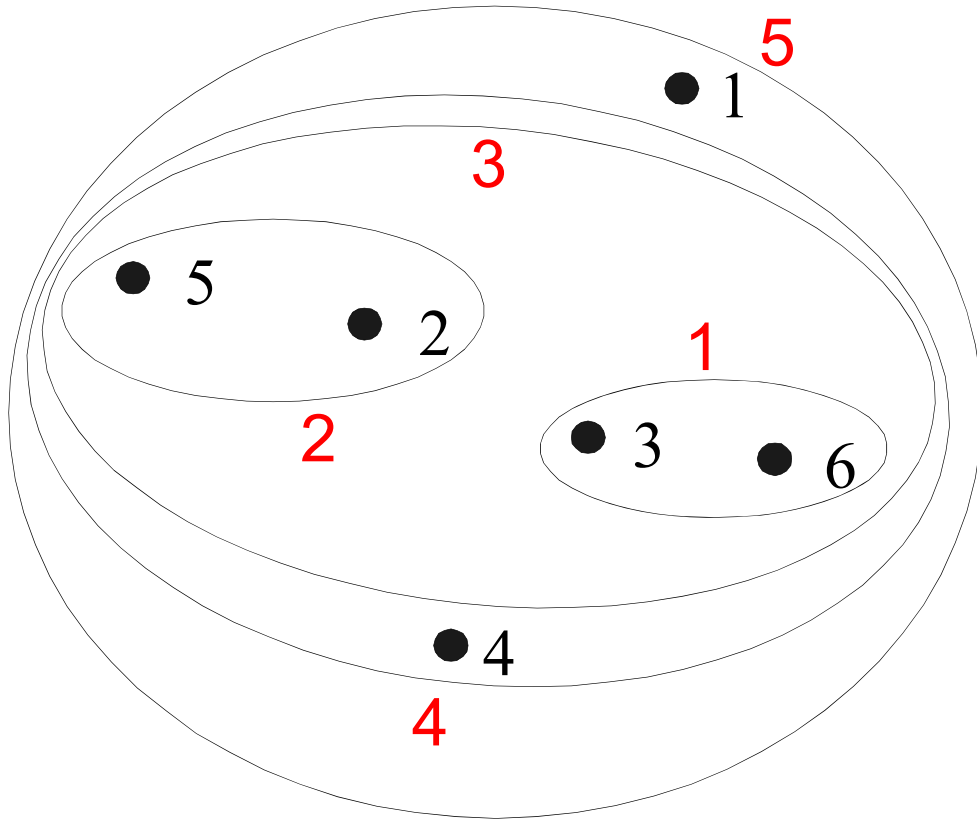
1. Single Link
2. Complete Link
3. Group Average
4. Distance Between Centroids

Cluster Similarity: Single Link

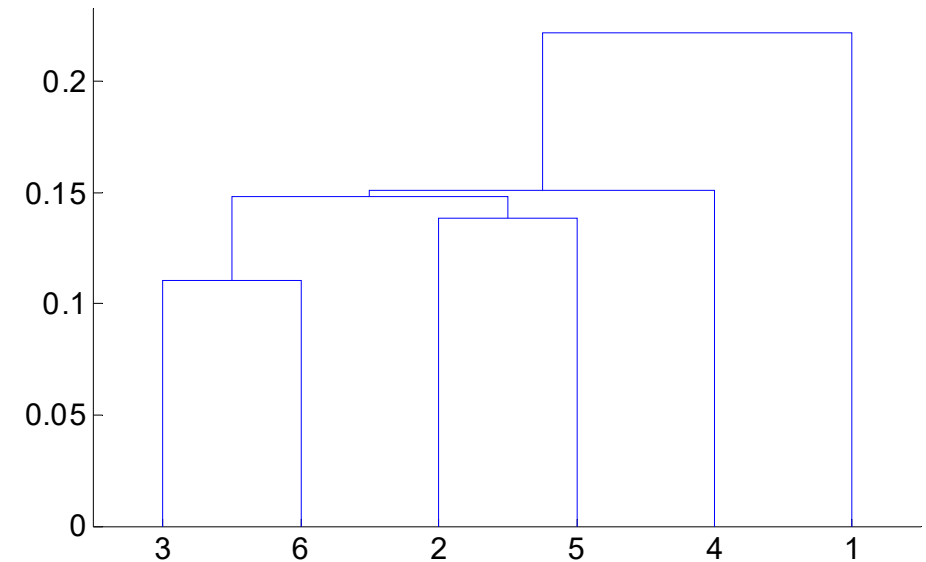


- Similarity of two clusters is based on the **two most similar (closest) points** in the different clusters
- Determined by one pair of points, i.e. by one link in the proximity graph.

