

clustering unsupervised learning

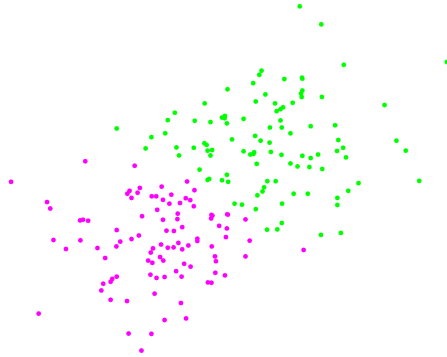
DTL SU @ AU

kristoffer l nielbo
kln@cas.au.dk
github.com/kln-courses/tmgu17
tmgu17.slack.com

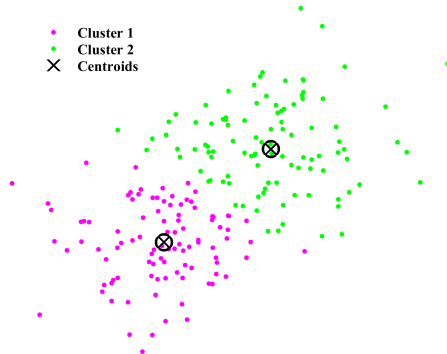
DAI|IMC|AARHUS UNIVERSITY



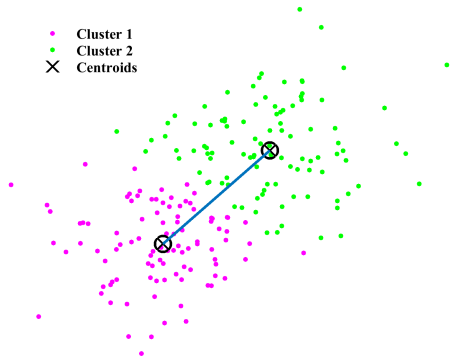
Implicit assumption that we study differences in variables (e.g., terms) between homogeneous objects (e.g., documents)



Systematic differences between objects result in non-random subsets that are often ignored



Cluster analysis: partitions data into homogeneous subsets using inter-object similarity/distance measures



Minimize distance between the centroid and points within each cluster
Maximize distance centroids and points between clusters

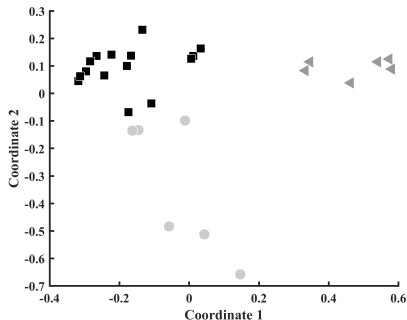
$$\begin{array}{c|cccc} & t_1 & t_2 & \dots & t_k \\ \hline d_1 & f_{d_1, t_1} & & & \\ d_2 & & f_{d_2, t_2} & & \\ \dots & & & \dots & \\ d_n & & & & f_{d_n, t_n} \end{array} \Rightarrow \begin{array}{c|cccc} & d_1 & d_2 & \dots & d_n \\ \hline d_1 & 0 & & & \\ d_2 & & 0 & & \\ \dots & & & \dots & \\ d_n & & & & 0 \end{array}$$

$$C = \{d_{1,C_1}, d_{2,C_2}, \dots, d_{n,C_k}\} \text{ where } k \leq n$$

Convert our matrix of n documents measured on k terms to a matrix of inter-document similarity and then apply a clustering method to the similarity/distance matrix

Either because we want **conceptually meaningful groups** of documents (or terms) that share common characteristics *or* because we want **useful groups** that abstract from the individual documents (summarization or compression)

Clustering for **understanding** or **utility**



Principle Component Analysis of the DTM is often used for visualization purpose

Agglomerative hierarchical clustering

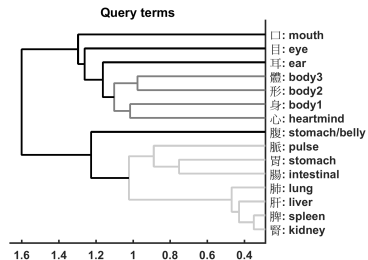
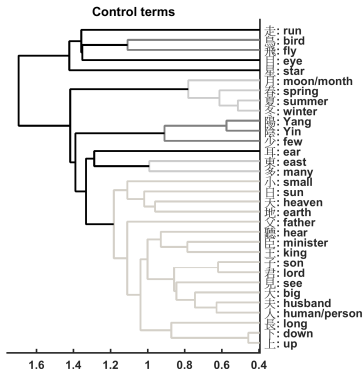
Set of clustering methods that starts with each document as a single cluster and then repeatedly merge the two closest clusters until a single, all encompassing cluster remains (alternate methods use divisive clustering)

Hierarchical clustering produce nested clusters that are organized in a tree-like structure (visualized with a dendrogram)

-
- | | |
|----|---|
| 1. | compute proximity matrix |
| 2. | repeat |
| 3. | merge the closest two clusters |
| 4. | update the proximity matrix to reflect the distance between
the new clusters and the original clusters |
| 5. | until only one cluster remains |
-

To compute the proximity between groups of data points a particular technique is chosen (e.g. MIN, MAX, group average)

Dualism in Ancient Chinese Litt.



With hierarchical clustering you cut or prune the tree at some level to define clusters.