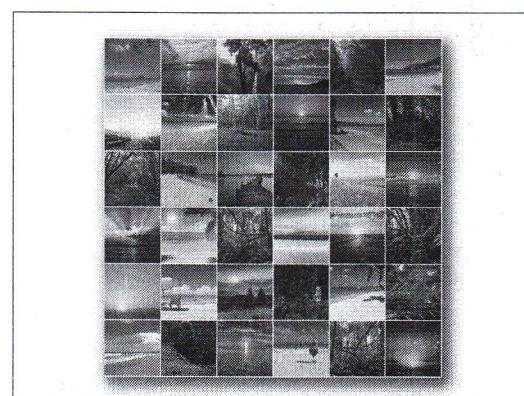


POLITECNICO DI MILANO

Clustering
Data Mining and Text Mining

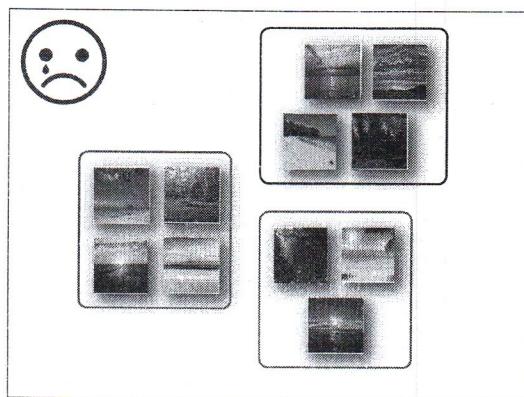
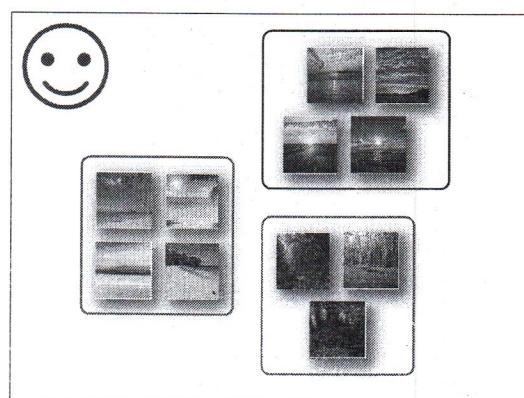
can be applied to any type of data
(we don't need, for instance, the notion
of "transaction" as we did with
associations rules)

Prof. Pierluca Lanzi



Clustering searches for "natural"
grouping/structure in un-labeled data

there is
no target



Clustering algorithms group a collection of data points into "clusters" according to some distance measure

Data points in the same cluster should have a small distance from one another

Data points in different clusters should be at a large distance from one another

low distance = high similarity

Clustering Applications

- Marketing
 - Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- Land use
 - Identification of areas of similar land use in an earth observation database
- Insurance
 - Identifying groups of motor insurance policy holders with a high average claim cost
- City-planning
 - Identifying groups of houses according to their house type, value, and geographical location
- Earthquake studies
 - Observed earthquake epicenters should be clustered along continent faults

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

What is Cluster Analysis?

- A cluster is a collection of data objects
 - Similar to one another within the same cluster
 - Dissimilar to the objects in other clusters
- Cluster analysis
 - Given a set data points try to understand their structure
 - Finds similarities between data according to the characteristics found in the data
 - Groups similar data objects into clusters
 - It is unsupervised learning since there is no predefined classes

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

What Is Good Clustering?

- A good clustering consists of high-quality clusters with
 - High intra-class similarity
 - Low inter-class similarity
- The quality of a clustering result depends on both
 - The similarity measure used by the method and
 - Its implementation (the algorithms used to find the clusters)
- The quality of a clustering method is also measured by its ability to discover some or all the hidden patterns
- Evaluation
 - Various measures of intra/inter cluster similarity
 - Manual inspection
 - Benchmarking on existing labels

→ different similarity measures give different results
 → influences a lot the quality of the clustering

→ also the implementation matters!

we can consider a labeled data (e.g. ins dataset), we remove the labels and we perform the clustering. if we obtain clusters similar to the existing labels then we can conclude that labels reflect an (underlying) existing structure of the data

Measure the Quality of Clustering

- Dissimilarity/Similarity metric
 - Similarity expressed in terms of distance function, typically a metric, $d(i, j)$
 - Definitions of distance functions are usually very different for interval-scaled, Boolean, categorical, ordinal ratio, and vector variables
 - Weights can/should be associated with different variables based on applications and data semantics
- Cluster quality measure
 - Separate from distance, there is a "quality" function that measures the "goodness" of a cluster
 - It is hard to define "similar enough" or "good enough" as the answer is typically highly subjective

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

"There are 2 major groups, in each there are 3 sub-groups, in one of the 3 there are n sub-subgroups, ..."

Clustering Strategies

- Hierarchical vs point assignment
- Numeric and/or symbolic data
- Deterministic vs. probabilistic
- Exclusive vs. overlapping
- Hierarchical vs. flat
- Top-down vs. bottom-up

"This is a good customer" vs.
"This is a good customer at 80%"

We can start from all the data and then divide or we can start from one datum and then aggregate

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Data Structures

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
...

Data Matrix

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

Dis/Similarity Matrix

$$\begin{bmatrix} 0 & d(2,1) & d(3,1) & \dots & d(n,1) \\ d(2,1) & 0 & d(3,2) & \dots & d(n,2) \\ d(3,1) & d(3,2) & 0 & \dots & \dots \\ \vdots & \vdots & \vdots & \ddots & 0 \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

Squared matrix $n \times n$
(n = number of data)

this for instance is the distance between the observation 2 and the observation 3

if the distance is 0 it's because we're looking at the distance between an element and itself

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Distance and Similarity Measures

independently on the algorithm we're going to use a distance/similarity measure (in fact the dist/sim. is almost \perp to the algorithm)

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Distance Measures

- Given a space and a set of points on this space, a distance measure $d(x,y)$ maps two points x and y to a real number, and satisfies three axioms

- $d(x,y) \geq 0$
- $d(x,y) = 0$ if and only $x=y$
- $d(x,y) = d(y,x)$
- $d(x,y) \leq d(x,z) + d(z,y)$

Prof. Pierluca Lanzi POLITECNICO DI MILANO

What Similarity Measure?

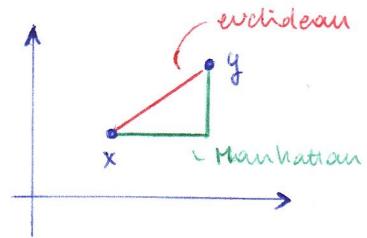
- Euclidean distance is the typical function used to compute the similarity between two examples

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- Another popular metric is city-block (Manhattan) metric, distance is the sum of absolute differences

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sum_{i=1}^n |x_i - y_i|$$

Prof. Pierluca Lanzi POLITECNICO DI MILANO



Jaccard Distance

- Jaccard distance is a measure of how dissimilar two sets are. Examples are represented by binary variables and are viewed as representations of set
- Jaccard distance is defined as

$$d(x, y) = 1 - J(x, y)$$

- J is the Jaccard similarity which is computed as the number of

$$J(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

- Which can also be interpreted as the percentage of identical attributes

we're considering boolean vectors (or itemsets) and we're making intersection over union without considering multiple copies of elements (actually they don't have to be boolean)

Example:

person 1 buy: {ZA, B, 3C}
 person 2 buy: {A, 3C, D}

$$|x \cap y| = |\{3C\}| = 1$$

$$|x \cup y| = |\{A|ZA, B, 3C, D\}| = 4$$

$$J(x, y) = \frac{1}{4} \rightarrow d(x, y) = \frac{3}{4}$$

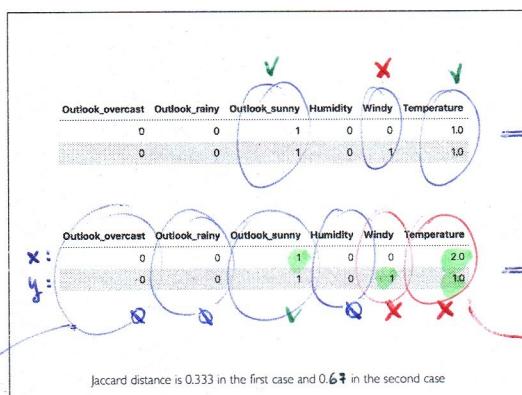
However JACCARD IS TYPICALLY USED FOR BINARY DATA!

Examples of Jaccard Distance

- Consider the following example consisting of two examples

Outlook_overcast	Outlook_rainy	Outlook_sunny	Humidity	Windy	
0	0	1	0	0	
0	0	1	0	1	

- They are at Jaccard distance of 0.5 since the first example "contains" only Outlook_sunny whereas the second example "contains" both Outlook_sunny and Windy so intersection has size one while union has size two. The Jaccard distance is 0.5 so the similarity is 1-0.5 that is 0.5



Unlike the Hamming distance, Jaccard considers 0 as a NON EXISTING VALUE

Hamming Distance

- Hamming distance between two vectors is the number of components in which they differ
- Or equivalently, given the number of variables p , and the number m of matching components, we define

$$d(x, y) = \frac{p - m}{p}$$

We need values ≠ 0 to count it in the union and we need a MATCH to count it in the intersection

Examples of Hamming Distance

- Consider the same example

Outlook_overcast	Outlook_rainy	Outlook_sunny	Humidity	Windy	
x: 0	0	0	1	0	0
y: 0	0	0	1	0	1

Hamming

- They are at Jaccard distance of 0.2 since all their attribute values are the same except for Windy so they differ for one value over 5, that is 0.2

$$d(x, y) = \frac{1}{5} = 0.2$$

Outlook_overcast	Outlook_rainy	Outlook_sunny	Humidity	Windy	Temperature
0	0	1	0	0	1.0
0	0	1	0	1	1.0

Outlook_overcast	Outlook_rainy	Outlook_sunny	Humidity	Windy	Temperature
0	0	1	0	0	2.0
0	0	1	0	1	1.0

Hamming distance is 1/6 in the first case and 1/3 in the second case.

very useful in text mining

Cosine Distance

22

- The cosine distance between x, y is the angle that the vectors to those points make

$$d(x, y) = \arccos \frac{\sum_1^n x_i y_i}{\sqrt{\sum_i^n x_i^2} \sqrt{\sum_i^n y_i^2}}$$

- This angle will be in the range 0 to 180 degrees, regardless of how many dimensions the space has.
- Example: given $x = (1, 2, -1)$ and $y = (2, 1, 1)$ the angle between the two vectors is 60

Prof Pierluca Longi POLITECNICO DI MILANO

Cosine Similarity

23

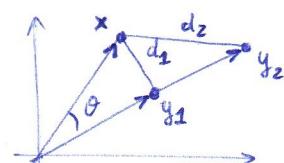
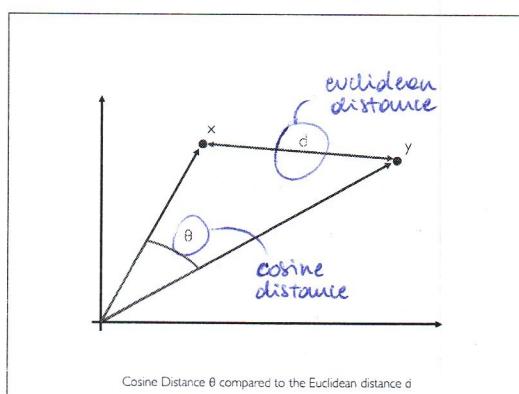
- The cosine similarity between x, y is simply computed as,

$$s(x, y) = \frac{\sum_1^n x_i y_i}{\sqrt{\sum_i^n x_i^2} \sqrt{\sum_i^n y_i^2}} = \text{normalized projection of one vector onto the other}$$

(we cannot do 1-dissimilarity to obtain similarity in this case)

Prof Pierluca Longi POLITECNICO DI MILANO

There are 3 (normalized) vectors: the projection of 2 on 3 will be higher than the projection of 1 on 3 \Rightarrow 2 and 3 are more similar than 1 and 3 (and the cosine distance confirms)



y_1 and y_2 are the same point but scaled: $d_1 \neq d_2$, $\theta_1 = \theta_2 \equiv \theta$ (the Euclidean distance between x and y_1 is different from the Euclidean distance between x and y_2 , meanwhile the cosine distance remains the same)

Edit Distance

25

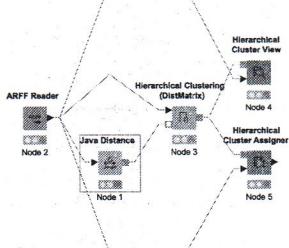
- The distance between a string $x=x_1x_2\dots x_n$ and $y=y_1y_2\dots y_m$ is the smallest number of insertions and deletions of single characters that will transform x into y .
- Alternatively, the edit distance $d(x, y)$ can be computed as the longest common subsequence (LCS) of x and y and then,

$$d(x, y) = |x| + |y| - 2|\text{LCS}|$$

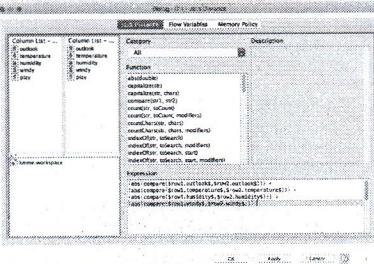
- Example
 - The edit distance between $x=abcde$ and $y=acfdeg$ is 3 (delete b, insert f, insert g), the LCS is acde which is coherent with the previous result

Prof Pierluca Longi POLITECNICO DI MILANO

Some tools allow users to define their own distance function



KNIME workflow applying hierarchical clustering with a custom distance function implemented in Java



The definition of the custom distance function used in the previous KNIME workflow.

Other tools provide only the usual distance function (Euclidean, Manhattan, etc.)

