

Multi-dimensional index structures

- ... cluster objects when the individual dimensions (features) are known...
- ...what if we do not have explicit features??

Maria Luisa Sapino (BDM 2018)

Multi-dimensional index structures

- ... cluster objects when the individual dimensions (features) are known...
- ...what if we do not have explicit features??
 - ...we may have a “black-box program” which can compare two objects

Maria Luisa Sapino (BDM 2018)

Multi-dimensional index structures

- ... cluster objects when the individual dimensions (features) are known...
- ...what if we do not have explicit features??
 - ...we may have a “black-box program” which can compare two objects
 - ...we may have “users” evaluating the similarity of objects

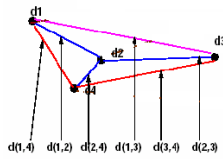
Maria Luisa Sapino (BDM 2018)

Multi-dimensional index structures

- ... cluster objects when the individual dimensions (features) are known...
- ...what if we do not have explicit features??
 - ...we may have a "black-box program" which can compare two objects
 - ...we may have "users" evaluating the similarity of objects
- we need a clustering scheme that does not need explicit features!!!!

Maria Luisa Sapino (BDM 2018)

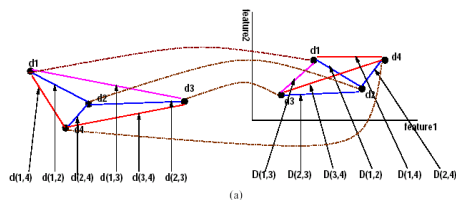
What if we do not have features??



- We know the distances, but
- we do not have explicit features
 - distances are not metric...

Maria Luisa Sapino (BDM 2018)

Multi Dimensional Scaling



$$\begin{array}{lll}
 d(1, 2) \approx D(1, 2) & d(1, 3) \approx D(1, 3) & d(1, 4) \approx D(1, 4) \\
 d(2, 3) \approx D(2, 3) & d(2, 4) \approx D(2, 4) & d(3, 4) \approx D(3, 4)
 \end{array}$$

Maria Luisa Sapino (BDM 2018)

MDS

- The criterion for the mapping is to minimize *stress*

$$stress = \sqrt{\frac{\sum_{i,j} (d_{ij} - d'_{ij})^2}{\sum_{i,j} d_{ij}^2}}$$

- Start with a (random) configuration of points with low dimensions
- Apply some form of steepest descent iteratively to minimize the stress.
 - move objects
 - add dimensions

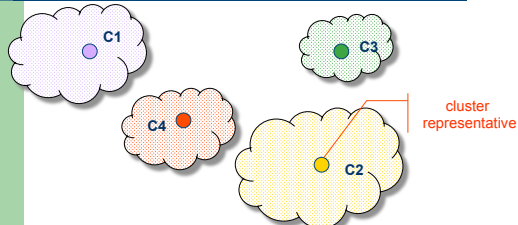
Maria Luisa Sapino (BDM 2018)

MDS

- How to map the query?

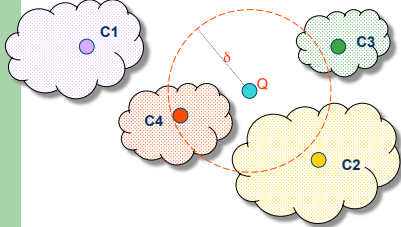
Maria Luisa Sapino (BDM 2018)

Use of clusters (prune search space)



Maria Luisa Sapino (BDM 2018)

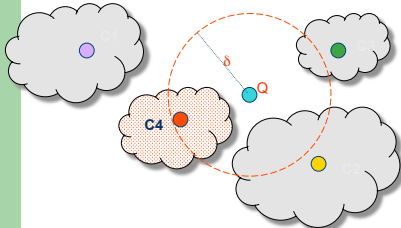
Use of clusters (prune search space)



- ...given a query

Maria Luisa Sapino (BDM 2018)

Use of clusters (prune search space)



- ...eliminate clusters based on their representatives

Maria Luisa Sapino (BDM 2018)

Clustering methods

- Sound methods:
 - need a fixed document-to-document similarity matrix
- Iterative methods:
 - use document vectors iteratively

Maria Luisa Sapino (BDM 2018)

Outline of sound methods

- Find the similarity of each object pair

Maria Luisa Sapino (BDM 2018)

Outline of sound methods

- Find the similarity of each object pair
- Setup a threshold
 - $\text{sim}(o1,o2) < T$ (objects are very different!)
 - $\text{sim}(o1,o2) \geq T$ (objects are comparable)