Example: Single Link

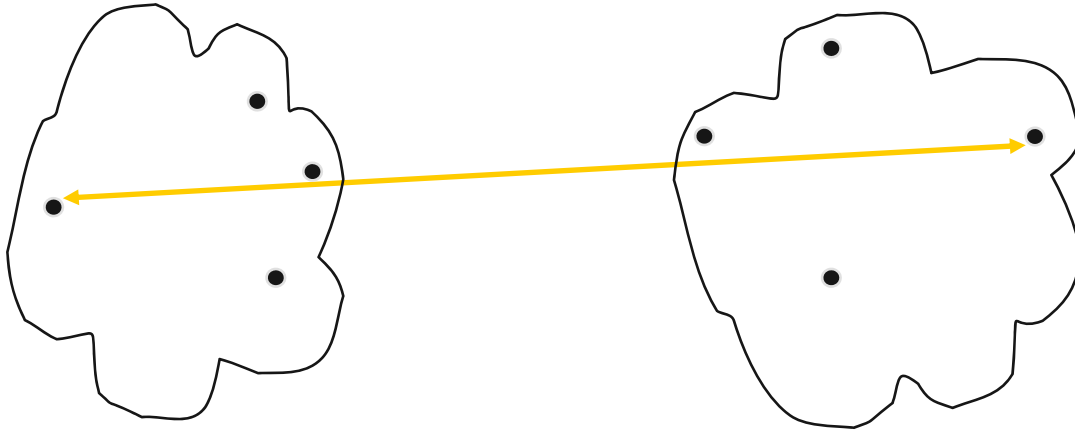


Nested Clusters



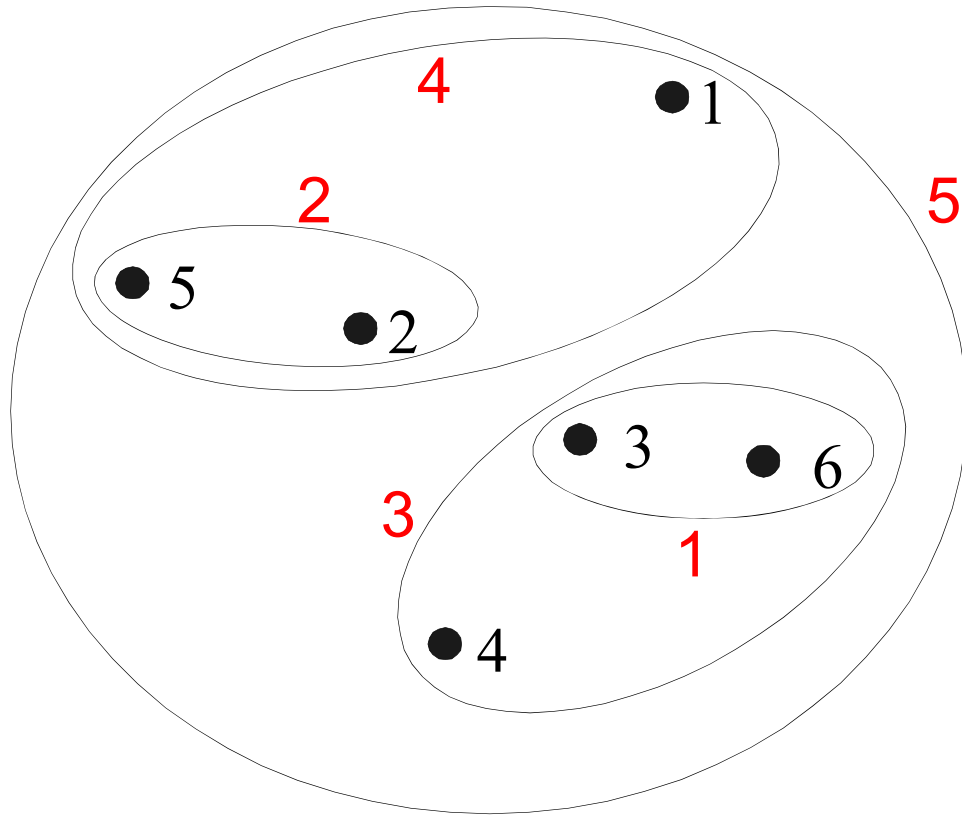
Dendrogram

Cluster Similarity: Complete Linkage

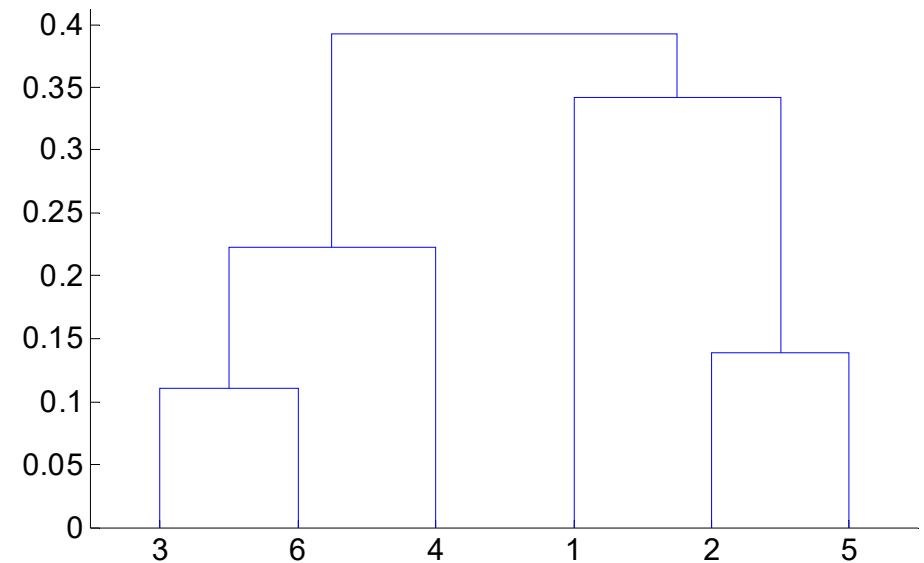


- Similarity of two clusters is based on the **two least similar (most distant) points** in the different clusters
- Determined by all pairs of points in the two clusters

Example: Complete Linkage