When a custom function is not an option,
we can try to transform the data

Normalization

Normalization and Other Issues

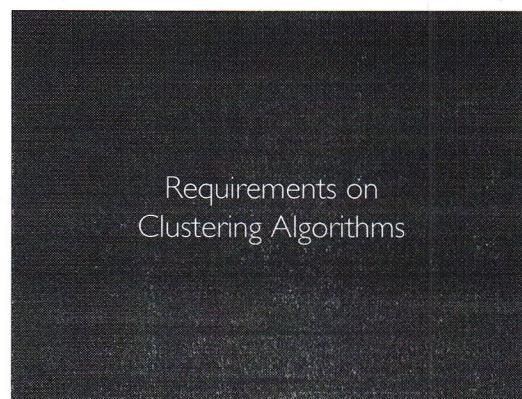
- Different attributes are measured on different scales, thus we often need to normalize them, using for instance range normalization or standard score normalization

$$x'_i = \frac{x_i - \min_i x_i}{\max_i x_i - \min_i x_i}$$

$$x'_i = \frac{x_i - \mu}{\sigma}$$

- For nominal attributes, the distance either 0 (the value is the same) or 1 (the value is different)
- Missing values are usually assumed to be maximally distant (given normalized attributes)

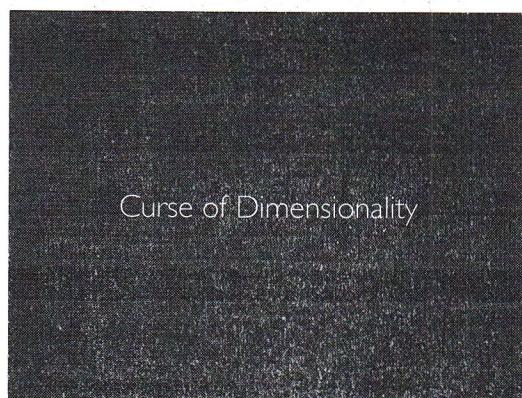
Prof. Pierluca Lanzi | 31 | POLITECNICO DI MILANO



Requisites for Clustering Algorithms

- Scalability
- Ability to deal with different types of attributes
- Ability to handle dynamic data
- Discovery of clusters with arbitrary shape
- Minimal requirements for domain knowledge to determine input parameters
- Able to deal with noise and outliers
- Insensitive to order of input records
- High dimensionality
- Incorporation of user-specified constraints
- Interpretability and usability

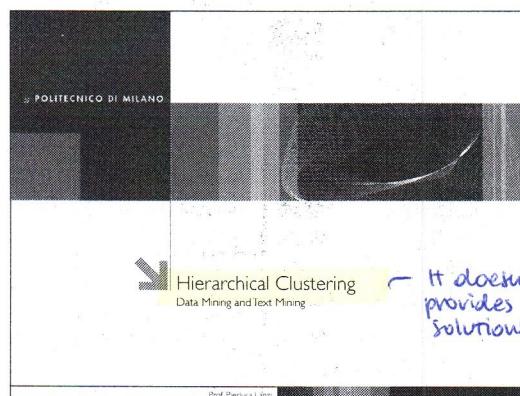
Prof. Pierluca Lanzi | 33 | POLITECNICO DI MILANO



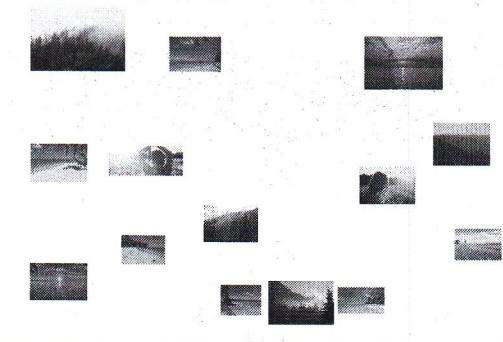
Curse of Dimensionality

in high dimensions, almost all pairs of points are equally far away from one another

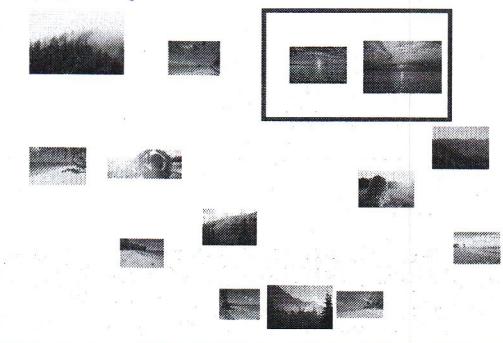
almost any two vectors are almost orthogonal



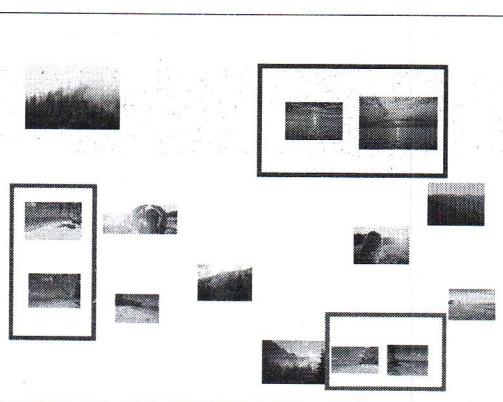
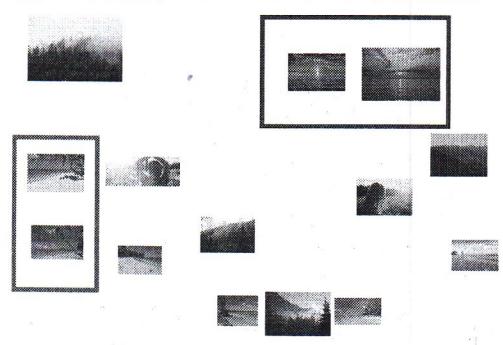
At the beginning we have n clusters.



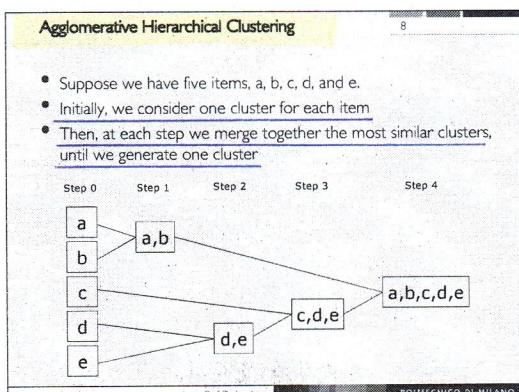
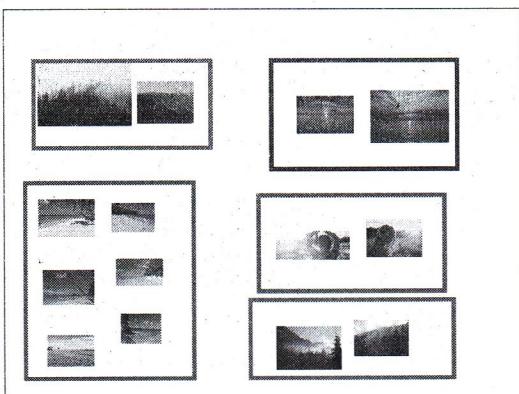
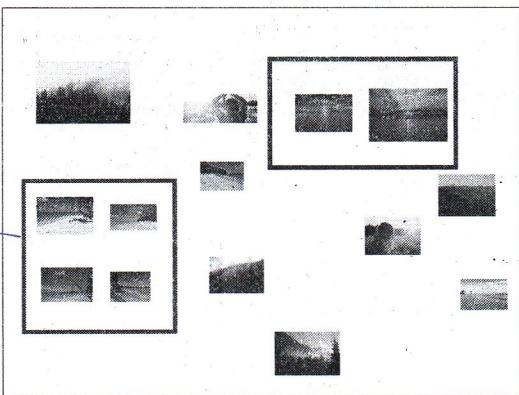
We put together the two more similar. Now we have $n-1$ clusters.



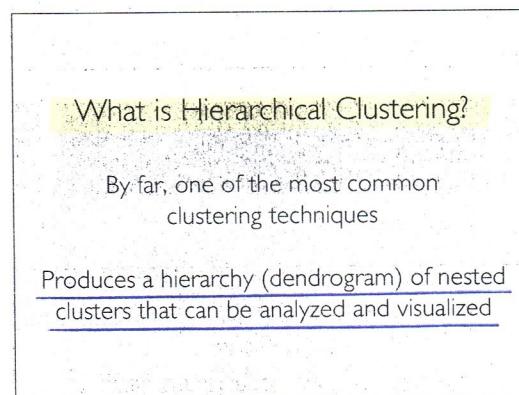
We iterate, now we have $n-2$ clusters.

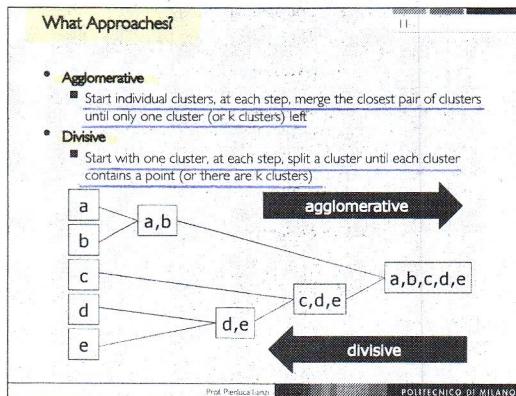


In this step we unified two already existing groups (before we used a distance for one element vs. another; now we're using a distance between groups? how do we do it?)



Clustering	Partition
C ₀	{A} {B} {C} {D} {E}
C ₁	{A, B} {C} {D} {E}
C ₂	{A, B} {C} {D, E}
C ₃	{A, B} {C, D, E}
C ₄	{A, B, C, D, E}

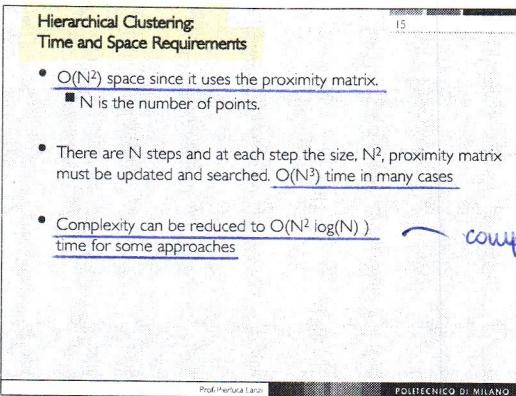
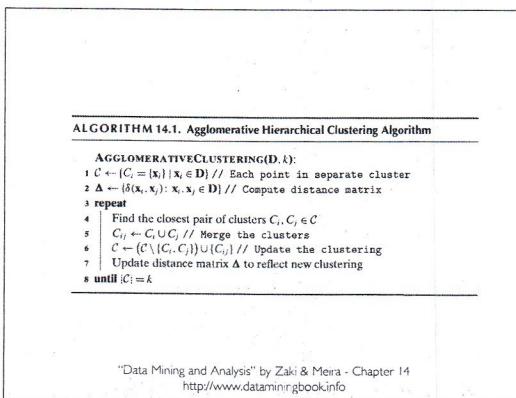
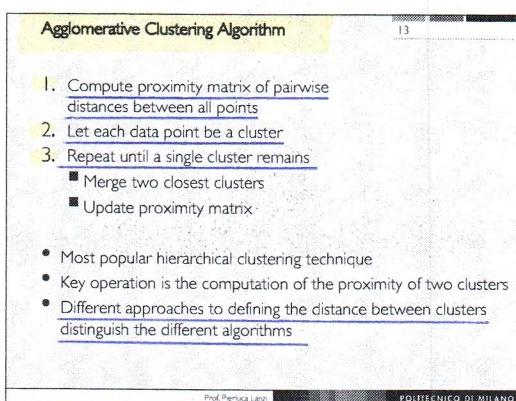
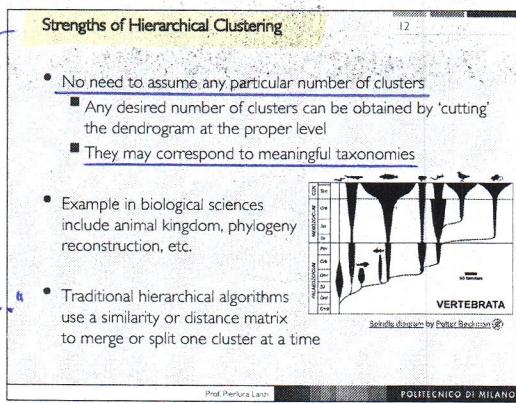




by giving a set of solutions instead of just one we give more informations about the data:

"there are k macrogroups, one of these macrogroups could be partitioned in m subgroups and another of these macrogroups can be partitioned in l subgroups,"

we're saying more about the structure of the data



Efficient Implementation

16

- Compute the distance between all pairs of points [$O(N^2)$]
- insert the pairs and their distances into a priority queue to find the min in one step [$O(N^2)$]
- When two clusters are merged, we remove all entries in the priority queue involving one of these two clusters [$O(N \log N)$]
- Compute all the distances between the new cluster and the remaining clusters [$O(N \log N)$]
- Since the last two steps are executed at most N time, the complexity of the whole algorithm is $O(N^2 \log N)$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

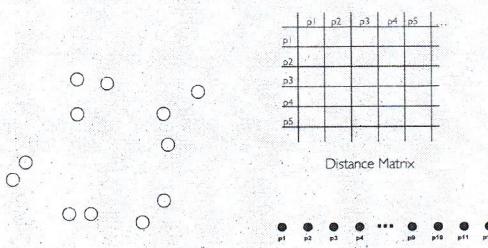
Distance Between Clusters

It's simple to evaluate the distance between two points, but how about the distance between a point and a cluster? How about the distance between two clusters?

Initial Configuration

18

- Start with clusters of individual points and the distance matrix



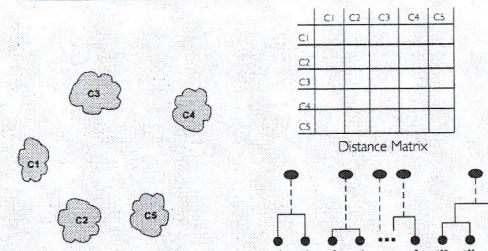
Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Intermediate Situation

19

- After some merging steps, we have some clusters



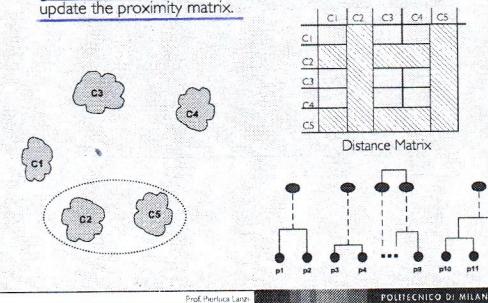
Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Intermediate Situation

20

- We want to merge the two closest clusters (C2 and C5) and update the proximity matrix.