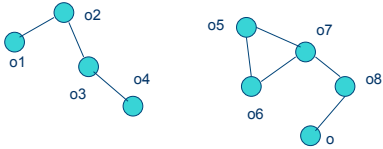
Maria Luisa Sapino (BDM 2018)

Outline of sound methods

- Find the similarity of each object pair
- Setup a threshold
 - $\text{sim}(o1,o2) < T$ (objects are very different!)
 - $\text{sim}(o1,o2) \geq T$ (objects are comparable)
- Create a graph which represents object similarities
 - Each pair of objects that are comparable is connected with an edge

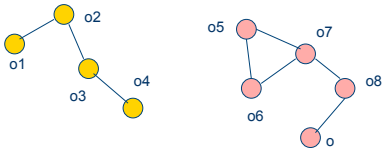
Maria Luisa Sapino (BDM 2018)

Example collection



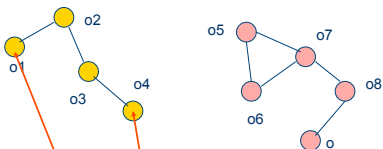
Maria Luisa Sapino (BDM 2018)

Connected components



Maria Luisa Sapino (BDM 2018)

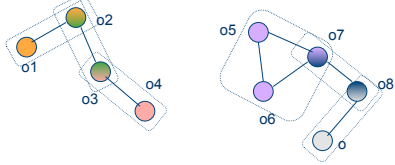
Connected components



- ...are o1 and o4 really similar????

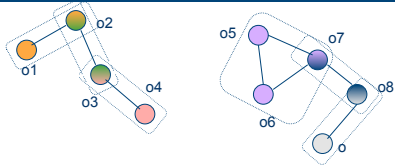
Maria Luisa Sapino (BDM 2018)

Clique...



Maria Luisa Sapino (BDM 2018)

Clique...



- Clusters are overlapping!
- Costlier to compute (NP-complete)

Maria Luisa Sapino (BDM 2018)

...or single-pass iterative method

- choose an object, and make it a cluster

Maria Luisa Sapino (BDM 2018)

...or single-pass iterative method

- choose an object, and make it a cluster
- choose another object, o
 - find the closest cluster, c
 - if $\text{dist}(o,c) < T$ add o to c
 - else o is a new cluster

Maria Luisa Sapino (BDM 2018)

...or single-pass iterative method

- choose an object, and make it a cluster
- choose another object, o
 - find the closest cluster, c
 - if $\text{dist}(o,c) < T$ add o to c
 - else o is a new cluster
- repeat until all objects are processed

Maria Luisa Sapino (BDM 2018)

...or single-pass iterative method

- choose an object, and make it a cluster
- choose another object, o
 - find the closest cluster, c
 - if $\text{dist}(o,c) < T$ add o to c
 - else o is a new cluster
- repeat until all objects are processed

also called the "leader" algorithm

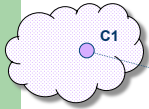
Maria Luisa Sapino (BDM 2018)

....how do we compute distance?



Maria Luisa Sapino (BDM 2018)

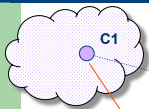
....how do we compute distance?



- use a cluster representative...

Maria Luisa Sapino (BDM 2018)

....how do we compute distance?

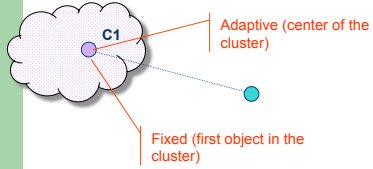


Fixed (first object in the cluster)

- use a cluster representative...

Maria Luisa Sapino (BDM 2018)

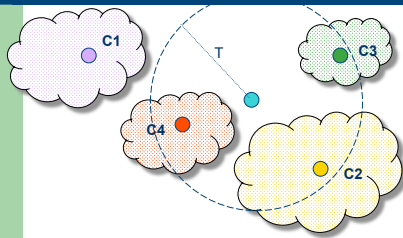
....how do we compute distance?



- use a cluster representative...

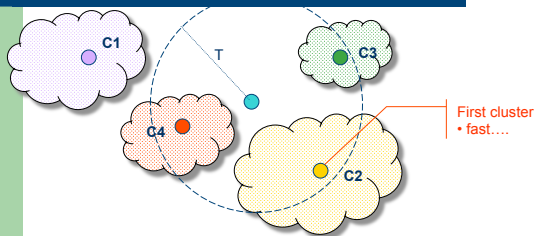
Maria Luisa Sapino (BDM 2018)

....alternatives...



