Information Retrieval (Text Clustering)

Fabio Aiolli

http://www.math.unipd.it/~aiolli

Dipartimento di Matematica Pura ed Applicata Università di Padova

Anno Accademico 2009/2010

1

2

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval -2009/10

What is clustering?

Clustering: the process of grouping a set of objects into classes of similar objects

- The commonest form of unsupervised learning
 - Unsupervised learning = learning from raw data, as opposed to supervised data where a classification of examples is given
- A common and important task that finds many applications in IR and other places
- Not only Document Clustering (e.g. terms)

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval -2009/10















ground truth



$$RI = \frac{A+D}{A+B+C+D}$$

Compare with standard Precision and Recall.

 $P = \frac{A}{A+B}$

 $R = \frac{A}{A+C}$

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval -2009/10

Rand Index example: 0.68

| Number of points | Same Cluster in clustering | Different Clusters in clustering |
|---|----------------------------------|--|
| Same class in ground truth | 20 | 24 |
| Different classes in ground truth | 20 | 72 |

F. Aiolli - Information Retrieval - 2009/10

17











Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval -2009/10





□ 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 docdoc 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.

| Dip. | di Matematica | а |
|------|---------------|---|
| Pura | ed Applicata | |

F. Aiolli - Information Retrieval -2009/10

Seed Choice

| Some seeds can result in poor convergence rate, or convergence to sub-optimal | A B O O O O D E | C O O F |
|--|--|--|
| clusterings. Select good seeds using a heuristic (e.g., doc least similar to any existing mean) Try out multiple starting points Initialize with the results of another method. | In the above, if with B and E as you converge to and {D,E,F} If you start with you converge to {A,B,D,E} {C,F} | you start centroids {A,B,C} n D and F |

29









The dendrogram

The y-axis of the dendogram represents the combination similarities, i.e. the similarities of the clusters merged by a the horizontal lines for a particular y

Assumption: The merge operation is monotonic, i.e. if s₁,..,s_{k-1} are successive combination similarities, then

 $s_1 \ge s_2 \ge ... \ge s_{k-1}$ must hold

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval -2009/10

Hierarchical Agglomerative Clustering (HAC)

Starts with each doc in a separate cluster

then repeatedly joins the closest pair of clusters, until there is only one cluster.

The history of merging forms a binary tree or hierarchy.

39













Summarizing

| Single-link | Max sim of any two points | O(N ²) | Chaining effect |
|-------------------|---------------------------------|------------------------|-----------------------|
| Complete-link | Min sim of any two points | O(N ² logN) | Sensitive to outliers |
| Centroid | Similarity of centroids | O(N ² logN) | Non monotonic |
| Group- average | Avg sim of any two points | O(N ² logN) | ОК |

F. Aiolli - Information Retrieval -2009/10