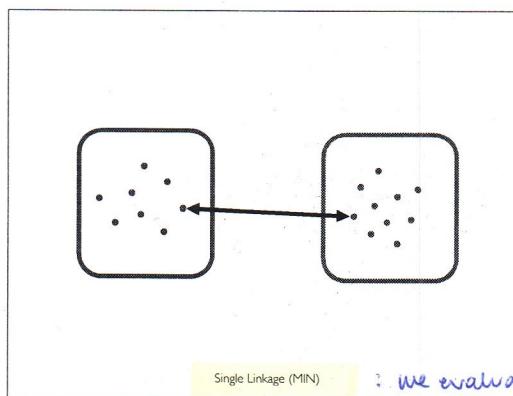
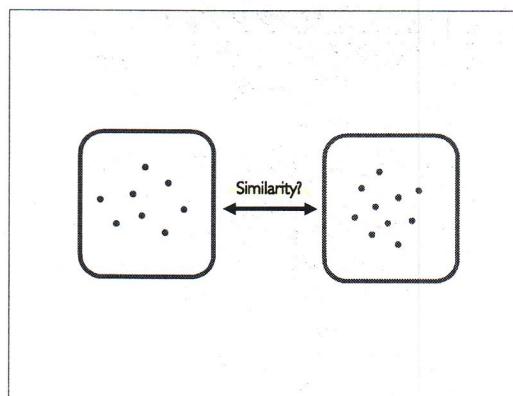
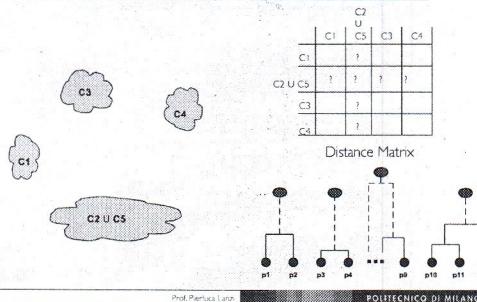
Prof. Pierluca Lanzi

POLITECNICO DI MILANO

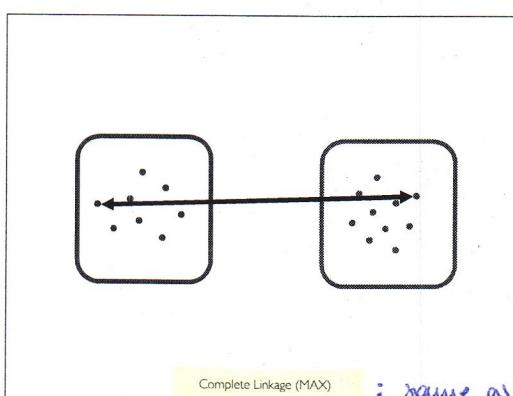
After Merging

21

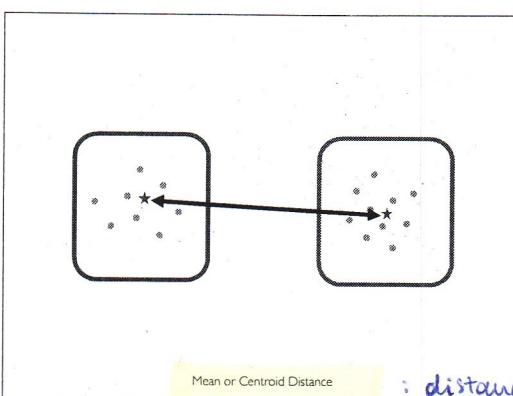
- The question is "How do we update the proximity matrix?"



: we evaluate all the distances and we consider as distance between the two groups the minimum distance

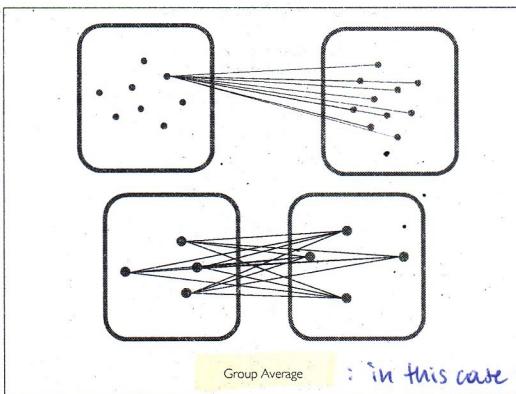


: same as before but with the maximum distance



: distance between the two averages
(averages are each the average of one group)

How are they different?
For example: if the minimum distance is > 0 then the two groups are separable, if the maximum distance is > 0 the two groups can still be overlapped



: in this case we do the average of the distances between every point of one group and every point of the other

Typical Alternatives to Calculate the Distance Between Clusters 27

- Single link (or MIN)**
 - smallest distance between an element in one cluster and an element in the other; i.e., $d(C_i, C_j) = \min(t_{ij}, t_{ji})$
- Complete link (or MAX)**
 - largest distance between an element in one cluster and an element in the other; i.e., $d(C_i, C_j) = \max(t_{ij}, t_{ji})$
- Mean Distance**
 - distance between the centroids of two clusters, i.e., $d(C_i, C_j) = d(\mu_i, \mu_j)$ where μ_i and μ_j are the cluster means (or centroids)
- Group average**
 - average distance between an element in one cluster and an element in the other; i.e., $d(C_i, C_j) = \text{avg}(d(t_{ij}, t_{ji}))$
- ...

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Example 28

- Suppose we have five items, a, b, c, d, and e.
- We want to perform hierarchical clustering on five instances following an agglomerative approach
- First we compute the distance or similarity matrix
- D_{ij} is the distance between instance "i" and "j"

$$D = \begin{pmatrix} 0.0 & & & & \\ 2.0 & 0.0 & & & \\ & 6.0 & 5.0 & 0.0 & \\ & 10.0 & 9.0 & 4.0 & 0.0 \\ & 9.0 & 8.0 & 5.0 & 3.0 & 0.0 \end{pmatrix} \quad \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix}$$

The first merge is between the elements 1 and 2 and it's performed at distance 2.0
→ the two rows and the two columns have to be merged in one row and one column
The new element is (1,2). What is the distance between (1,2) and the other elements?

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Example 29

- Group the two instances that are closer
- In this case, a and b are the closest items ($D_{2,1}=2$)
- Compute again the distance matrix, and start again.
- Suppose we apply single-linkage (MIN), we need to compute the distance between the new cluster {1,2} and the others

$$\begin{aligned} d(12,3) &= \min[d_{13}, d_{23}] = d_{23} = 5.0 \\ d(12,4) &= \min[d_{14}, d_{24}] = d_{24} = 9.0 \\ d(12,5) &= \min[d_{15}, d_{25}] = d_{25} = 8.0 \end{aligned}$$

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Example 30

- The new distance matrix is,
- At the end, we obtain the following dendrogram

$$D = \begin{pmatrix} 0.0 & & & & \\ 5.0 & 0.0 & & & \\ 9.0 & 4.0 & 0.0 & & \\ 8.0 & 5.0 & (3.0) & 0.0 & \\ (1,2) & 3 & 4 & 5 & \end{pmatrix}$$

the dendrogram says the merges + the distances at which the merges were done

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Determining the Number of Clusters

Hierarchical clustering generates
a set of N possible partitions

Which one should I choose?

Ideally a good clustering should
data partition points so that ...

Data points in the same cluster should have
a small distance from one another

} measure of COHESION

Data points in different clusters should be at
a large distance from one another.

} measure of SEPARATION

How can we evaluate
the quality of a clustering solution?

3 possible approaches

Clustering Quality Measures

- Internal Validation Measures
 - Employ criteria that are derived from the data itself
 - For instance, intraclass and intercluster distances to measure cluster cohesion and separation
 - Cohesion evaluates how similar are the points in the same cluster
 - Separation, how far apart are the points in different clusters
- External Validation Measures
 - Use prior or expert-specified knowledge about the clusters
 - For example, we cluster the Iris data using the four input variables, then we evaluate the clustering using known class labels
 - Employs criteria that are not inherent to the dataset
- Relative Validation Measures
 - Aim to directly compare different solutions, usually those obtained via different parameter settings for the same algorithm.

$WSS \rightarrow 0$ as $k \rightarrow N$

$BSS = 0$ if $k=1$
and BSS starts
increasing as k
increases

Internal Measures: Cohesion and Separation 36

- **Cohesion measures how closely related are objects in a cluster**
 - Within-cluster sum of squares
$$WSS(C) = \sum_{i=1}^k \sum_{x_j \in C_i} d(x_j, \mu_i)^2$$
 - where μ_i is the centroid of cluster C_i (in case of Euclidean spaces)
- **Separation measures how well separated a cluster is from other clusters**
 - Between-cluster sum of squares
$$BSS(C) = \sum_{i=1}^k |C_i| d(\mu, \mu_i)^2$$
 - where μ is the centroid of the whole dataset (= mean of the whole dataset)

Prof. Renato Carloni POLITECNICO DI MILANO

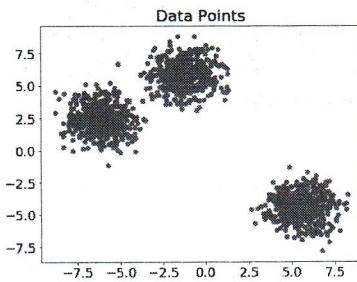
(μ_i = mean of the points of group i)

(= mean of the whole dataset)

Evaluation of Hierarchical Clustering using Knee/Elbow Analysis

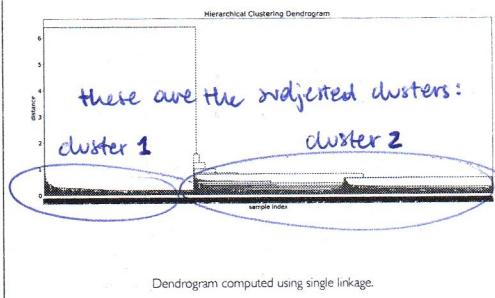
plot the WSS and BSS for every clustering and look
for a knee in the plot that shows a significant
modification in the evaluation metrics

Suppose that we have:

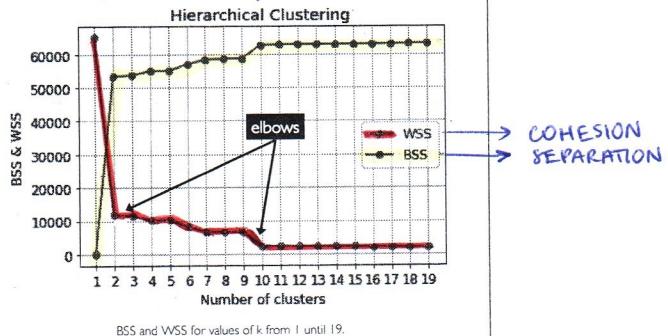


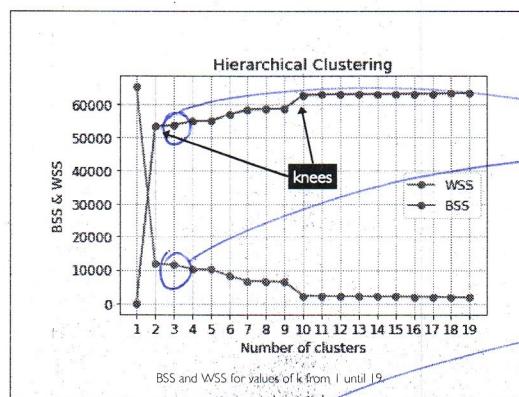
Example data generated using the make_blobs function of Scikit-Learn

The dendrogram that comes out is:



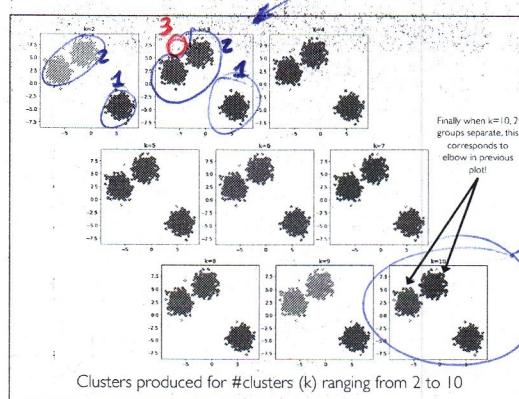
The knee/elbow analysis:





knowing the data we know that there are 3 clusters. That's why we're a little bit surprised to see only 2 groups. What's happening in groups when we choose 2 or 3 clusters?

It creates a one-point cluster (not really interesting)



($K=10$)
When we step from 9 to 10 we finally creates 3 supergroups. The problem is the small clusters. We can either eliminate them or we can merge the one-point clusters to the major clusters.

How can we know that this is the scenario without analyzing the plots? We can see the numerosity of the clusters. If from $k=2$ to $k=3$ the numerosity is $[500, 500] \rightarrow [500, 499, 1]$ we already know what's going on.

Knee/elbow plots might not tell the whole story!

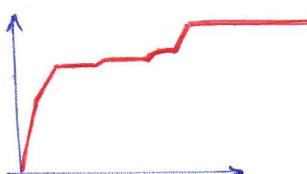
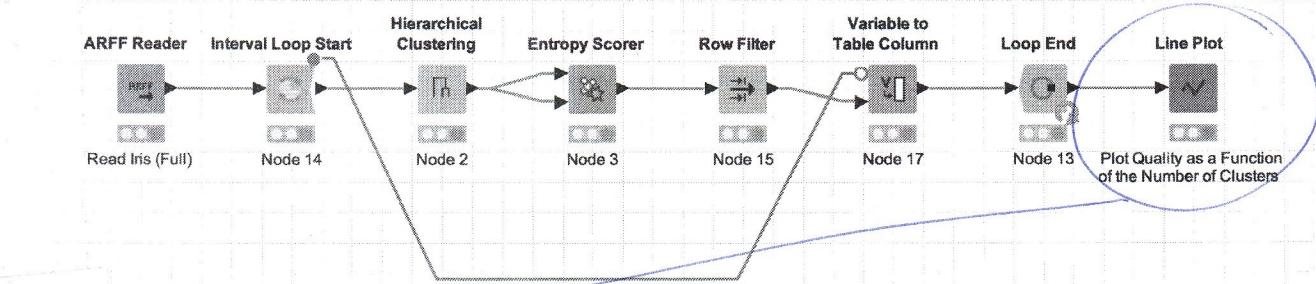
In the example, the plots suggested several choices of clustering

We should use the elbow as a guide but test multiple solutions around that value

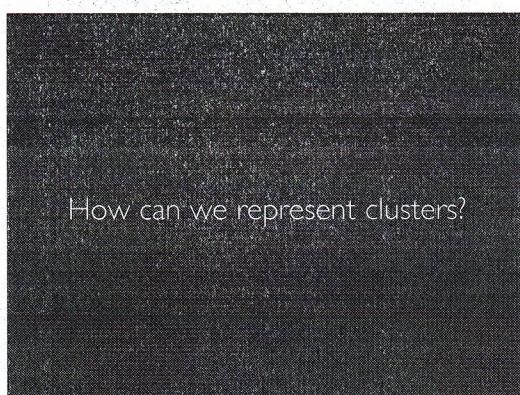
We should analyze different solutions around the knees/elbow

We should also analyze the merges to eliminate the ones due to noise

Examples using KNIME



Computing cluster quality from one to 20 clusters using the entropy scorer



Euclidean vs Non-Euclidean Spaces

- Euclidean Spaces
 - We can identify a cluster using for instance its centroid (e.g. computed as the average among all its data points)
 - Alternatively, we can use its convex hull
- Non-Euclidean Spaces
 - We can define a distance (jaccard, cosine, edit)
 - If we cannot compute a centroid, we can introduce medoid
- Medoid
 - Existing data point that we take as a cluster representative
 - Point that minimizes the sum of distances to all other points in cluster
- Alternatives to Medoid
 - Choose a point that minimize the maximum distance to another point or the sum of the squares of the distances to the other points in the cluster

48

Prof. Pierluca Lanzì POLITECNICO DI MILANO

We consider a point and we evaluate the distance between this point and all the others (we nm). We repeat for the whole cluster: the point for which the nm of distances is minimum is the MEDOID.