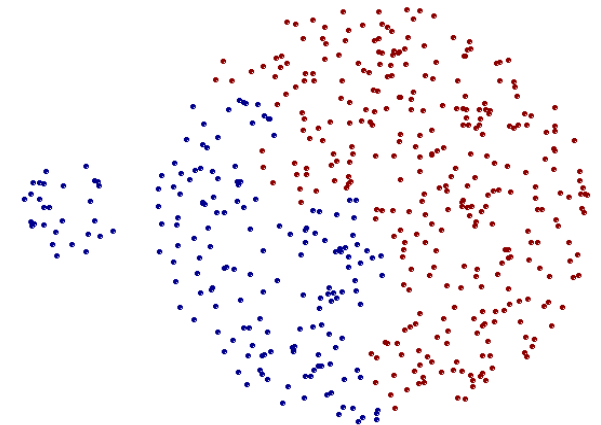
Nested Clusters



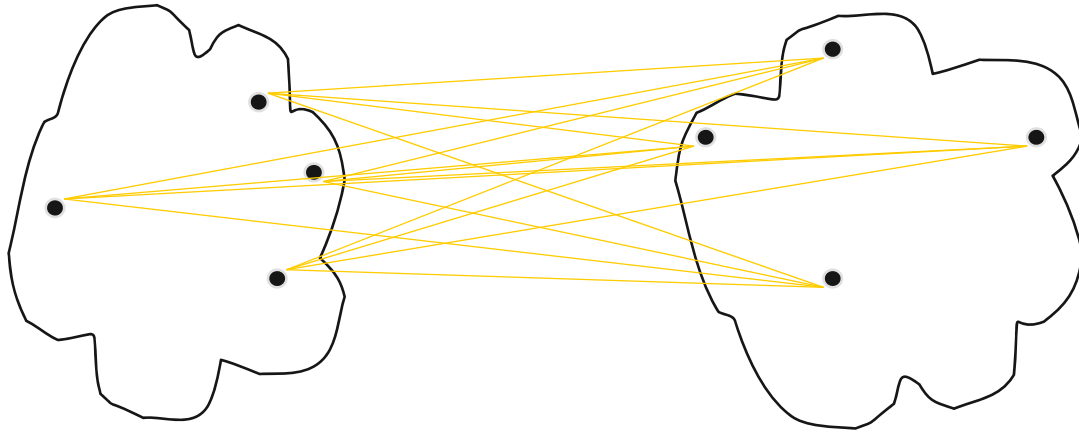
Dendrogram

Single Link vs. Complete Linkage

- Single Link:
 - Strength: Can handle non-elliptic shapes
 - Limitation: Sensitive to noise and outliers
- Complete Linkage:
 - Strength: Less sensitive to noise and outliers
 - Limitation: Biased towards globular clusters
 - Limitation: Tends to break large clusters, as decisions can not be undone.



