



Machine Learning
– winter term 2016/17 –

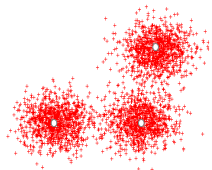
Chapter 05: Clustering

Prof. Adrian Ulges
Masters “Computer Science”
DCSM Department
University of Applied Sciences RheinMain

Unsupervised Learning = Learning without Labels images from [2], [1]



- ▶ **Clustering**: discover coherent groups of samples
- ▶ **Dimensionality reduction**: compressing samples
- ▶ **Itemset mining**: finding frequent substructures in the data
- ▶ **Anomaly detection**: detecting outliers in the data

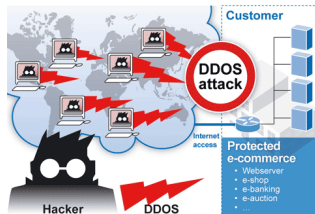


Customers Who Bought This Item Also Bought

LOOK INSIDE!

slide:ology: The Art and Science of Creating Great... by Nancy Duarte
★★★★☆ (98)
\$23.09

The Naked Presenter: Delivering Powerful Present... by Garr Reynolds
\$16.49



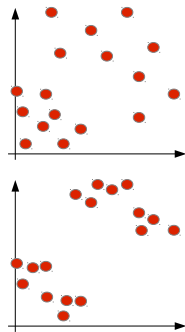


1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
4. Expectation Maximization
5. Document Clustering
6. Agglomerative Clustering

Clustering: Definition



- ▶ Clustering (or *cluster analysis*) is an unsupervised learning problem (*remember: **samples only**, no labels*)
- ▶ The challenge is to discover coherent subgroups (or *clusters*) of samples
- ▶ **Difference to classification:** In clustering, we try to *find* the classes and assign samples to them



Challenges

1. Often, it is unclear by which criterion to cluster (*example: cluster users, but by which demographic attributes?*)
2. Cluster granularity is unclear a priori

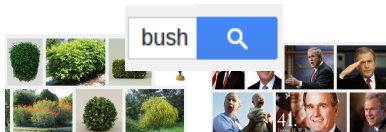
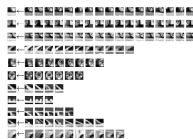
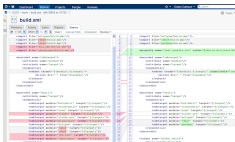
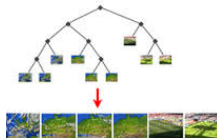


Clustering: Applications images from [4], [3]



Clustering has numerous **applications** in various areas

- ▶ market research
- ▶ life sciences
- ▶ information retrieval
- ▶ computer vision
- ▶ social networks
- ▶ data mining



Example: Demographic Clustering on YouTube [8]



Video X



tags T

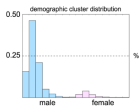


User-Cluster U

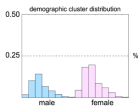


$P(T|X)$

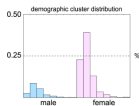
$P(U|T)$



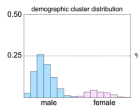
counterstrike,
skateboarding,
worldofwarcraft,
darth-vader,
simpsons, soccer



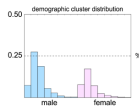
singing, cake,
cooking, choir, food,
baby, kitchen, cats,
dancing, dogs



horse, anime,
cheerleading, kiss,
gymnastics, cake,
riding, dancing,
videoblog



obama, mccain,
georgewbush, court,
interview,
press-conference,
airplane-flying, riot



americas-got-talent,
cats, cartoon,
origami, piano,
muppets,
commercial, tornado





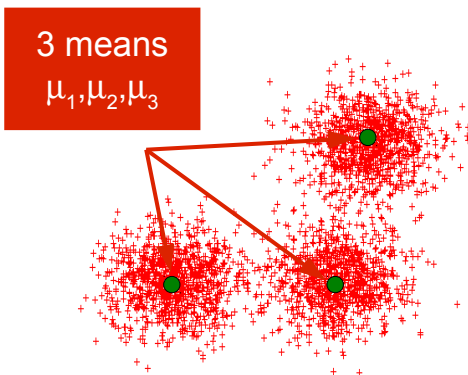
1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
4. Expectation Maximization
5. Document Clustering
6. Agglomerative Clustering