Maria Luisa Sapino (BDM 2018)

....alternatives...



Maria Luisa Sapino (BDM 2018)

....alternatives...

Closest cluster
• compact clusters....

Maria Luisa Sapino (BDM 2018)

....alternatives...

smallest cluster
• spread the wealth
• high entropy

Maria Luisa Sapino (BDM 2018)

What if we do not have a threshold??

- Find a minimum spanning tree of the input graph
 - $O(N^2)$ edges to $O(N)$ edges

Maria Luisa Sapino (BDM 2018)

What if we do not have a threshold??

- Find a minimum spanning tree of the input graph
 - $O(N^2)$ edges to $O(N)$ edges
- Remove all edges longer than the average of their neighbors
 - threshold is determined based on the neighborhood

Maria Luisa Sapino (BDM 2018)

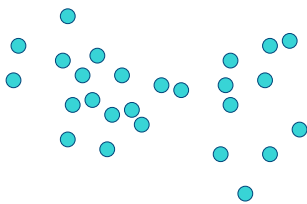
What if we do not have a threshold??

- Find a minimum spanning tree of the input graph
 - $O(N^2)$ edges to $O(N)$ edges
- Remove all edges longer than the average of its neighbors
 - threshold is determined based on the neighborhood
- Apply connected-components or clique..

Maria Luisa Sapino (BDM 2018)

Max-a-min

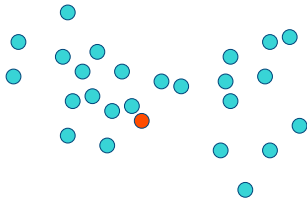
- The number of clusters is known, r (say 4)



Maria Luisa Sapino (BDM 2018)

Max-a-min

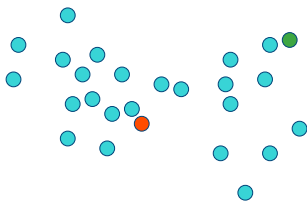
- Choose a random leader



Maria Luisa Sapino (BDM 2018)

Max-a-min

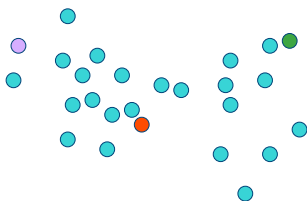
- Find the furthest point to the leader



Maria Luisa Sapino (BDM 2018)

Max-a-min

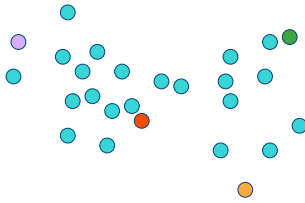
- Find the furthest point to the 2 leaders



Maria Luisa Sapino (BDM 2018)

Max-a-min

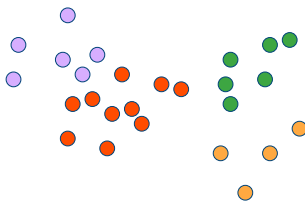
- Find the furthest point to the 3 leaders



Maria Luisa Sapino (BDM 2018)

Max-a-min

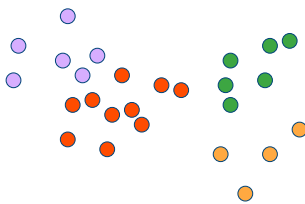
- Assign points to closest leader..



Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

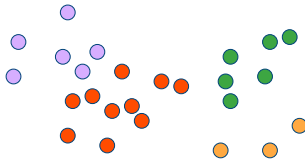
- Minimize a “global cost function”



Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

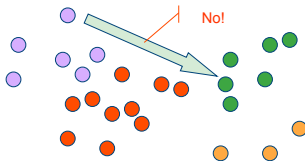


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

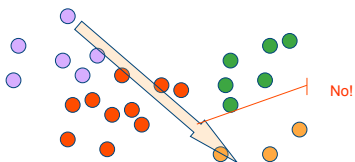


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

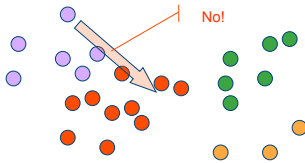


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

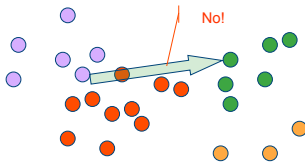


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

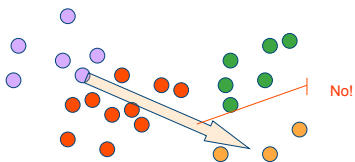


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

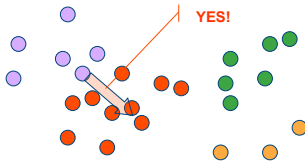


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

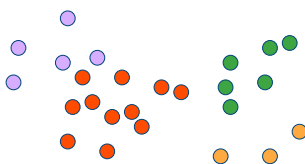


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”