Cluster Similarity: Group Average

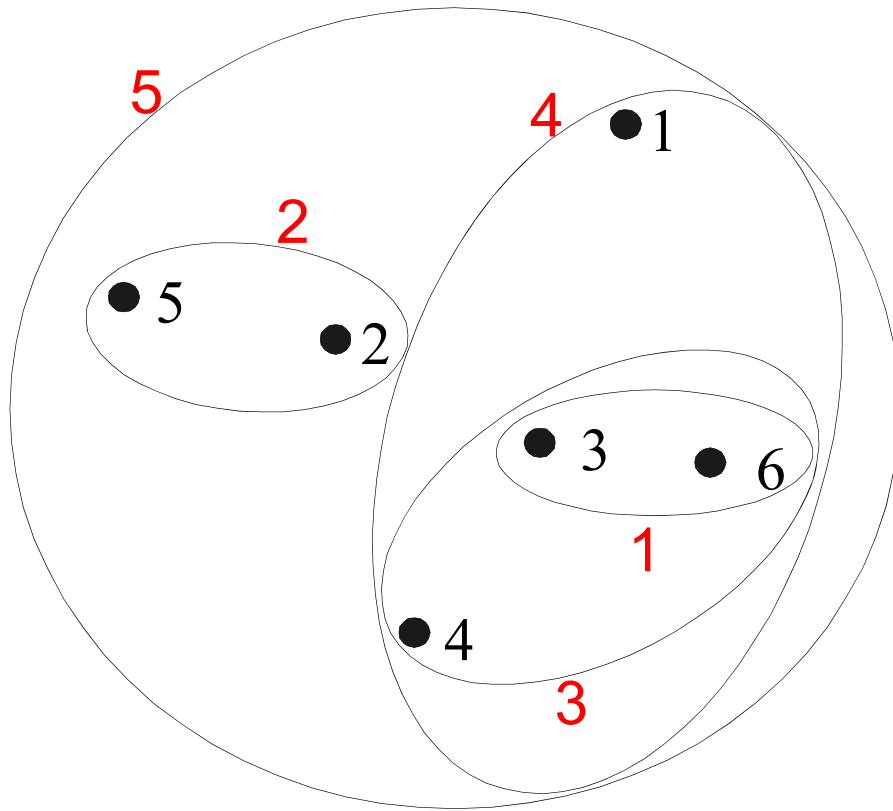


- Proximity of two clusters is the average of pair-wise proximity between all points in the two clusters.

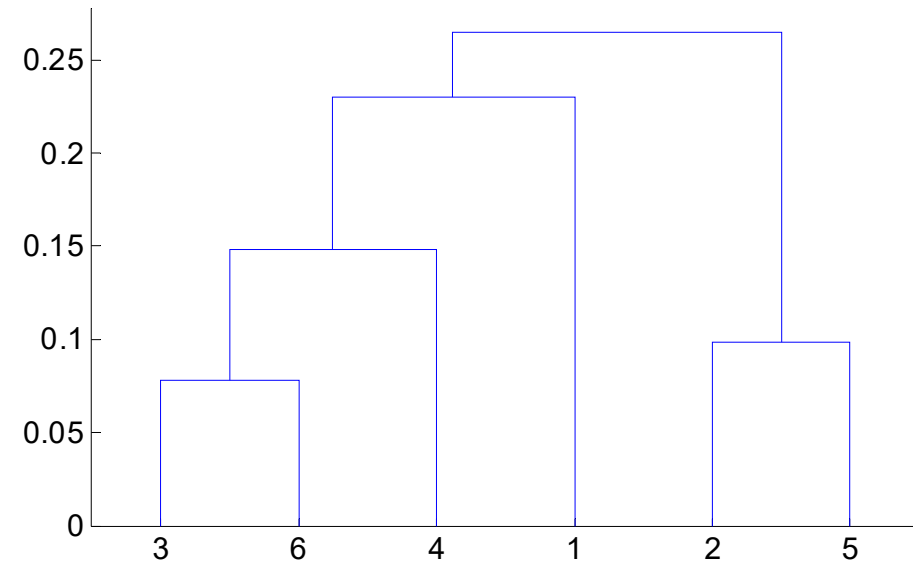
$$\text{proximity}(\text{Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i \in \text{Cluster}_i \\ p_j \in \text{Cluster}_j}} \text{proximity}(p_i, p_j)}{|\text{Cluster}_i| * |\text{Cluster}_j|}$$