Clustering: K-Means



We start with the “first choice” clustering algorithm: **K-Means**

- ▶ Given: samples $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ We assume that samples are clustered around K centers (the “ K means”) $\mu_1, \dots, \mu_K \in \mathbb{R}^d$
- ▶ Each sample \mathbf{x}_i belongs to a mean $k(i)$
- ▶ The clusters are **spheres** of **identical size**



K-Means: Approach

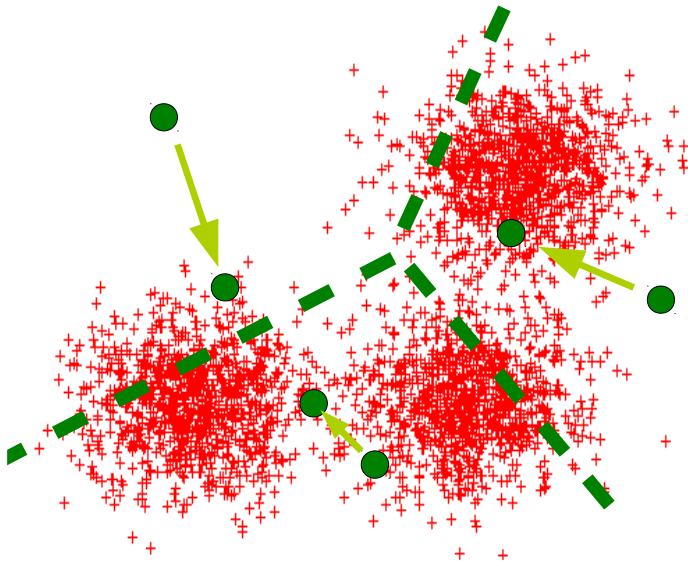


When trying to determine the clusters / the means, we face a **chicken-egg problem**

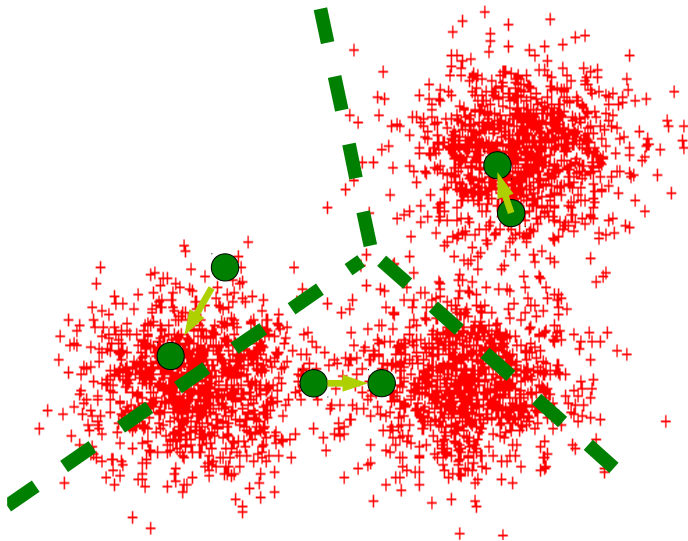
- ▶ If we knew the clusters, we could easily determine the means
(by averaging all samples of a cluster)
- ▶ If we knew the means, we could determine the clusters
(by assigning each sample to its closest mean)
- ▶ Approach (**interleaved optimization**): Alternately, fix the clusters/means and estimate the other

```
1 function KMEANS( $\mathbf{x}_1, \dots, \mathbf{x}_n, K$ )
2   initialize  $\mu_1, \dots, \mu_K$  by random sampling from  $\mathbf{x}_1, \dots, \mathbf{x}_n$ 
3   repeat
4     for  $i = 1, \dots, n$ : // assign each sample to its closest cluster
5        $k(i) := \arg \min_{k=1, \dots, K} \|\mathbf{x}_i - \mu_k\|$ 
6     for  $k = 1, \dots, K$ : // re-estimate each cluster's mean
7        $X_k := \{\mathbf{x}_i \mid k(i) = k\}$ 
8        $\mu_k := \frac{1}{|X_k|} \sum_{\mathbf{x} \in X_k} \mathbf{x}$ 
9   until  $k(1), \dots, k(n)$  do not change
10  return  $\mu_1, \dots, \mu_K$ 
```

