

models and algorithms #2

MGMT|from text analysis to actionable knowledge

kristoffer l nielbo
kln@cas.au.dk

DTL|IMC|AARHUS UNIVERSITY-CAS

the goal of **statistical learning** is to build a machine that can learn from data and automatically make the right decisions

unsupervised

identify latent classes in the data → lack theoretical 'ground truth'

supervised

infer mapping between data & class-information → theoretical 'ground truth'

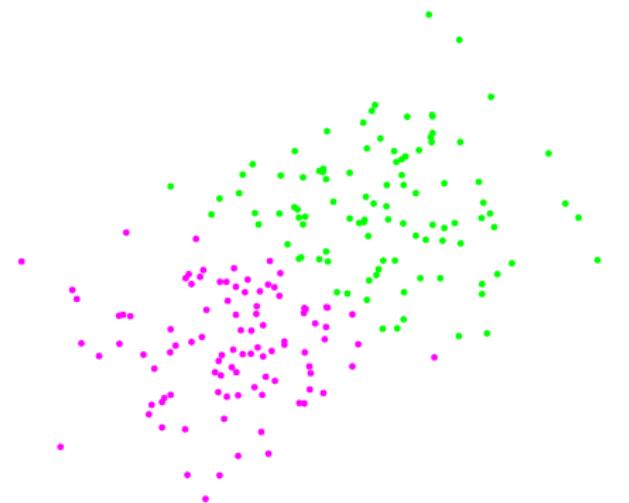
unsupervised learning|clustering

groups of features or latent variables often identify non-random subsets in a collection of documents

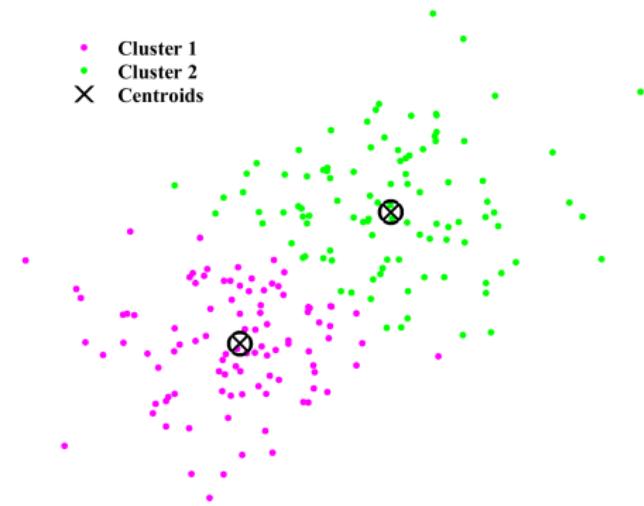
these subsets or ‘clusters’ can be used for understanding document content (e.g., thematic analysis) and computational utility (e.g., identification of prototype documents)



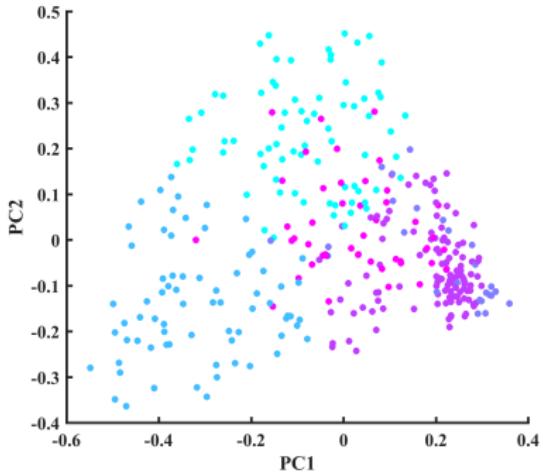
implicit assumption that we study differences in variables (e.g., terms) between homogeneous objects (e.g., documents)



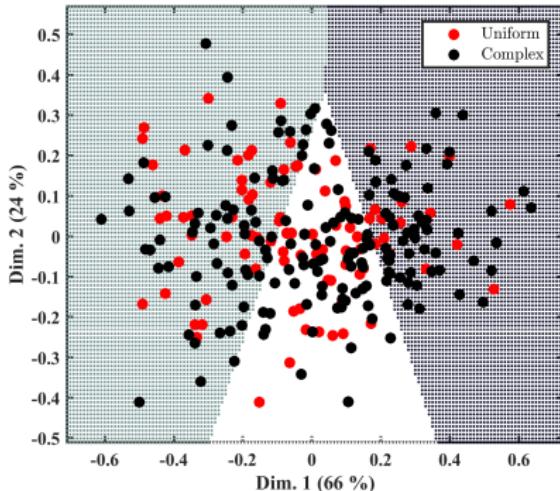
systematic differences between documents result in sub-corpora (e.g., genre, author characteristics, community effects)