- Compromise between single and complete link
 - Strength: Less sensitive to noise and outliers than single link
 - Limitation: Biased towards globular clusters

Example: Group Average



Nested Clusters



Dendrogram

Hierarchical Clustering: Problems and Limitations

- Different schemes have problems with one or more of the following:
 1. Sensitivity to noise and outliers
 2. Difficulty handling non-elliptic shapes
 3. Breaking large clusters
- High Space and Time Complexity
 - $O(N^2)$ space since it uses the proximity matrix.
 - N is the number of points.
 - $O(N^3)$ time in many cases
 - There are N steps and at each step the size N^2 proximity matrix must be updated and searched.
 - Complexity can be reduced to $O(N^2 \log(N))$ time in some cases.
 - Workaround: Apply hierarchical clustering to a random sample of the original data (<10,000 examples).

Agglomerative Clustering in RapidMiner

**Creates
hierarchical clustering**

**Flattens clustering to a given
number of clusters**

The screenshot displays the RapidMiner workflow editor. On the left, a workflow is shown with two main nodes: 'Clustering' and 'Flatten Clustering'. The 'Clustering' node is highlighted with an orange border and a green checkmark. It has an input port labeled 'exa' and two output ports labeled 'clu' and 'exa'. The 'Flatten Clustering' node also has a green checkmark and is connected to the 'Clustering' node. It has an input port labeled 'hie' and two output ports labeled 'fla' and 'exa'. The 'Flatten Clustering' node has three additional output ports labeled 'res'. On the right, the 'Parameters' panel is open, showing the configuration for the 'Clustering (Agglomerative Clustering)' node. The parameters are: 'mode' set to 'SingleLink', 'measure types' set to 'MixedMeasures', and 'mixed measure' set to 'MixedEuclideanDista...'. The top toolbar shows a zoom level of 100% and various icons for zooming and navigating the workflow.

5. Proximity Measures

- So far, we have seen different clustering algorithms all of which rely on proximity (distance, similarity, ...) measures.
- Now, we treat proximity measures in more detail.
- A **wide range of different measures** is used depending on the requirements of the application.
- Similarity
 - Numerical measure of how alike two data objects are.
 - Often falls in the range $[0,1]$
- Dissimilarity
 - Numerical measure of how different are two data objects
 - Minimum dissimilarity is often 0
 - Upper limit varies

5.1 Proximity of Single Attributes

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$, where n is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d = p - q $	$s = -d, s = \frac{1}{1+d} \text{ or } s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