Hierarchical Clustering: Problems and Limitations

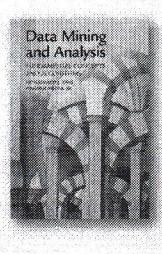
- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
 - Sensitivity to noise and outliers
 - Difficulty handling different sized clusters and convex shapes
 - Breaking large clusters
- Major weakness of agglomerative clustering methods
 - They do not scale well: time complexity of at least $O(n^2)$, where n is the number of total objects
 - They can never undo what was done previously

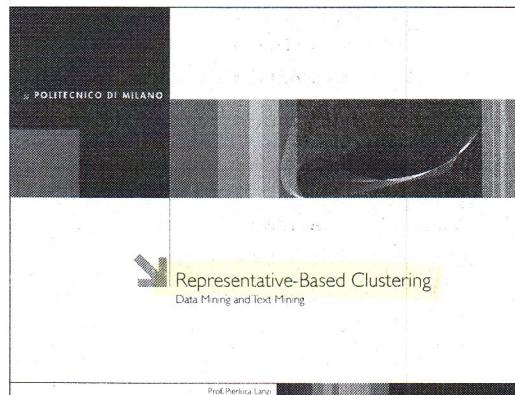
50

Prof. Pierluca Lanzì POLITECNICO DI MILANO

Study Material

- "Data Mining and Analysis" by Zaki & Meira
 - Chapter 14
- <http://www.dataminingbook.info>





For each cluster there is a point (or a parameter vector) that summarizes it

Representative-Based Algorithms

- Given a dataset of N instances, and a desired number of clusters k , this class of algorithms generates a partition C of N in k clusters $\{C_1, C_2, \dots, C_k\}$
- For each cluster there is a point that summarizes the cluster
- The common choice being the mean of the points in the cluster

$$\mu_i = \frac{1}{n_i} \sum_{x_j \in C_i} x_j$$

where $n_i = |C_i|$ and μ_i is the centroid

Prof. Pier Luca Lanzi POLITECNICO DI MILANO

Brute-force Approach

Generate all the possible clustering $C = \{C_1, C_2, \dots, C_k\}$

Then, select the best one

Unfortunately, there are $O(k^N/k!)$ possible partitions

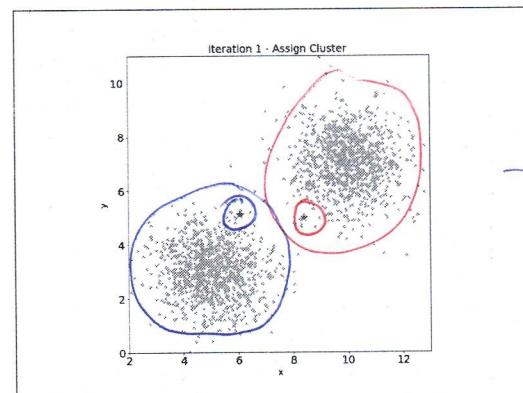
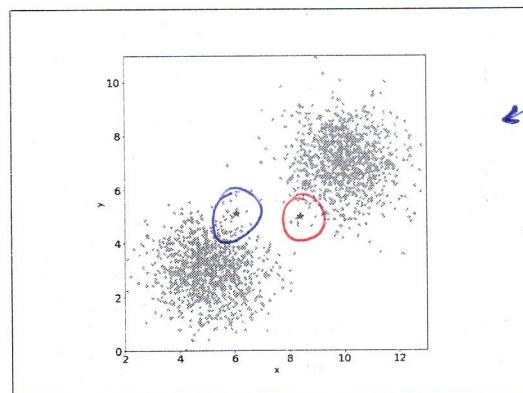
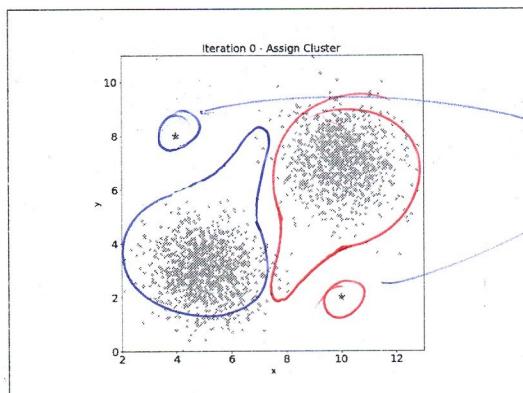
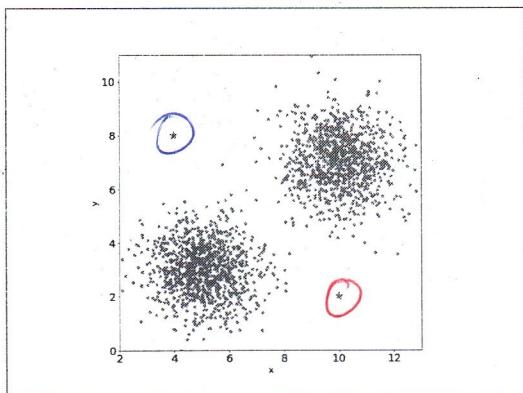
K-Means algorithm

K-Means Algorithm

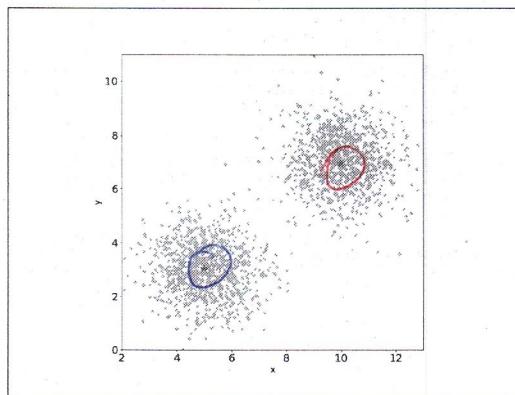
- Greedy iterative approach to find a clustering that minimizes the SSE objective:

$$\text{SSE}(C) = \sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

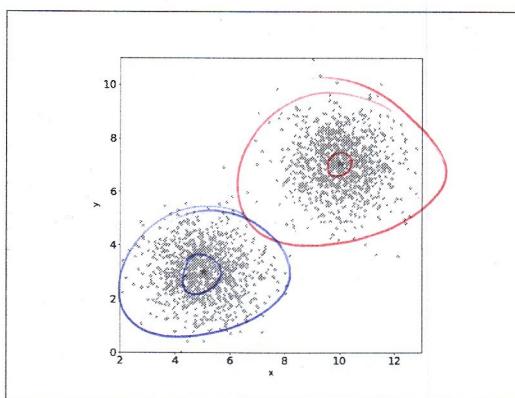
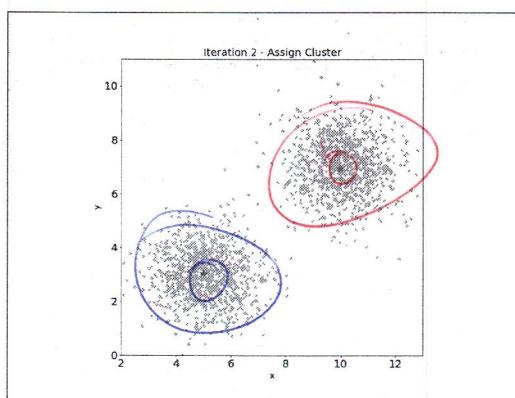
- The goal of the clustering process is thus to find

$$C^* = \arg \min_C \text{SSE}(C)$$


now that we changed (updated) centroids, also the clusters will change (update)



again we update the clusters



At some points we reach some convergence (even if we may have some oscillating points)

Notice that K-means is shape-biased! It's biased towards classes that are spheres in n-dimension

ALGORITHM 13.1. K-means Algorithm

```

K-MEANS ( $D, k, \epsilon$ ):
1  $i = 0$ 
2 Randomly initialize  $k$  centroids:  $\mu_1^0, \mu_2^0, \dots, \mu_k^0 \in \mathbb{R}^d$ 
3 repeat
4    $t \leftarrow t + 1$ 
5    $C_j \leftarrow \emptyset$  for all  $j = 1, \dots, k$ 
6   // Cluster Assignment Step
7   foreach  $x_i \in D$  do
8      $j^* \leftarrow \operatorname{argmin}_j \left\{ \|x_i - \mu_j^{t-1}\|^2 \right\}$  // Assign  $x_i$  to closest centroid
9      $C_{j^*} \leftarrow C_{j^*} \cup \{x_i\}$ 
10  // Centroid Update Step
11  foreach  $j = 1 \text{ to } k$  do
12     $\mu_j^t \leftarrow \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$ 
13 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$ 

```

"Data Mining and Analysis" by Zaki & Meira - Chapter 13
<http://www.dataminingbook.info>

K-Means Algorithm

15

- Most widely known representative-based algorithm
- Assumes an Euclidean space but sometimes it can be extended to the non-Euclidean case
- Employs a greedy iterative approach that minimizes the SSE objective. Accordingly it can converge to a local optimal instead of a globally optimal clustering.

Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

- Cluster assignment takes $O(nkd)$ time since, for each of the n points, it computes its distance to each of the k clusters, which takes d operations in d dimensions
- The centroid re-computation step takes $O(nd)$ time because it adds at total of n d -dimensional points
- Assuming that there are t iterations, the total time for K-means is given as $O(tnkd)$.
- In terms of the I/O cost it requires $O(t)$ full database scans, because we have to read the entire database in each iteration.

Centroid initialization

- **Solution 1**
 - Pick points that are as far away from one another as possible.
- **Solution 2**

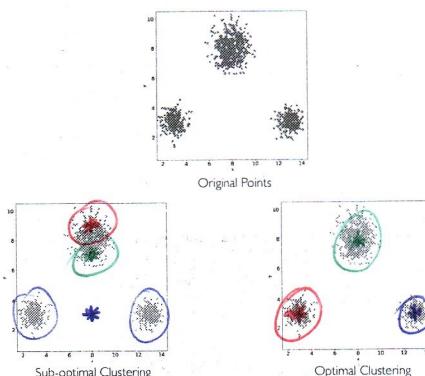
```

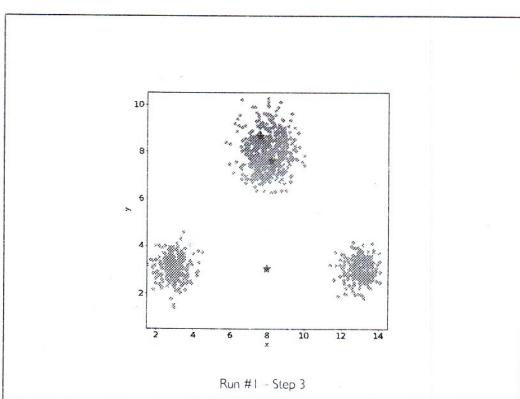
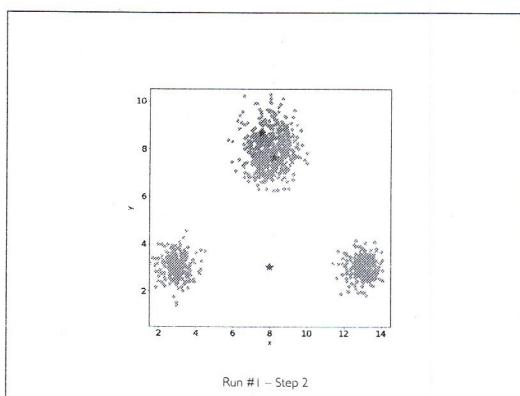
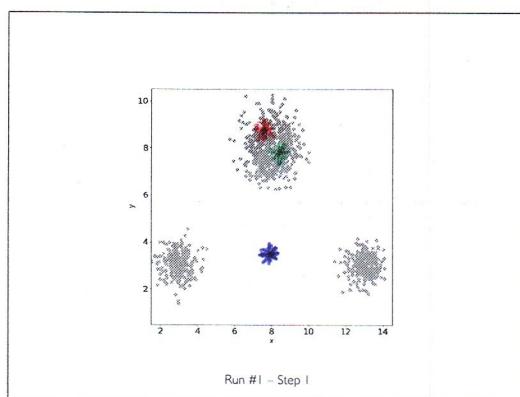
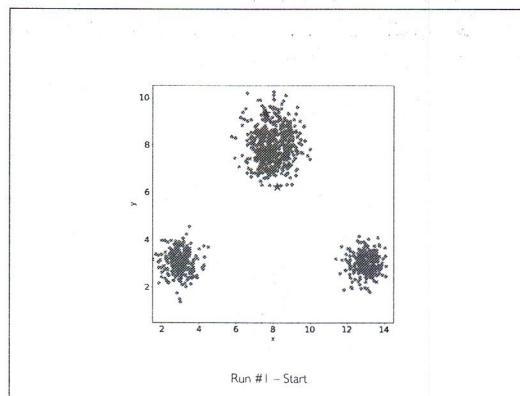
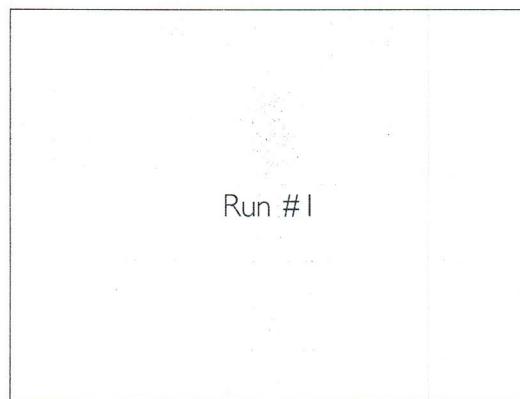
    Pick the first point at random;
    WHILE there are fewer than k points DO
        Add the point whose minimum distance
        from the selected points is as large as
        possible;
    END;
```
- **Solution 3**
 - Cluster a sample of the data, perhaps hierarchically, so there are k clusters. Pick a point from each cluster, perhaps that point closest to the centroid of the cluster.

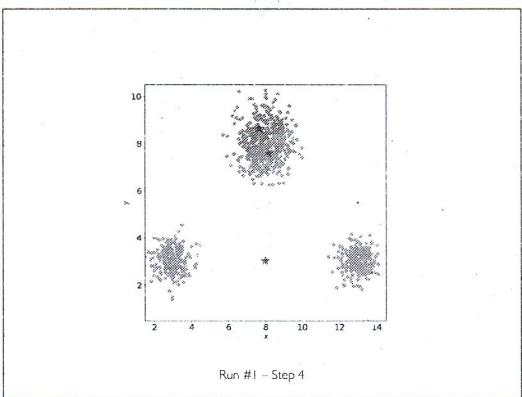
we pick randomly one point (1), then we pick the most-far point from (1), then we pick the most-far point from both, then ..

take a sample of the data, apply hierarchical clustering, cut at k (given by user) and consider the centroids of the clusters generated by hierarchical

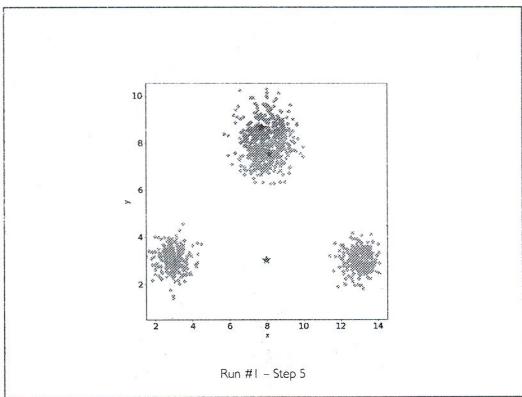
Initial centroids matters!