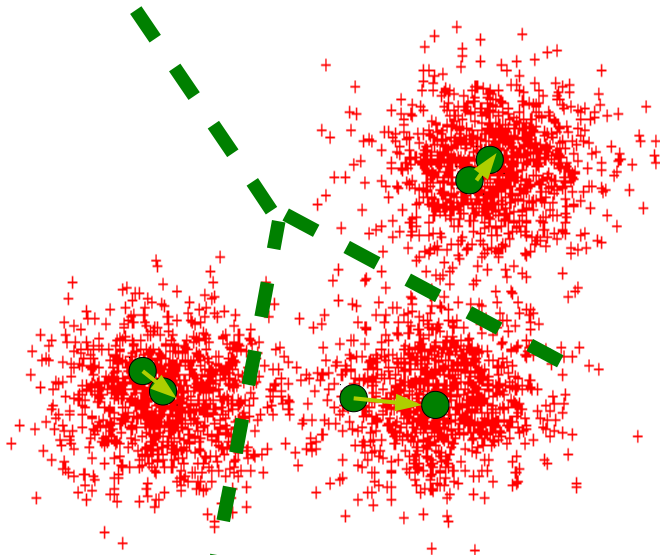
K-Means: Example (Step 1)



K-Means: Example (Step 2)



K-Means: Example (Step 3...)



K-Means: Properties



- ▶ K-Means corresponds to a local optimization of the sum of squared errors

$$E(\mu_1, \dots, \mu_K) = \sum_{i=1}^n (\mathbf{x}_i - \mu_{k(i)})^2$$

- ▶ Computational effort: $O(K \cdot n \cdot d)$ per iteration. The number of iterations is often moderate.
- ▶ Convergence is guaranteed.

Proof of Convergence

K-Means: Properties



Proof of Convergence (cont'd)

K-Means: Properties



Proof of Convergence (cont'd)

K-Means: Properties



Does K-Means always lead to the same results?

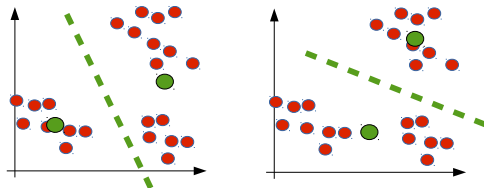
No: K-Means is a local search method!

- ▶ **Problem 1:** The order of means can be permuted

$$\mu_1 = (0, 0), \mu_2 = (1, 1), \mu_3 = (5, 3)$$

$$\mu_1 = (5, 3), \mu_2 = (0, 0), \mu_3 = (1, 1)$$

- ▶ **Problem 2:** The resulting means can be completely different
- ▶ **Approach:** Restart multiple times, and keep the result with minimal error E .
- ▶ During the algorithm, **empty clusters** may occur. **Approach:** Reinitialize the corresponding center randomly and continue.

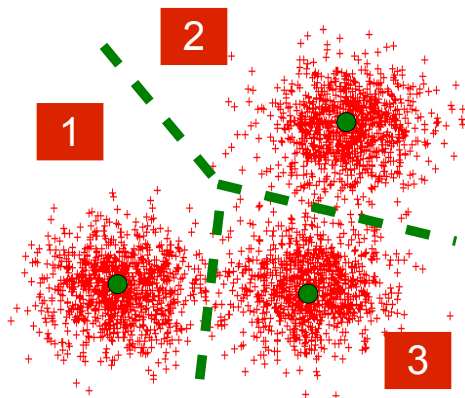


K-Means: Properties (cont'd)