Similarity and dissimilarity for simple attributes

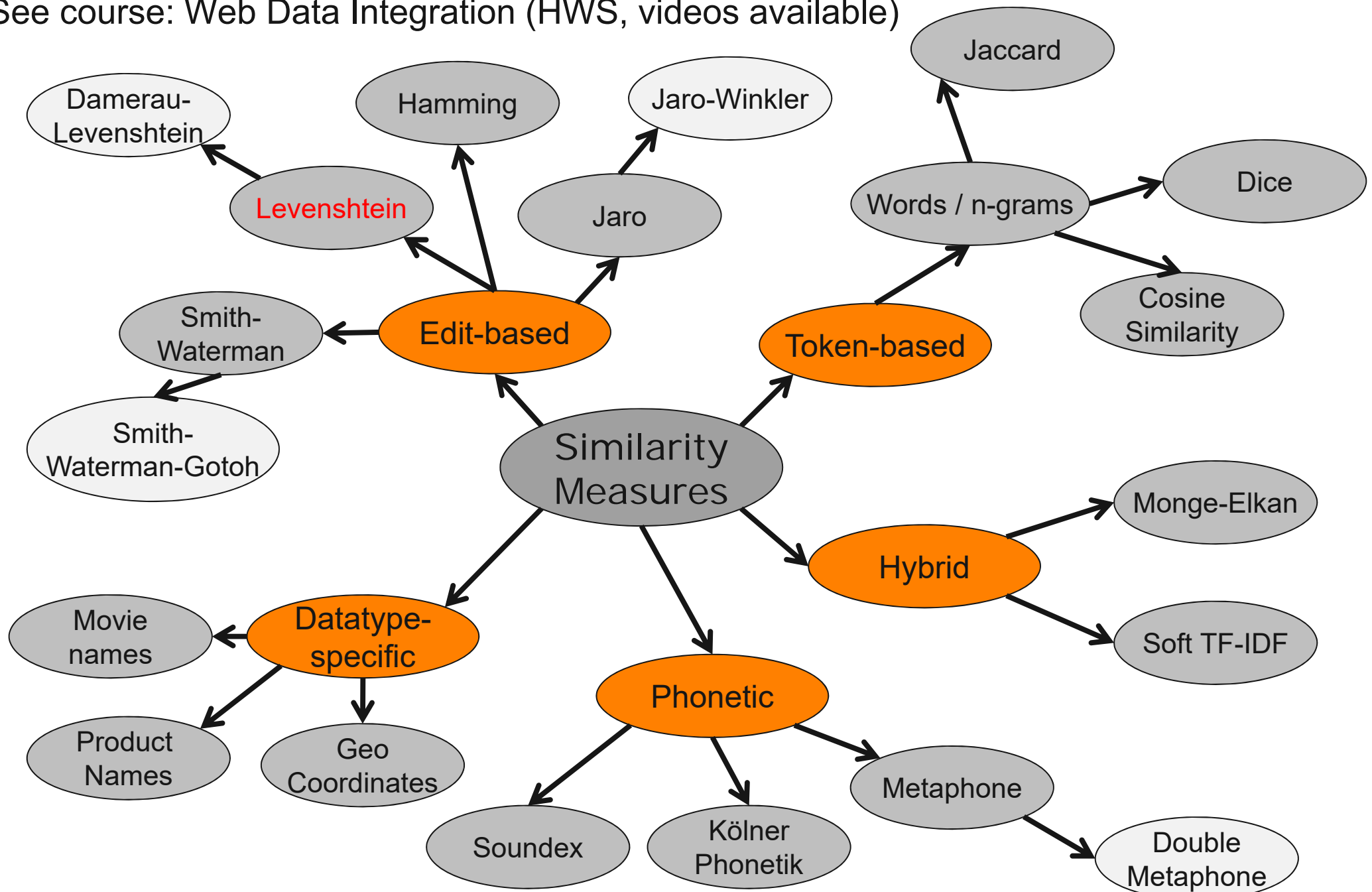
p and q are attribute values for two data objects

Levenshtein Distance

- Measures the dissimilarity of **two strings**.
- Measures the minimum number of edits needed to transform one string into the other.
- Allowed edit operations:
 1. insert a character into the string
 2. delete a character from the string
 3. replace one character with a different character
- Examples:
 - $\text{levensthein}(\text{'Table'}, \text{'Cable'}) = 1$ (1 Substitution)
 - $\text{levensthein}(\text{'Table'}, \text{'able'}) = 1$ (1 Deletion)

Further String Similarity Measures

See course: Web Data Integration (HWS, videos available)



5.2 Proximity of Multidimensional Data Points

- All measures discussed so far cover the proximity of single attribute values
- But we usually have data points with many attributes
 - e.g., age, height, weight, sex...
- Thus, we need proximity measures for data points
 - taking multiple attributes/dimensions into account

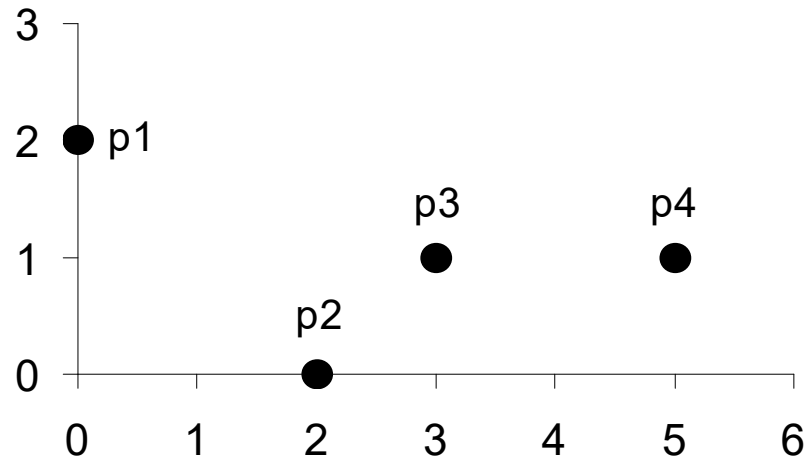
Euclidean Distance

Definition

$$\mathit{dist} = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Where n is the number of dimensions (attributes) and p_k and q_k are the k^{th} attributes of data points p and q .