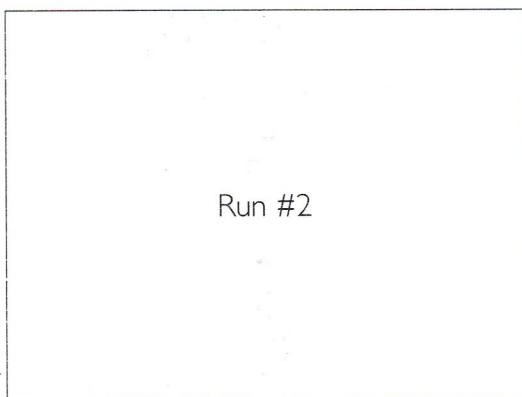




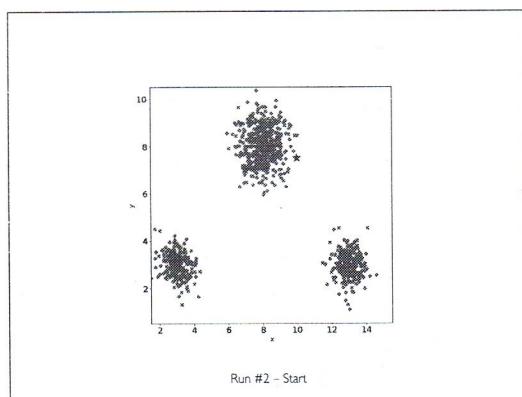
Run #1 – Step 4



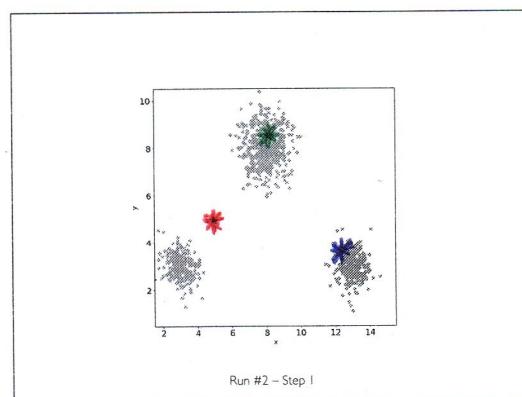
Run #1 – Step 5



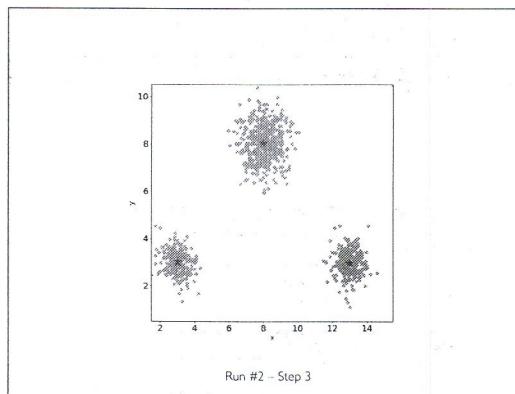
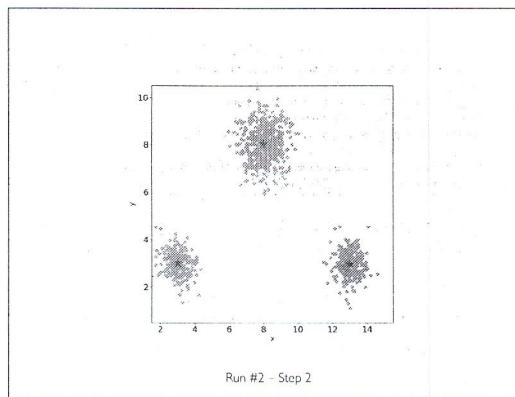
Run #2



Run #2 – Start



Run #2 – Step 1



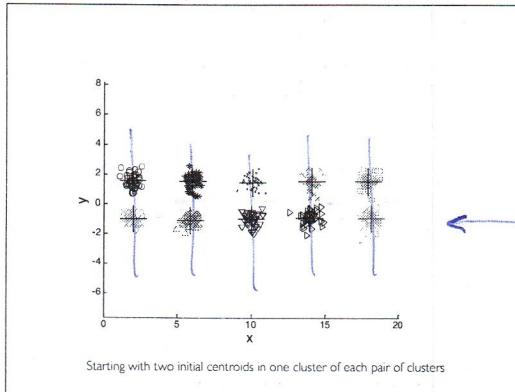
Why Selecting the Best Initial Centroids is Difficult?

- If there are K 'real' clusters then the chance of selecting one centroid from each cluster is small.
- Chance is relatively small when K is large
- If clusters are the same size, n , then

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{n^K}$$

- For example, if $K = 10$, then probability = $10!/10^{10} = 0.00036$
- Sometimes the initial centroids will readjust themselves in 'right' way, and sometimes they don't
- Consider an example of five pairs of clusters

Prof. Pier Luca Lanzi POLITECNICO DI MILANO

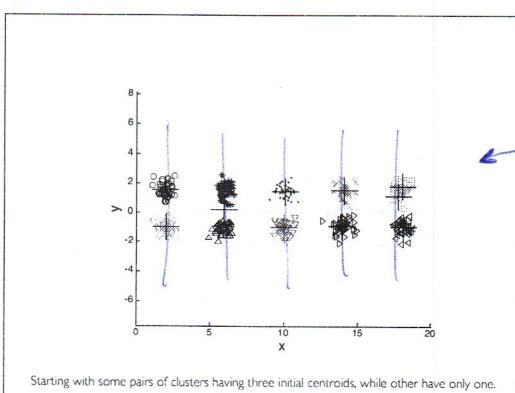


BISECTING K-MEANS

We start with 2 clusters. We take the worst cluster and we add a "seed". Then we continue.

Consider for example:

we put 2 centroids in each one ✓



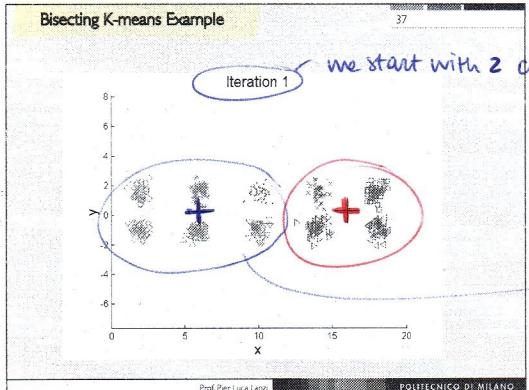
we put 2 in some, 3 in others and 1 in others ✗

Since we don't know which one will be good for us we have to establish a rigorous method.
let's try the bisection k-means!

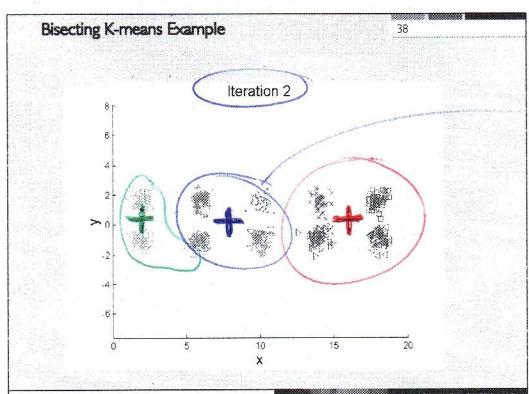


Algorithm 3 Bisecting K-means Algorithm.

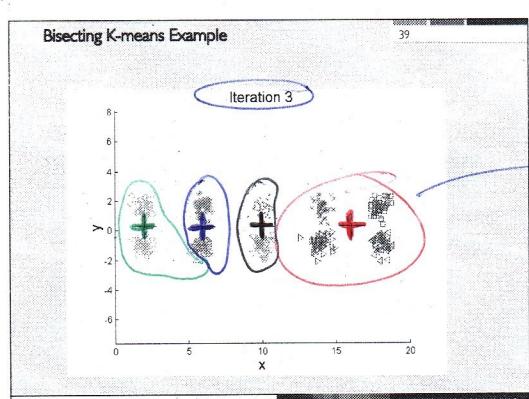
- 1: Initialize the list of clusters to contain the cluster containing all points.
- 2: repeat
- 3: Select a cluster from the list of clusters
- 4: for $i = 1$ to *number_of_iterations* do
- 5: Bisect the selected cluster using basic K-means
- 6: end for
- 7: Add the two clusters from the bisection with the lowest SSE to the list of clusters.
- 8: until Until the list of clusters contains K clusters



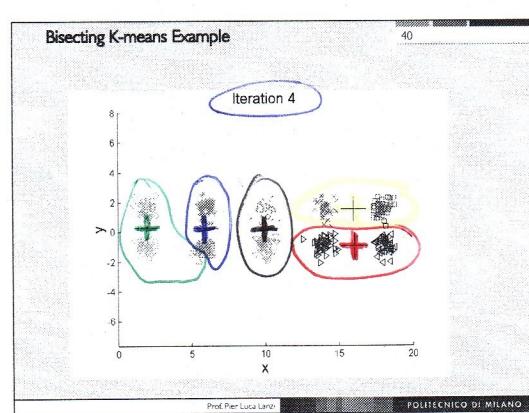
this cluster is worse in terms
of error (it less cohesive)
→ we add a seed to this
group (the other remains
unchanged)



→ this is again the less cohesive,
we put another seed



→ now this is the less cohesive.
We add a seed.

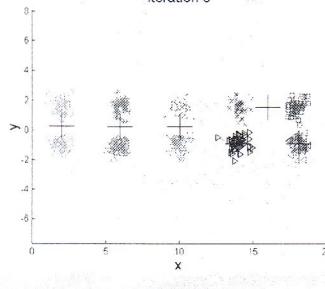


we go on.

Bisecting K-means Example

41

Iteration 5



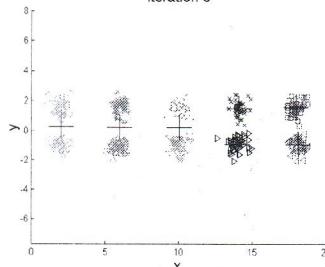
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Bisecting K-means Example

42

Iteration 6



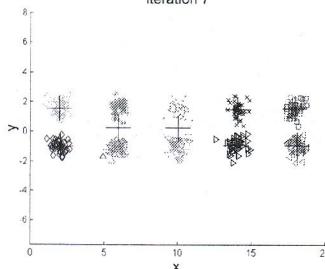
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Bisecting K-means Example

43

Iteration 7



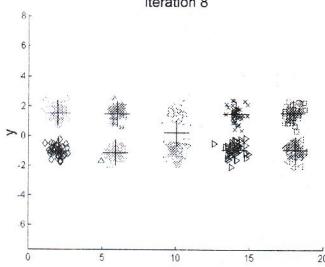
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Bisecting K-means Example

44

Iteration 8



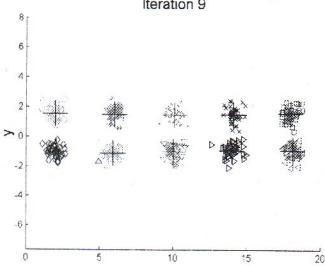
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Bisecting K-means Example

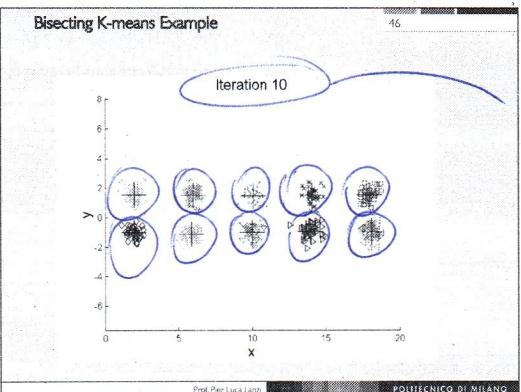
45

Iteration 9



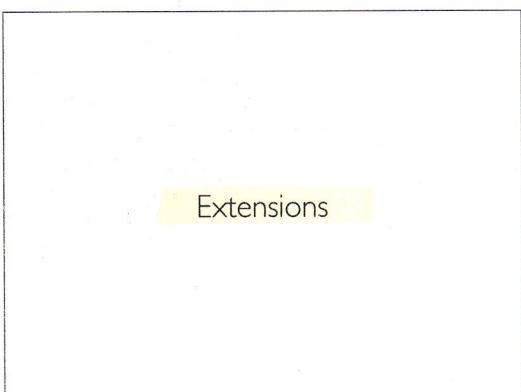
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO



We stop here because if we add a seed the SSE does not improve
 \rightarrow 10 clusters
 (doesn't improve or it improves so little)

- Dealing with the Initial Centroids Issue
- 47
- Multiple runs, helps, but probability is not on your side
 - Sample and use another clustering method (hierarchical?) to determine initial centroids
 - Select more than k initial centroids and then select among these initial centroids
 - Postprocessing
 - Bisecting K-means, not as susceptible to initialization issues
- Prof. Pier Luca Lanzi POLITECNICO DI MILANO



- Updating Centers Incrementally
- 49
- In the basic K-means algorithm, centroids are updated after all points are assigned to a centroid
 - An alternative is to update the centroids after each assignment (incremental approach)
 - Each assignment updates zero or two centroids
 - More expensive
 - Introduces an order dependency
 - Never get an empty cluster
 - Can use "weights" to change the impact
- Prof. Pier Luca Lanzi POLITECNICO DI MILANO

- Pre-processing and Post-processing
- 50
- Pre-processing
 - Normalize the data
 - Eliminate outliers
 - Post-processing
 - Eliminate small clusters that may represent outliers
 - Split 'loose' clusters, i.e., clusters with relatively high SSE
 - Merge clusters that are 'close' and that have relatively low SSE
 - These steps can be used during the clustering process
- Prof. Pier Luca Lanzi POLITECNICO DI MILANO

because distances are involved
 method, we do it
 this is based on the method

- A few variants of the k-means which differ in
 - Selection of the initial k means
 - Dissimilarity calculations
 - Strategies to calculate cluster means
- Handling categorical data: k-modes
 - Replacing means of clusters with modes
 - Using new dissimilarity measures to deal with categorical objects
 - Using a frequency-based method to update modes of clusters
 - A mixture of categorical and numerical data: k-prototype method

Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Limitations

Limitations of K-means

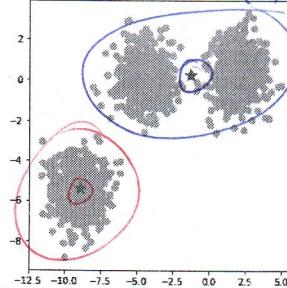
53

- K-means has problems when clusters are of differing
 - Sizes
 - Densities
 - Non-globular shapes
- K-means has also problems when the data contains outliers.
- we have to decide k a priori

