

## Partitioning Algorithms

- Partitioning method: Construct a partition of n documents into a set of K clusters
- $\Box$  Given: a set of documents and the number K
- Find: a partition of K clusters that optimizes the chosen partitioning criterion
  - Globally optimal: exhaustively enumerate all partitions
  - Effective heuristic methods: K-means and K-medoids algorithms







since each vector is assigned to the closest centroid.



### Time Complexity

- Computing distance between two docs is O(m) where m is the dimensionality of the vectors.
- Reassigning clusters: O(Kn) distance computations, or O(Knm).
- Computing centroids: Each doc gets added once to some centroid: O(nm).
- Assume these two steps are each done once for *I* iterations: O(IKnm).

Dip. di Matematica Pura ed Applicata

# Time Complexity

	So, k-means is linear in all relevant factors (iterations, number of clusters, number of documents, and dimensionality of the space)	
	But M>100.000 !!!	
	Docs are sparse but centroids tend to be dense -> distance computation is time consuming	
	Effective heuristics can be defined for making centroid-doc distance computation as efficient as doc- doc distance computation	
	K-medoids is a variant of k-means that compute medoids (the docs closest to the centroid) instead of centroids as cluster centers.	
Din	di Matematica E Aiolli - Sistemi Informativi	29

ip. di Matematica Pura ed Applicata

Pura ed Applicata

2006/2007

## Seed Choice

Results can vary based on random seed selection.	Example showing sensitivity to seeds
Some seeds can result in poor convergence rate, or convergence to sub-optimal	ABCOOCO OOCO DEF
<ul> <li>Clusterings.</li> <li>Select good seeds using a heuristic (e.g., doc least similar to any existing mean)</li> <li>Try out multiple starting points</li> <li>Initialize with the results of</li> </ul>	In the above, if you start with B and E as centroids you converge to {A,B,C} and {D,E,F} If you start with D and F you converge to {A,B,D,E} {C,F}
another method.	

2006/2007







#### Considerations

- Finding good seeds is even more critical for EM than for k-means (EM is prone to get stuck in local optima)
- Therefore (as in k-means) an initial assignment is often computed by another algorithm
- □ If the model of the data is correct EM algorithm finds the correct structure
- Hardly a document collection can be considered generated by a simple mixture model
- At least, model based clustering allows for analysis and adaptations

Dip. di Matematica Pura ed Applicata 36