Example: Euclidean Distance



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

Normalization

- Attributes should be normalized so that all attributes can have **equal impact** on the computation of distances.
- Consider the following pair of data points
 - \mathbf{x}_i : (0.1, 20) and \mathbf{x}_j : (0.9, 720).

$$\text{dist}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(0.9 - 0.1)^2 + (720 - 20)^2} = 700.000457$$

- The distance is almost completely dominated by $(720-20) = 700$.
- Solution: Normalize attributes to all have a common value range, for instance [0,1].

Normalizing Attribute Values in Rapidminer

The screenshot displays the Rapidminer software interface. On the left, a workflow is visible in the 'Process' tab, consisting of three connected nodes: 'Retrieve', 'Normalize', and 'Clustering'. The 'Normalize' node is highlighted with an orange border. On the right, the 'Parameters' tab for the 'Normalize' node is open. The parameters are as follows:

- create view:** ☐
- attribute filter type:** subset (dropdown menu)
- attributes:** Select Attributes... (button)
- invert selection:** ☐
- include special attributes:** ☐
- method:** range transformation (dropdown menu)
- min:** 0.0 (text input)
- max:** 1.0 (text input)

Similarity of Binary Attributes

- Common situation is that objects, p and q , have only binary attributes.
 - Products in shopping basket
 - Courses attended by students
- We compute similarities using the following quantities:
 - M_{11} = the number of attributes where p was 1 and q was 1
 - M_{00} = the number of attributes where p was 0 and q was 0
 - M_{01} = the number of attributes where p was 0 and q was 1
 - M_{10} = the number of attributes where p was 1 and q was 0

Symmetric Binary Attributes

- A binary attribute is **symmetric** if both of its states (0 and 1) have equal importance, and carry the same weights, e.g., male and female
- Similarity measure: **Simple Matching Coefficient**

$$SMC(\mathbf{x}_i, \mathbf{x}_j) = \frac{M_{11} + M_{00}}{M_{01} + M_{10} + M_{11} + M_{00}}$$

Number of matches / number of all attributes values

Asymmetric Binary Attributes

- Asymmetric: If one of the states is more important than the other.
 - By convention, state 1 represents the more important state.
 - 1 is typically the rare or infrequent state.
 - Examples: Shopping basket, word vector
- Similarity measure: **Jaccard Coefficient**

$$J(\mathbf{x}_i, \mathbf{x}_j) = \frac{M_{11}}{M_{01} + M_{10} + M_{11}}$$

Number of 11 matches / number of not-both-zero attributes values

SMC versus Jaccard: Example

$p = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$
 $q = 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$

example interpretation:
p bought item 1
q bought item 7 and 10

$M_{11} = 0$ (the number of attributes where p was 1 and q was 1)

$M_{00} = 7$ (the number of attributes where p was 0 and q was 0)

$M_{01} = 2$ (the number of attributes where p was 0 and q was 1)

$M_{10} = 1$ (the number of attributes where p was 1 and q was 0)

$$SMC = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0 + 7) / (2 + 1 + 0 + 7) = 0.7$$

$$J = (M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$

SMC versus Jaccard: Question

- Which of the two measures would you use
- ...for a dating agency?
 - hobbies
 - favorite bands
 - favorite movies
 - ...
- ...for the Wahl-O-Mat?
 - (dis-)agreement with political statements
 - recommendation for voting



Using Weights to Combine Similarities

- You may not want to treat all attributes the same.
 - Use weights w_k which are between 0 and 1 and sum up to 1.
 - Weights are set according to the importance of the attributes.
- Example: **Weighted Euclidean Distance**

$$\text{dist}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_r(x_{ir} - x_{jr})^2}$$