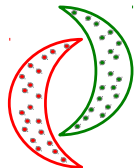
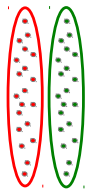
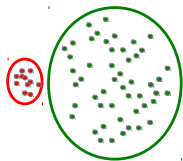


Given a clustering result μ_1, \dots, μ_K , we can assign new samples \mathbf{x} to clusters (this is called **vector quantization**):

$$k(\mathbf{x}) = \arg \min_k \|\mathbf{x} - \mu_k\|$$



K-Means: Discussion





1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
4. Expectation Maximization
5. Document Clustering
6. Agglomerative Clustering

Choosing K : Model Selection



“Model selection is the task of selecting a statistical model from a set of candidate models, given data.”

(en.wikipedia.org)

Here: Model Selection = Choosing K

- ▶ K too small (*undersegmentation*): clusters too diverse
- ▶ K too high (*oversegmentation*): too many parameters, clusters too fine-grain
- ▶ Choosing the 'wrong' K leads to **instable results**

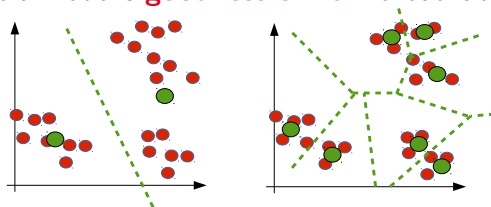
Approach 1: External Benchmark

- ▶ Sometimes, clustering is just one processing step of a **larger system**, and we can benchmark that larger system
- ▶ **Example**: User clustering for advertising
(\rightarrow benchmark by click-through-rate)

Approach 2: Cluster Validation



Goal: measure a model's **goodness-of-fit** without labels



Example: The **Bayes' Information Criterion (BIC)**

1. The clusters should be **compact** (*small error E*)
 2. The model should be simple, i.e. have only **few parameters**
- ▶ Let θ be the model parameters to learn, and let $\#\theta$ be their number (e.g., in K-Means: $\#\theta = K \cdot d$)
 - ▶ Test different values of K , and pick this one:

$$K^* = \arg \min_K -2 \ln \left(p(\mathbf{x}_1, \dots, \mathbf{x}_n | \theta) \right) + \#\theta \cdot \ln(n)$$

BIC for K-Means: Derivation



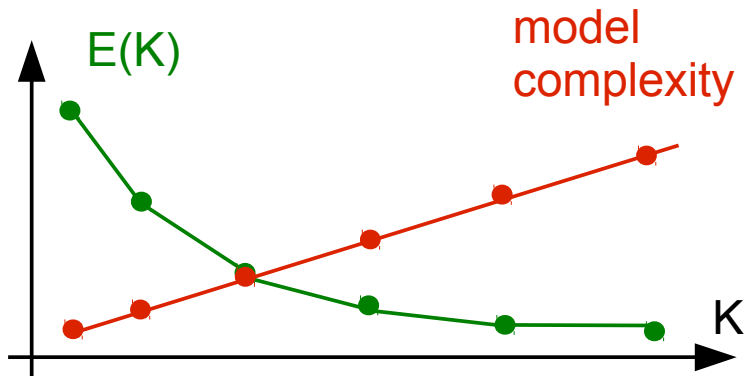
BIC for K-Means: Derivation



The Bayes Information Criterion



$$K^* = \arg \min_K \underbrace{\sum_{i=1}^n (x_i - \mu_{k(i)})^2}_{E(K)} + \underbrace{K \cdot d \cdot \ln(n)}_{\text{model complexity}}$$



Selecting K : Search Strategies

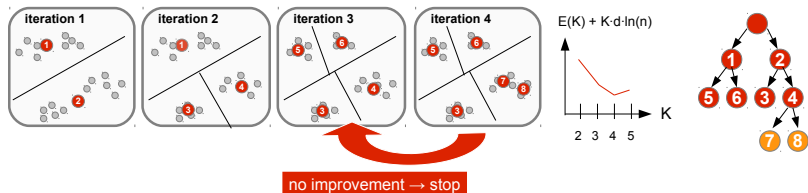


Approach 1: **Naive**

- ▶ test values for K in a reasonable range.
- ▶ For every K , re-run clustering and evaluate (**expensive!**)