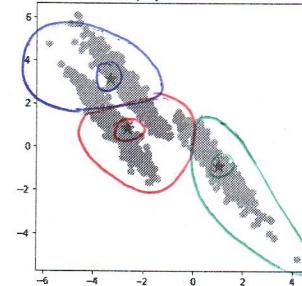
Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

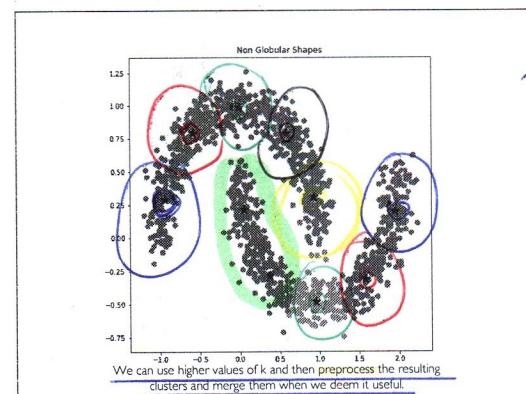
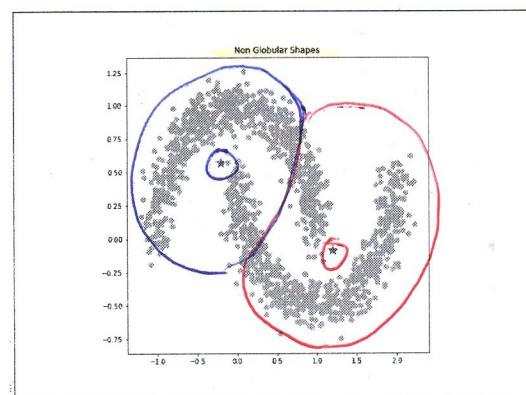
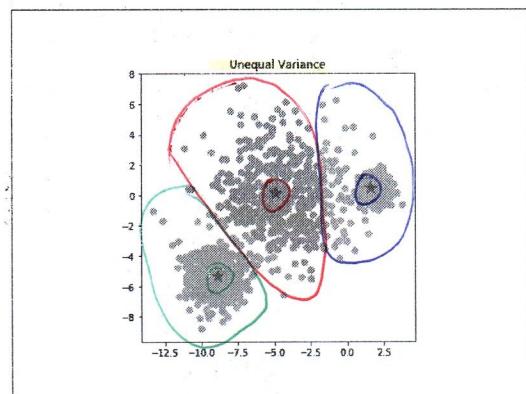
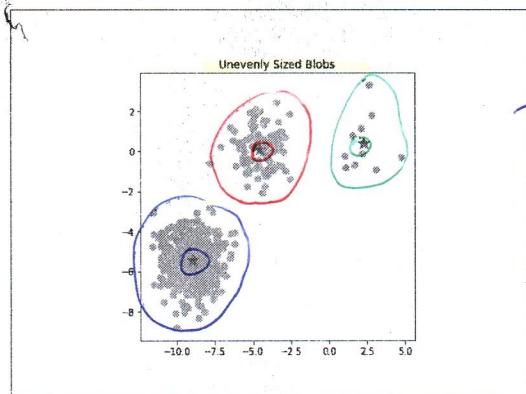
Incorrect Number of Blobs (K)



Anisotropically Distributed Blobs



The clusters try to remain spherical



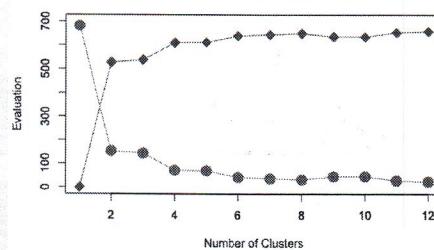
k-Means Clustering Summary	
• Strength	<ul style="list-style-type: none"> ■ Relatively efficient ■ Often terminates at a local optimum ■ The global optimum may be found using techniques such as: deterministic annealing and genetic algorithms
• Weakness	<ul style="list-style-type: none"> ■ Applicable only when mean is defined, then what about categorical data? ■ Need to specify k, the number of clusters, in advance ■ Unable to handle noisy data and outliers ■ Not suitable to discover clusters with non-convex shapes
• Advantages	<ul style="list-style-type: none"> ■ Simple, understandable ■ Items automatically assigned to clusters
• Disadvantages	<ul style="list-style-type: none"> ■ Must pick number of clusters before hand ■ All items forced into a cluster ■ Too sensitive to outliers

- A few variants of the k-means which differ in
 - Selection of the initial k means
 - Dissimilarity calculations
 - Strategies to calculate cluster means
- Handling categorical data: k-modes
 - Replacing means of clusters with modes
 - Using new dissimilarity measures to deal with categorical objects
 - Using a frequency-based method to update modes of clusters
- A mixture of categorical and numerical data: k-prototype method

How do we choose k?

(internal measure of performance)

Within/Between Cluster Sum-of-square



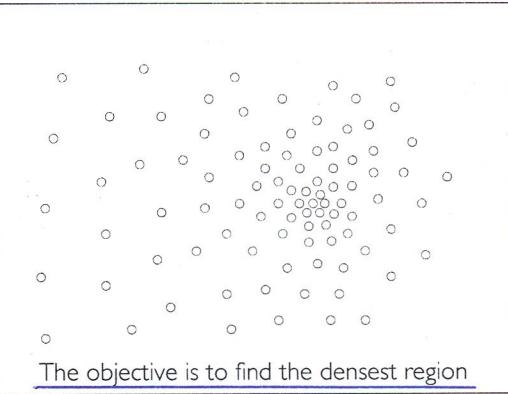
Mean Shift Clustering

Mean shift is an iterative, nonparametric, and versatile algorithm

It searches for the mode (i.e., the point of highest density) of a data distribution

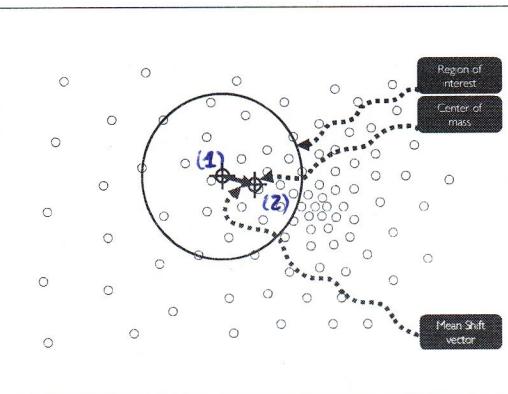
Not really a clustering procedure. This finds points with maximal density (clouds with maximal density)

Underlying intuition

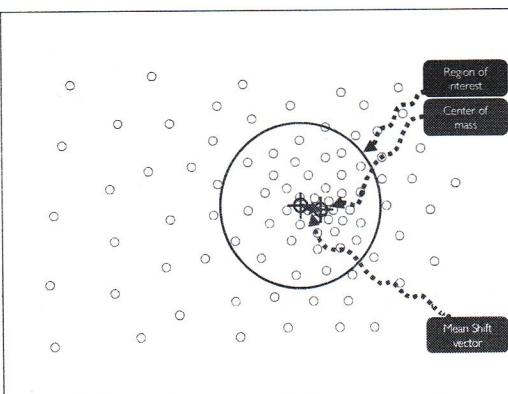
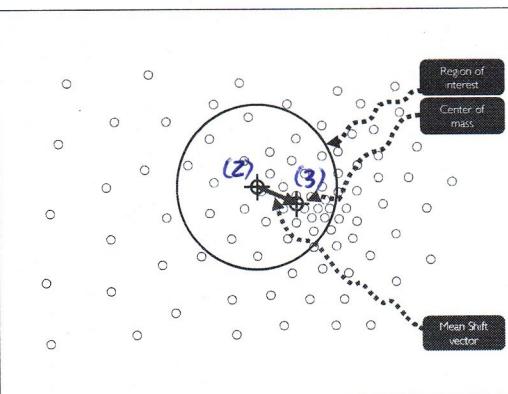


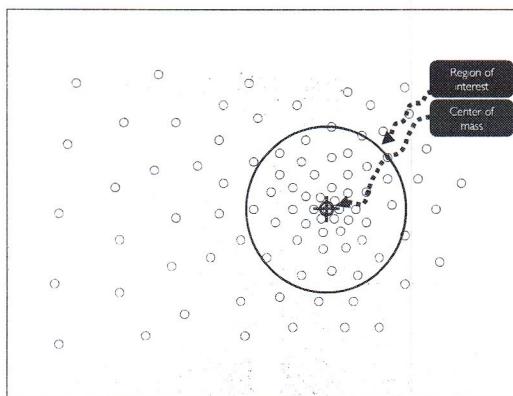
The objective is to find the densest region

We take a random point (1) and we consider its surroundings. We move to the center of mass of the whole surroundings:
(1) \rightarrow (2)

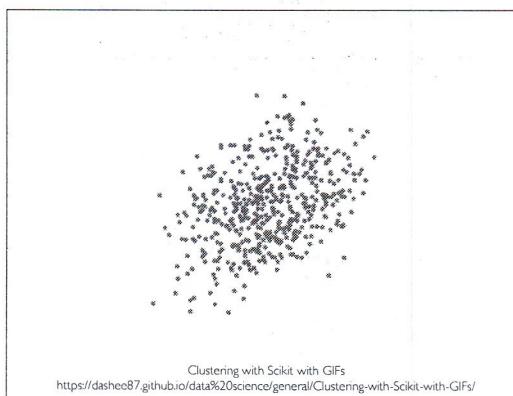


Again : (2) \rightarrow (3)





At some point the mean shift vector will be extremely small
→ we stop (and we found the point with the highest density)



Mean Shift Algorithm Pseudocode

1. Choose a search window size (the bandwidth)
2. Choose the initial location of the search window
3. Compute the mean location in the search window
4. Center the search window at the location computed in Step 3
5. Repeat Steps 3 and 4 until convergence

73

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Non-Parametric Density Estimation

- We assume that the data points are sampled from an underlying probability density function (PDF)

Assumed Underlying PDF Real Data Samples

74

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Kernel Density Estimation - Examples

- Epanechnikov Kernel

$$K_E(x) = \begin{cases} c(1 - ||x||^2) & ||x|| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$
- Uniform Kernel

$$K_U(x) = \begin{cases} c & ||x|| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$
- Normal/Gaussian Kernel

$$K_N(x) = c \cdot e^{-0.5||x||^2}$$

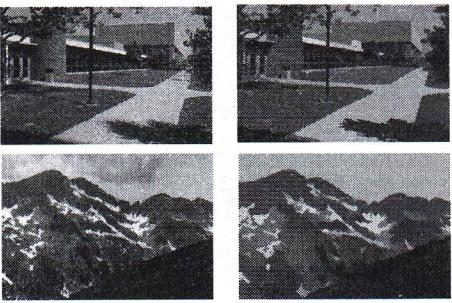
$$K_G\left(\frac{x - x_i}{h}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-x_i)^2}{2h^2}}$$

75

Prof. Pierluca Lanz | POLITECNICO DI MILANO

When we have to evaluate the center of mass we can actually use other methods (we don't necessarily have to use the mean)

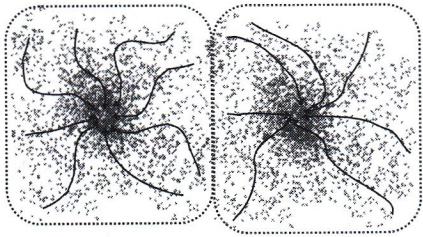
Mean Shift Segmentation
(kind of image-based clustering)



D. Comaniciu and P. Meer, Mean Shift: A Robust Approach toward Feature Space Analysis, PAMI 2002.

Mean Shift Clustering

Trajectories that lead to the same mode
(attraction basin) should be in the same cluster



Clustering with Scikit with GIFs
<https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/>

Mean Shift Clustering - Pseudocode

81

1. Choose kernel and bandwidth
2. For each point
 - a. Center a window on that point
 - b. Compute the mean of the data in the search window
 - c. Center the search window at the new mean location
 - d. Repeat (b,c) until convergence
3. Assign points that lead to nearby modes to the same cluster

Prof. Perucà Lanz

POLITECNICO DI MILANO

K-Means & Mean Shift Clustering

82

- Number of clusters and shape
 - k-means assumes that the number of clusters is already known and the clusters are shaped spherically (or elliptically)
 - Mean shift does not assume anything about number of clusters. Instead, the number of modes give the number of clusters
 - Also, since mean shift is based on density estimation, it can handle arbitrarily shaped clusters.
- Initialization and outliers
 - k-means is sensitive to initialization
 - Mean shift is robust to initializations. Typically, mean shift is run for each point or sometimes points are selected uniformly from the feature space
 - Similarly, k-means is sensitive to outliers while Mean Shift is less sensitive.
- Time complexity
 - k-means is fast and has a time complexity $O(knT)$ where k is the number of clusters, n is the number of points and T is the number of iterations.
 - Classic mean shift is computationally expensive with a time complexity $O(Tn^2)$.

Prof. Perucà Lanz

POLITECNICO DI MILANO

Mean Shift Clustering

83

- Strength
 - Does not assume any prior shape on data clusters
 - Can handle arbitrary feature spaces
 - Only the window size h to choose
 - The window size h (or bandwidth) has a physical meaning
- Weaknesses
 - The window size h (or bandwidth) selection is not trivial
 - Inappropriate window size can cause modes to be merged, or generate additional "shallow" modes
 - Adaptive window size can help
 - Not suitable for high-dimensional features

Prof. Perucà Lanz

POLITECNICO DI MILANO

Computational Complexity

84

- MeanShift clustering runs at $O(Tn^2)$, k-Means at $O(knT)$ where T is the number of iterations and n is the number of data points
- The functions to estimate the bandwidth scale particularly badly (see sklearn estimate_bandwidth documentation)

Prof. Perucà Lanz

POLITECNICO DI MILANO

Expectation Maximization

Expectation-Maximization (EM) Clustering

86

- k-means assigns each point to only one cluster (hard assignment)
- The approach can be extended to consider soft assignment of points to clusters, so that each point has a probability of belonging to each cluster
- We assume that each cluster C_i is characterized by a multivariate normal distribution and thus identified by
 - The mean vector μ_i
 - The covariance matrix Σ_i
- A clustering is identified by a vector of parameter θ defined as

$$\theta = \{\mu_i, \Sigma_i, P(C_i)\}$$

where $P(C_i)$ are the prior probability of all the clusters C_i , which sum up to one

Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

Expectation-Maximization (EM) Clustering

87

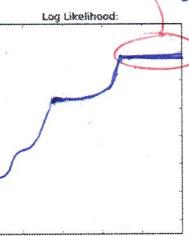
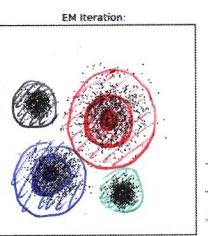
- The goal of maximum likelihood estimation (MLE) is to choose the parameters θ that maximize the likelihood, that is

$$\theta^* = \arg \max_{\theta} P(D|\theta)$$

- General idea
 - Starts with an initial estimate of the parameter vector
 - Iteratively rescores the patterns against the mixture density produced by the parameter vector
 - The rescored patterns are used to update the parameter updates
 - Patterns belonging to the same cluster, if they are placed by their scores in a particular component

Prof. Pier Luca Lanzi

POLITECNICO DI MILANO



Clustering with Scikit with GIFs
<https://dashee87.github.io/data%20science/general/Clustering-with-Skikit-with-GIFs/>

Example in One Dimension

EM in One Dimension - Initialization

90

- Consider a dataset consisting of a single attribute $X = \{x_1, \dots, x_n\}$
 - Clusters will be represented using univariate normal,
- $$f_i(x) = f(x|\mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left\{-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right\}$$
- Each cluster will be represented by the parameters $\{\mu_i, \sigma_i^2, P(C_i)\}$
 - We initialize the parameters as follows, the mean μ_i is selected uniformly at random, σ_i^2 is initialized to 1, the cluster $P(C_i)$ probability with $1/k$ (where k is the number of clusters).