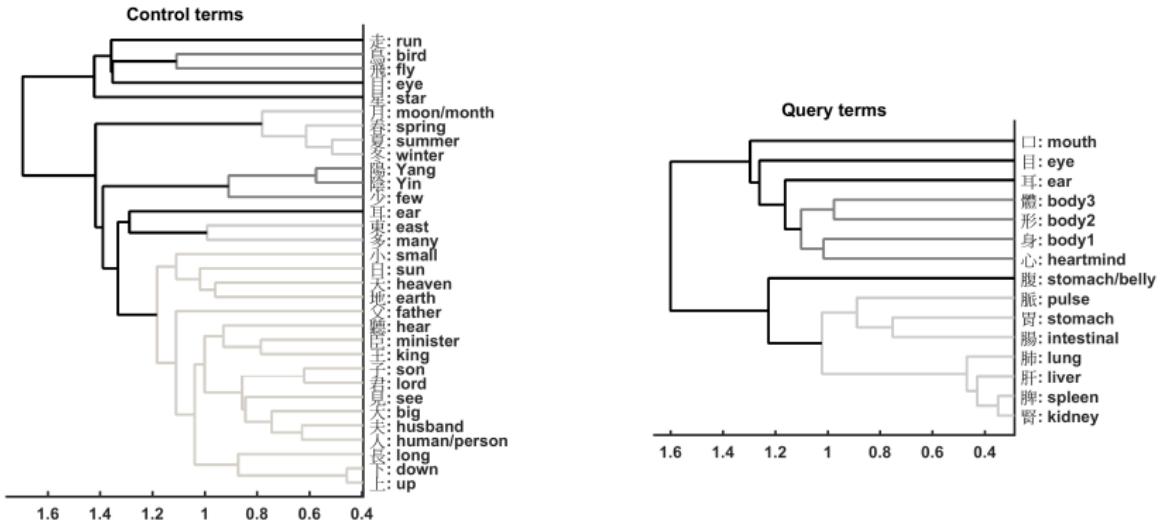
partition documents (or other objects) into homogeneous subsets using document similarity/distance



KEYWORDS	brother fight mujahidin ummah women	abu brother children love people	Al Qaeda AQAP Inspire issue magazine	America attack bomb people world	Allah* Islam message Muslim time

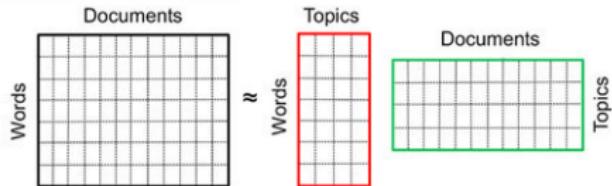


Causal relationship	LRD in <i>ads + articles</i>	LRD in <i>ads</i>	LRD in <i>articles</i>	no LRD
Uniform ($ads \Rightarrow article$)	cigarettes	\emptyset	LIFESTYLE sweater flanel cocker spaniels parakeet bikes	ENTERTAINMENT cinema restaurent
Complex ($ads \Leftrightarrow article$)	PRODUCE potatoes apples beans lettuce cauliflower	\emptyset	INTERIOR couch dressoir carpet	TECHNOLOGY tape-recorder radio television



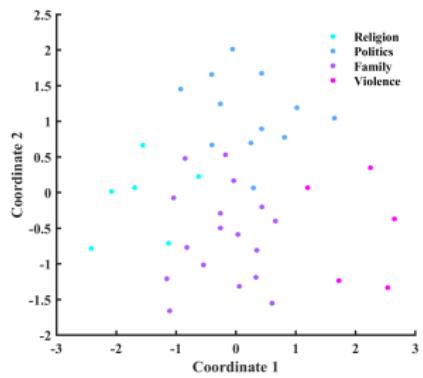
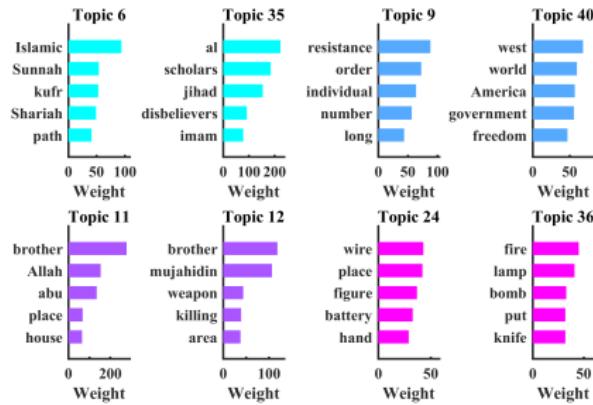
– with hierarchical clustering you cut or prune the tree at some level to define clusters

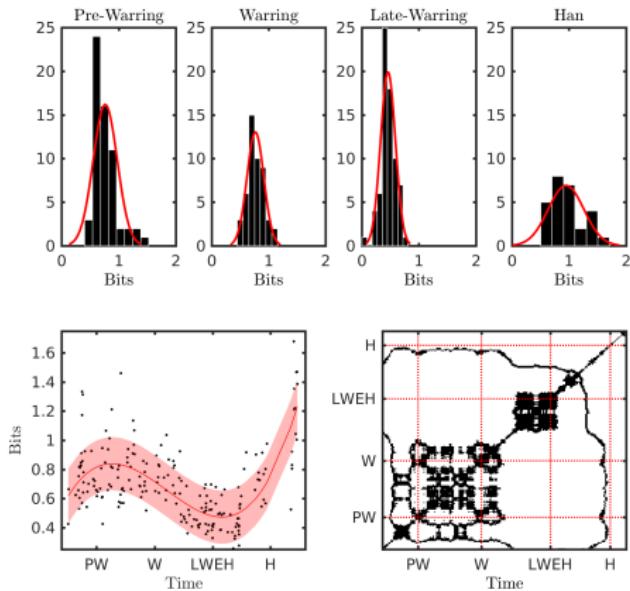
we can extract latent semantic variables that generated the documents by reverse engineering the process *words & doc* \Rightarrow *topics*



topic modeling a set of unsupervised mixed models where each document is more or less likely within each clusters

- ① **discover** thematic structure
- ② **annotate** documents
- ③ **use** the annotations to visualize, organize, summarize, ...





represent each document as a distribution on topics: $\theta_{d=1\dots M}$

semantic innovation \sim relative entropy between documents:

$$D_{KL}(P \parallel Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

supervised learning|classification

documents often come with associated metadata in the form of classes or labels (e.g., genres, departments, companies) → approximate a function that can map document features onto class information

train a classifier on labeled documents, test on unseen documents, and generalize to new instances

while clustering (unsupervised learning) searches for groups within the corpus, classification learns to map a collection of documents onto a categorical class values