Combining different Similarity Measures

Silk Workbench

Workspace: Cora

Editor: linkcora

Generate Links

Reference Links

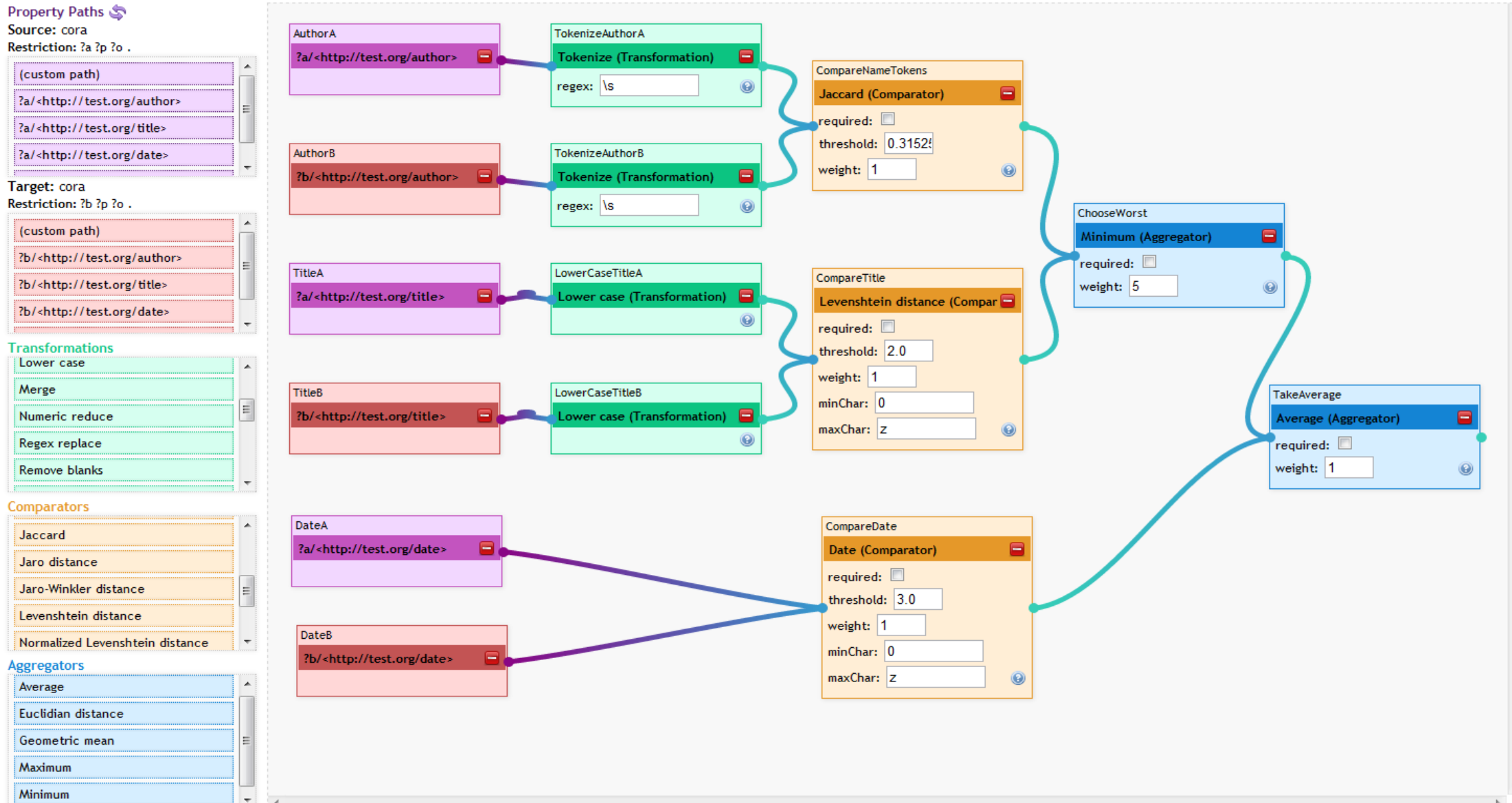
Learn

About

Export as Silk-LS

Help

Precision = 0.98 | Recall = 0.20 | F-measure = 0.33



How to choose a good Clustering Algorithm?

- “Best” algorithm depends on
 1. the analytical goals of the specific use case
 2. the distribution of the data
- Standardization of data, feature selection, distance function, and parameter settings have equally high influence on results.
- Due to these complexities, the common practice is to
 1. run several algorithms using different distance functions, feature subsets and parameter settings, and
 2. then visualize and interpret the results based on knowledge about the application domain as well as the goals of the analysis.

Literature Reference for this Slideset

Pang-Ning Tan, Michael Steinbach, Vipin Kumar:
Introduction to Data Mining.
Pearson / Addison Wesley.

Chapter 8: Cluster Analysis

Chapter 8.2: K-Means

Chapter 8.3: Agglomerative Hierarchical Clustering

Chapter 8.4: DBSCAN

Chapter 2.4: Measures of Similarity and Dissimilarity