Prof. Pier Luca Lanzi

POLITECNICO DI MILANO

EM in One Dimension – Expectation Step

91

- For each cluster, $1 \dots k$, we can use the current estimate of the parameters $\{\mu_i \ \sigma_i^2 \ P(C_i)\}$ to compute the posterior probabilities,

$$P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$

- We denote with

w_{ij} the values $P(C_i|x_j)$ = probability of point j to belong to cluster i
 w_i the weight vector for cluster C_i over all the n points,
 that is, $(w_{i,1} \dots w_{i,n})^T$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

EM in One Dimension – Maximization Step (1)

92

- Given all the posterior probability values w_{ij} the maximization step computes the maximum likelihood estimates of the cluster parameters $\{\mu_i \ \sigma_i^2 \ P(C_i)\}$

- The means are updated as,

$$\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot x_j}{\sum_{j=1}^n w_{ij}}$$

which can be rewritten as

$$\mu_i = \frac{\mathbf{w}_i^T \mathbf{X}}{\mathbf{w}_i^T \mathbf{1}}$$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

EM in One Dimension – Maximization Step (2)

93

- Similarly, for σ_i^2

$$\sigma_i^2 = \frac{\sum_{j=1}^n w_{ij} (x_j - \mu_i)^2}{\sum_{j=1}^n w_{ij}} \quad \text{or} \quad \sigma_i^2 = \frac{\mathbf{w}_i^T \mathbf{Z}_i^2}{\mathbf{w}_i^T \mathbf{1}}$$

where $\mathbf{Z}_i = (x_1 - \mu_i, \dots, x_n - \mu_i)^T$

- Finally,

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{\sum_{a=1}^k \sum_{j=1}^n w_{aj}} = \frac{\sum_{j=1}^n w_{ij}}{\sum_{j=1}^n 1} = \frac{\sum_{j=1}^n w_{ij}}{n}$$

- The expectation and maximization steps are repeated until convergence (e.g. until the means change a little between steps)

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Example with one variable
 (Example 13.4 from the textbook)

Example 13.4 from the Textbook (1)

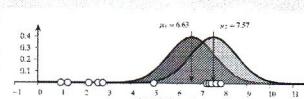
95

- Data Points

$x_1 = 1.0$	$x_2 = 1.3$	$x_3 = 2.2$	$x_4 = 2.6$	$x_5 = 2.8$
$x_6 = 5.0$	$x_7 = 7.3$	$x_8 = 7.4$	$x_9 = 7.5$	$x_{10} = 7.7$
$x_{11} = 7.9$				

- Initialization of the parameters

$\mu_1 = 6.63$	$\sigma_1^2 = 1$	$P(C_1) = 0.5$
$\mu_2 = 7.57$	$\sigma_2^2 = 1$	$P(C_2) = 0.5$



Prof. Pierluca Lanzi

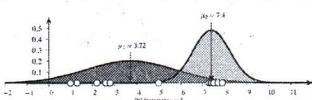
POLITECNICO DI MILANO

Example 13.4 from the Textbook (2)

96

- After one expectation-maximization iteration the parameters are

$$\begin{array}{lll} \mu_1 = 3.72 & \sigma_1^2 = 6.13 & P(C_1) = 0.71 \\ \mu_2 = 7.4 & \sigma_2^2 = 0.69 & P(C_2) = 0.29 \end{array}$$



Prof. Pierluca Lanzi

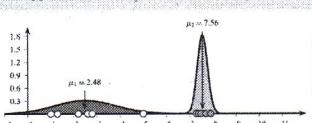
POLITECNICO DI MILANO

Example 13.4 from the Textbook (3)

97

- The procedure converges after five iterations to

$$\begin{array}{lll} \mu_1 = 2.48 & \sigma_1^2 = 1.69 & P(C_1) = 0.55 \\ \mu_2 = 7.56 & \sigma_2^2 = 0.05 & P(C_2) = 0.45 \end{array}$$



- EM Clustering returns probabilities, it does not assign points to clusters.
- In this example, we assigned elements to the cluster with the highest posterior probability

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

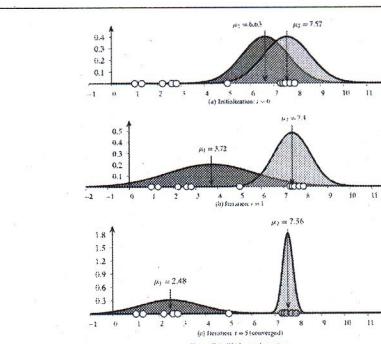
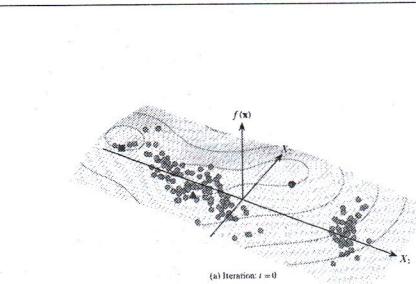


Figure 13.4: EM in one dimension

Notice: the number of clusters (k) is defined by the user

Example from Chapter 13 of the textbook

Example with two variables



(a) Iteration: $i = 0$

Figure 13.5 from Chapter 13 of the textbook

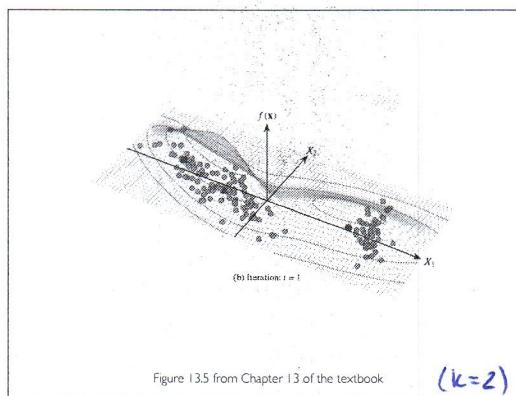


Figure 13.5 from Chapter 13 of the textbook

($k=2$)

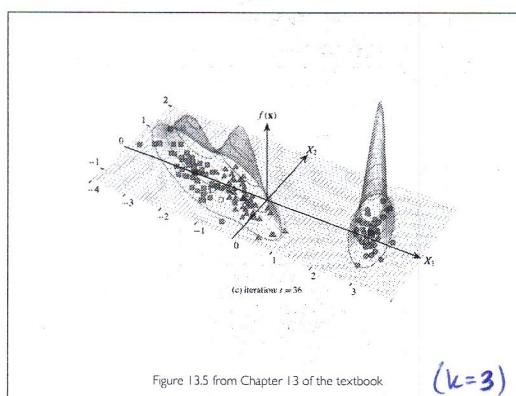
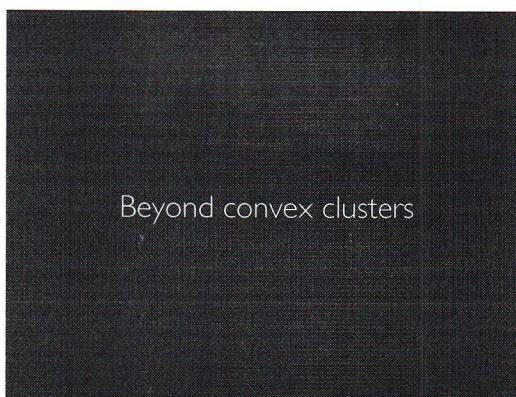


Figure 13.5 from Chapter 13 of the textbook

($k=3$)

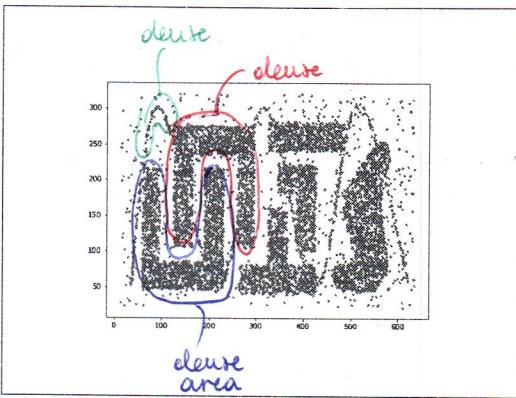


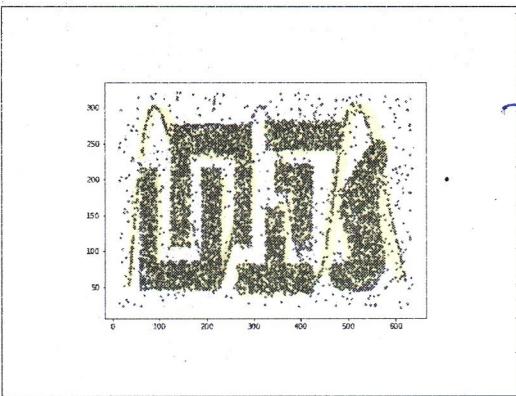
Representative-based clustering methods are suitable for finding ellipsoid-shaped clusters, or at best convex clusters

For nonconvex clusters, these methods have trouble finding the true clusters

Density-based methods can mine such nonconvex clusters

The idea is to find dense areas





We consider all the non-yellow points as noise, so we look for clusters only within the yellows

What is density-based clustering? 107

- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan
 - Need density parameters as termination condition

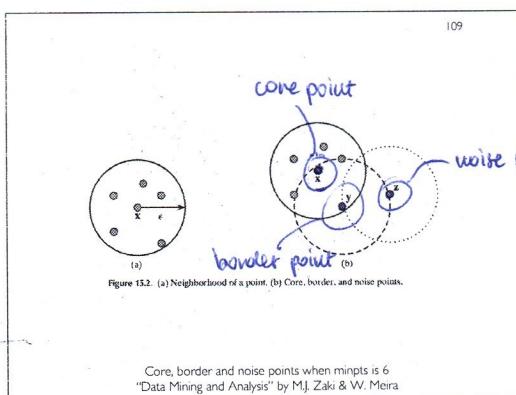
Prof. Pierluca Lanzi POLITECNICO DI MILANO

The user defines the definition of density (different meanings for densities lead to different results)

DBSCAN: Basic Concepts 108

- The neighborhood within a radius ϵ of a given object is called the ϵ -neighborhood of the object.
 $N_\epsilon(x) = \{y | \delta(x, y) \leq \epsilon\}$
- **Core Point (points in a dense area)**
 - If the ϵ -neighborhood of an object contains at least minpts objects, then the object is a core object
- **Border Point (points nearby a dense area)**
 - If its ϵ -neighborhood does not contain at least minpts, but it is inside the neighborhood of a core point
- **Noise Point (isolated points)**
 - If its ϵ -neighborhood does not contain at least minpts and it does not belong to the neighborhood of a core point

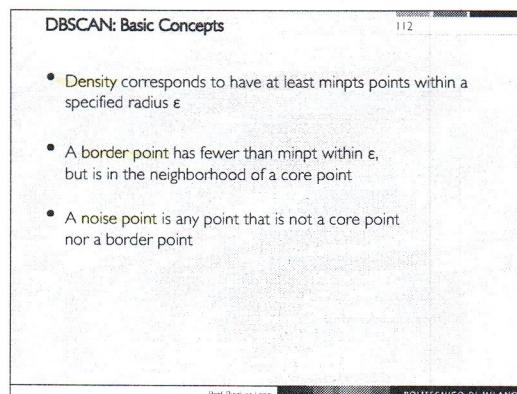
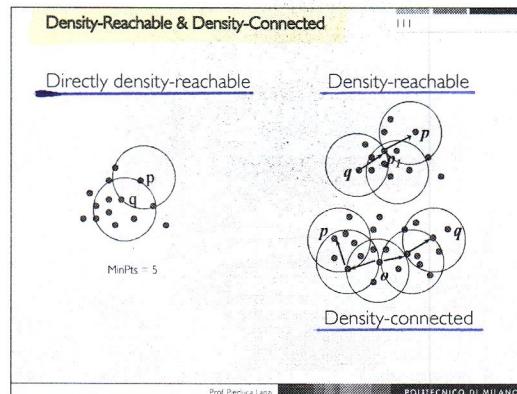
Prof. Pierluca Lanzi POLITECNICO DI MILANO



DBSCAN: Basic Concepts 110

- **Directly density reachable**
 - An object x is directly density-reachable from object y if x is within the ϵ -neighborhood of y and y is a core object
- **Density Reachable**
 - An object x is density-reachable from object y if there is a chain of objects x_1, \dots, x_n where $x_1=y$ and $x_n=x$ such that x_{i+1} is directly density reachable from x_i
- **Density Connected**
 - An object p is density-connected to q with respect to ϵ and MinPts if there is an object o such that both p and q are density reachable from o
- **Density-Based Cluster**
 - A density-based cluster is defined as a maximal set of density connected points.