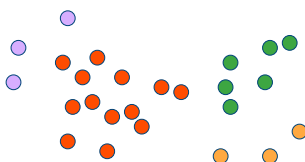


- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

K-means (iterative improvement)

- Minimize a “global cost function”



- ...each item is checked whether moving to another cluster would reduce global cost

Maria Luisa Sapino (BDM 2018)

What are the possible criteria??

- Compactness (minimize root-mean-square)



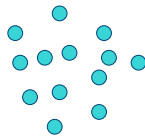
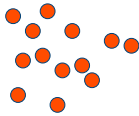
$$RMSE_m = \sqrt{\frac{1}{N} \sum_{i=0}^N [\hat{F}_m(z_i) - F_m(z_i)]^2}$$

$$RMSE = \frac{1}{M} \sum_{m=1}^M RMSE_m$$

Maria Luisa Sapino (BDM 2018)

What are the possible criteria??

- Evenly sized clusters (maximize entropy)



$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

Maria Luisa Sapino (BDM 2018)

What if we do not have distances???

- ...we need to learn from
 - user feedback or
 - user access patterns!!

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assumption:
 - each cluster can have maximum r objects
- ...keeps adapting the clusters to user access pattern

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a and o_b are in the same cluster C_j then
 - $\text{conf}(a,j)++$
 - $\text{conf}(b,j)++$

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a is in cluster C_i and o_b is in cluster C_j then

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a is in cluster C_i and o_b is in cluster C_j then
 - if $\text{conf}(a,i) > 1$ and $\text{conf}(b,j) > 1$ then
 - $\text{conf}(a,i)--$
 - $\text{conf}(b,j)--$

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a is in cluster C_i and o_b is in cluster C_j then
 - if $\text{conf}(a,i) == 1$ and $\text{conf}(b,j) == 1$ then
 - $\text{conf}(b,i) = 1$ (move o_b from C_j to C_i)
 - $\text{conf}(c,j) = 1$ (move some o_c from C_i to C_j)

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a is in cluster C_i and o_b is in cluster C_j then
 - if $\text{conf}(a,i) > 1$ and $\text{conf}(b,j) == 1$ then
 - there exists $\text{conf}(c,i) == 1$
 - $\text{conf}(b,i) = 1$ (move o_b from C_j to C_i)
 - $\text{conf}(c,j) = 1$ (move o_c from C_i to C_j)

Maria Luisa Sapino (BDM 2018)

Confidence clustering...

- Assign each object o_i to cluster C_j with random confidence ($0 \leq \text{conf}(i,j) \leq 10$)
- If user accesses o_a and o_b , then
 - if o_a is in cluster C_i and o_b is in cluster C_j then
 - if $\text{conf}(a,i) > 1$ and $\text{conf}(b,j) == 1$ then
 - there does not exist $\text{conf}(c,i) == 1$
 - $\text{conf}(a,i) --$

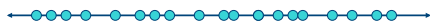
Maria Luisa Sapino (BDM 2018)

Adaptive clustering...

- What if we do not know the number of clusters????
- ...keeps adapting the clusters to user access pattern

Maria Luisa Sapino (BDM 2018)

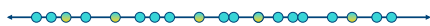
Adaptive clustering...



- start with a random assignment of objects to a line

Maria Luisa Sapino (BDM 2018)

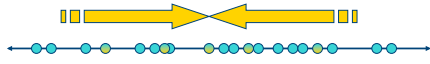
Adaptive clustering...



- If a set of objects are accessed together...

Maria Luisa Sapino (BDM 2018)

Adaptive clustering...



- ..pull the objects closer to their average point..

Maria Luisa Sapino (BDM 2018)

Adaptive clustering...



- ..choose another random set of objects..

Maria Luisa Sapino (BDM 2018)

Adaptive clustering...



- ..push these away from their average point..

Maria Luisa Sapino (BDM 2018)

Adaptive clustering...



- ...over time...
 - similar objects will come closer...
 - different objects will get apart...

Maria Luisa Sapino (BDM 2018)