Prof. Pierluca Lanzi POLITECNICO DI MILANO



ALGORITHM 15.1. Density-based Clustering Algorithm

```

DBSCAN( $D, \epsilon, \text{minpts}$ ):
1  $\text{Core} \leftarrow \emptyset$ 
2 foreach  $x_i \in D$  do // Find the core points
3   | Compute  $N_\epsilon(x_i)$ 
4   |  $id(x_i) \leftarrow \emptyset$  // cluster id for  $x_i$ 
5   | if  $|N_\epsilon(x_i)| \geq \text{minpts}$  then  $\text{Core} \leftarrow \text{Core} \cup \{x_i\}$ 
6    $k \leftarrow 0$  // cluster id
7   foreach  $x_i \in \text{Core}$ , such that  $id(x_i) = \emptyset$ 
8     |  $k \leftarrow k + 1$ 
9     |  $id(x_i) \leftarrow k$  // assign  $x_i$  to cluster id  $k$ 
10    | DENSITYCONNECTED( $x_i, k$ )
11  $C \leftarrow \{C_i\}_{i=1}^k$ , where  $C_i \leftarrow \{x \in D \mid id(x) = i\}$ 
12  $Noise \leftarrow \{x \in D \mid id(x) = \emptyset\}$ 
13  $Border \leftarrow D \setminus (\text{Core} \cup Noise)$ 
14 return  $C, Core, Border, Noise$ 

DENSITYCONNECTED( $x, k$ ):
15 foreach  $y \in N_\epsilon(x)$  do
16   |  $id(y) \leftarrow k$  // assign  $y$  to cluster id  $k$ 
17   | if  $y \in \text{Core}$  then DENSITYCONNECTED( $y, k$ )

```

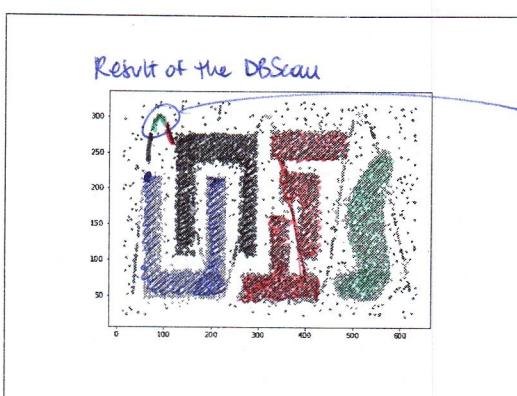
Density Based Scan

Complexity of DBScan

114

- DBSCAN needs to compute the ϵ -neighborhood for each point
- If the dimensionality is not too high this can be done efficiently using a spatial index structure in $O(n \log n)$
- If the dimensionality is high, it takes $O(n^2)$ to compute the neighborhood for each point.
- Once the ϵ -neighborhood has been computed the algorithm needs only a single pass over all the points to find the density connected clusters.
- The overall complexity of DBSCAN is between $O(n \log n)$ and in the worst-case.

Prof. Pierluca Lanzi POLITECNICO DI MILANO



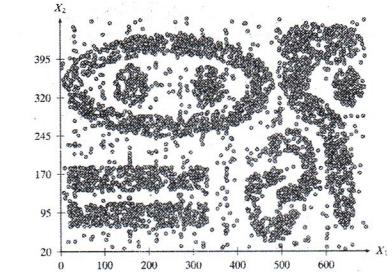


Figure 15.1 from Chapter 15 of the textbook "Data Mining and Analysis" by M.J. Zaki & W. Meira

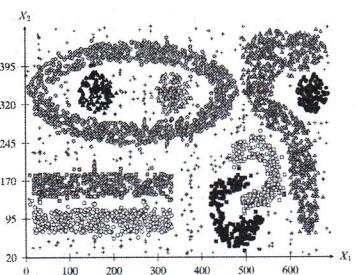


Figure 15.3 from Chapter 15 of the textbook "Data Mining and Analysis" by M.J. Zaki & W. Meira

Let's play with DBSCAN

<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

DBScan fails when applied
to data of varying density !

HDBSCAN

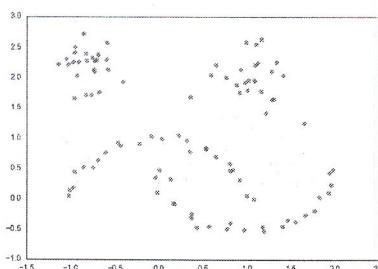
→ Hierarchical clustering
+ Density Based clustering

- Extends DBSCAN by converting it into a hierarchical clustering algorithm, and then using a technique to extract a flat clustering based on the stability of clusters
- It defines the "core distance" of a point as the maximum distance of a point to its k-th nearest neighbor or $\text{core}_k(x)$
- Thus, points in high density areas will have a low core distance, while point in low density areas will have a high core distance
- This provides a local inexpensive estimate of density

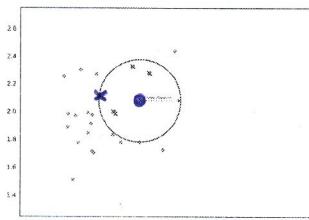
Mutual Reachability Distance

- This new distance metric is defined as,
 $d_{\text{reach}-k}(a, b) = \max\{\text{core}_k(a), \text{core}_k(b), d(a, b)\}$
using the core distance and the original distance metrics d
- This distance keeps dense points with lower core distance at the same distance while pushing away sparser points (at least at their core distance from any other point).
- Note that as with DBSCAN high values of k focus the algorithm toward very dense region

The goal is to remap the points in another space where dense points are kept at the same distance while sparse points are pushed away



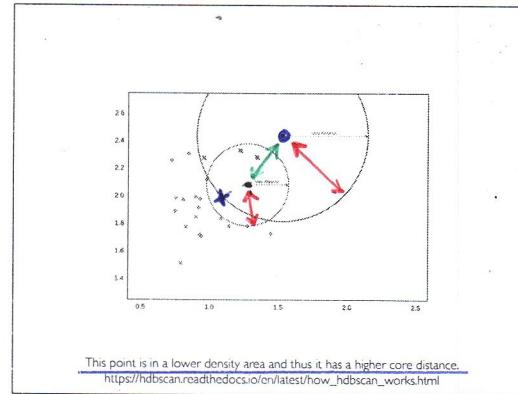
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html



Let's consider a k equal to 5 and plot the core distance for one point.
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

We consider the point and the 5 closest to it. The core distance of this point is the maximum of the 5 distances (the distance of the point with the far-ext point of the 5).

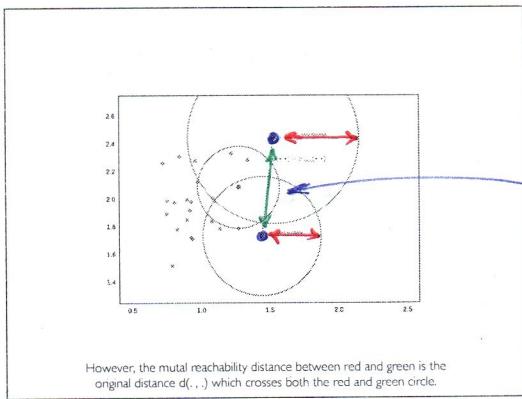
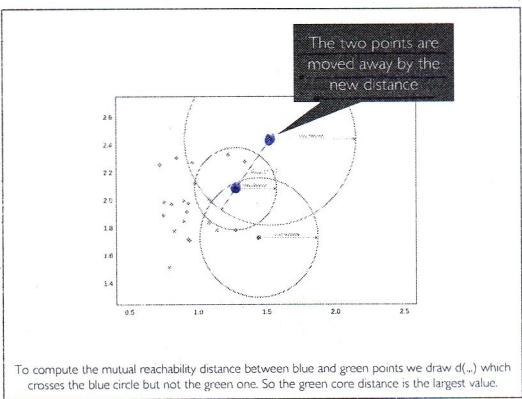
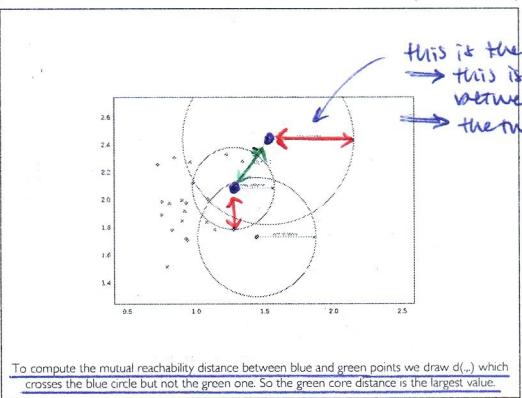
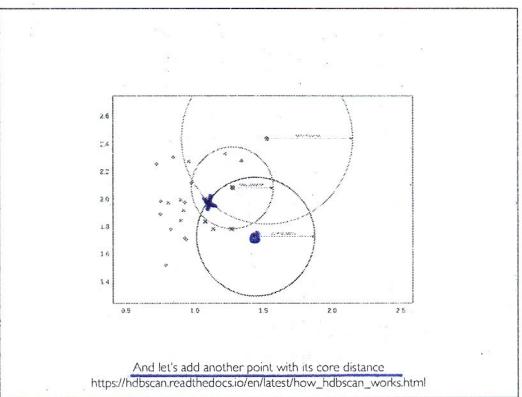
dense area \rightarrow small core dist.



This point is in a lower density area and thus it has a higher core distance.
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

(MUTUAL REACHABILITY DISTANCE)

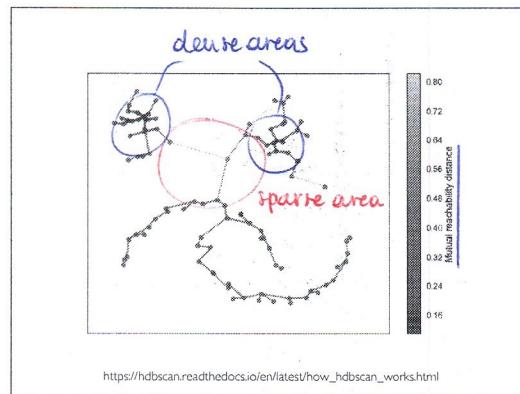
Now that we have two points. The reachability is the maximum between three distances : the two cores distances and the (original) distance between the two points



Next, Build the Minimum Spanning Tree using Prim's Algorithm 130

- Consider the data as a weighted graph
- Data points are vertices
- An edge between any two points with weight equal to the mutual reachability distance of those points
- Build the spanning tree one edge at a time by adding the lowest weight edge that connects the current tree to a vertex not yet in the tree

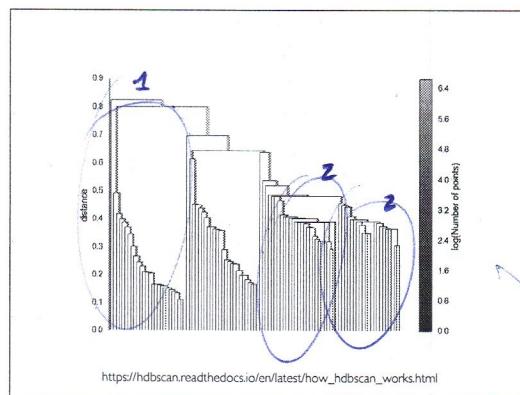
Once we know which points to push and which not



HDBSCAN converts the minimum spanning tree into a hierarchy of connected components

It sorts the edges of the tree by distance
(in increasing order)

Then iterate creating a new merged cluster for each edge.



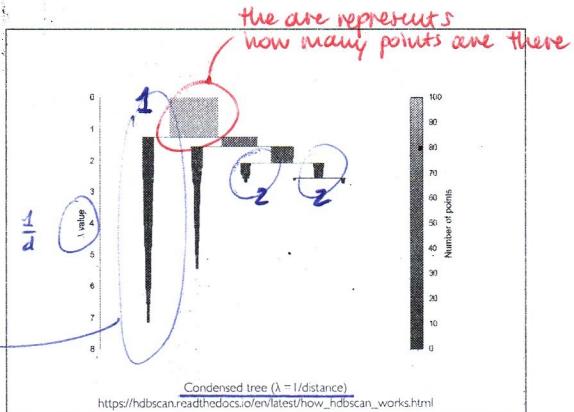
We have a hierarchy now
we want a set of flat clusters

- Condensing the Tree
- We introduce the notion of "minimum cluster size" threshold
 - We navigate the hierarchy and from the top and at each split we check the size of the merged clusters
 - If a cluster has fewer points than the threshold, then the smaller cluster is eliminated (the points have "fallen out of the cluster")
 - Otherwise both clusters are maintained.
 - Note that, which a point p 'fallen out of the cluster' and at what distance value that happened is stored as $\lambda_p = 1/distance$
 - At the end, a much smaller tree with a small number of nodes remains with node has data about how the size of the cluster at that node decreases over varying distance
- 135
- Prof. Pierluigi Lenzi
POLITECNICO DI MILANO

* This difference between the different meanings can be synthetized in the next (condensed) dendrogram

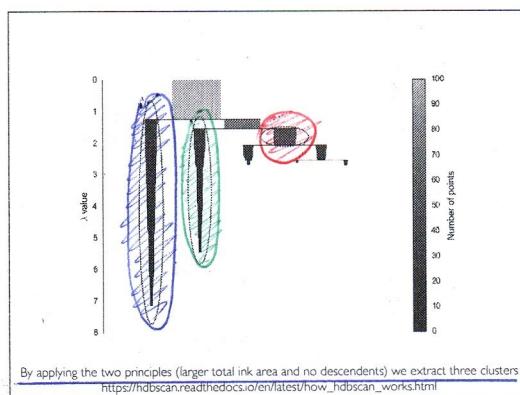
The idea: merges between clusters that are small are not so interesting. We want to find structures that last for a long time (structures that last during a lot of distances). In the structure 1 we see a structure that lasts a lot (we keep adding points to an existing structure). In the structure 2 we see a "desperate attempt" of merging points (small clusters merged without a proper criterium).

This is a clear representation of a structure that is present for a long time (unlike z)



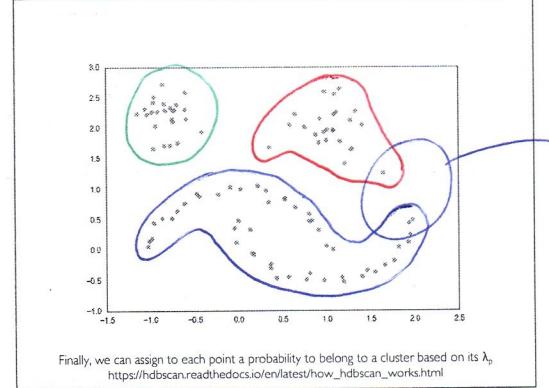
Extracting the Cluster

- Given the simplified dendrogram we want to select clusters that persist and have a longer lifetime since short lived clusters probably represent artifacts
- We may say that we want to choose those clusters that have the greatest area of ink in the plot. Also, we cannot select any a descendant of a cluster we chose.



Finally, we can assign to each point a probability to belong to a cluster based on its λ_p

HDBSCAN complexity is comparable and lower than DBSCAN



notice that these points are colored with a color which is in between red and blue since the color represents the probability of belonging to a class (is not a strict labeling)

These are only the basic ones

141

- Birch and Chameleon are other hierarchical clustering algorithms
- There are also probabilistic hierarchical clustering algorithms
- OPTICS, and DENCLU are interesting density-based algorithms
- STING and CLIQUE are grid-based methods, where CLIQUE is also a subspace clustering algorithm
- ...

Prof. Pierluca Lanzi
POLITECNICO DI MILANO

K-Means shouldn't be your first choice

SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>

K-Means probably shouldn't be your second choice either

SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>

Think hard about what "cluster" means for your application

SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>

Study Material

145

- "Data Mining and Analysis" by Zaki & Meira
 - Chapter 13
 - Chapter 15
- <http://www.dataminingbook.info>
- https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

Prof. Pierluca Lanzi
POLITECNICO DI MILANO