Approach 2: **Hierarchical Clustering** (*more efficient*)

- ▶ ... Iteratively, pick the largest cluster
- ▶ ... and apply K -Means to the samples in this cluster, obtaining K new clusters
- ▶ ... stop once the overall quality (e.g., BIC) stops improving
- ▶ We obtain a **tree** of clusters



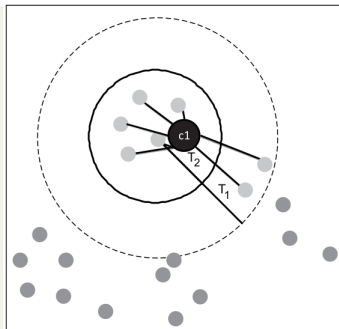
Selecting K : Canopy Clustering image from [7]



Approach 3: Canopy Clustering

- ▶ A **greedy strategy** to find (potentially suboptimal) clusters on large datasets
- ▶ We use it to estimate K and to initialize the means
- ▶ Canopy clusters can **overlap!**
- ▶ Canopy clustering uses **two thresholds**
 - ▶ T_1 (determines the number of clusters)
 - ▶ T_2 (determines the overlap of clusters) ($T_2 > T_1$)

```
1 function CLUSTER_CANOPY( $X := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ )
2    $C := \{\}$ 
3   while  $X \neq \{\}$ :
4     choose a random sample  $\mathbf{x} \in X$ 
5      $Y := \{\mathbf{y} \in X \mid \|\mathbf{y} - \mathbf{x}\| \leq T_1\}$ 
6      $Z := \{\mathbf{y} \in X \mid T_1 < \|\mathbf{y} - \mathbf{x}\| \leq T_2\}$ 
7      $C := C \cup \{\mathbf{x}\}$ 
8      $X := X \setminus Y$ 
9   return  $C$ 
```





1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
- 4. Expectation Maximization**
5. Document Clustering
6. Agglomerative Clustering

Expectation Maximization (EM)



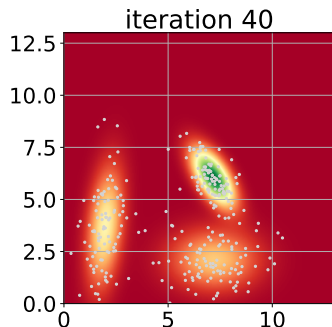
- ▶ We can overcome some of the above limitations by **generalizing K-Means**, resulting in a famous approach called **Expectation Maximization (EM)**

EM: Model

- ▶ We explain the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ by a **Gaussian mixture model**

$$\mathbf{x}_1, \dots, \mathbf{x}_n \sim \sum_{k=1}^K P_k \cdot p(\mathbf{x} | \mu_k, \Sigma_k)$$

where p is the **multivariate normal density** (*Chapter 3*), μ_1, \dots, μ_K are K centers, $\Sigma_1, \dots, \Sigma_K$ are K covariance matrices (the *shapes* of the clusters), and P_1, \dots, P_K are the cluster's proportions of the data (also called *priors*).



Expectation Maximization (EM)



Remarks

- ▶ In K-Means, we would have $P_1 = P_2 = \dots = P_K = \frac{1}{K}$ and

$$\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ 0 & \dots & \dots & 0 \\ 0 & \dots & 0 & \sigma^2 \end{pmatrix}$$

Approach

- ▶ We **rename** the two alternating K-Means steps

E-Step Re-assigning samples to clusters \rightarrow “Expectation-Step”

M-Step Re-estimating the cluster centers \rightarrow “Maximization-Step”

- ▶ We **modify** these steps a bit
 - ▶ **E-Step**: No hard assignment of samples to centers, but a **soft assignment** by computing the probability $P(k(i) = k | \mathbf{x}_i)$
 - ▶ **M-Step**: Do not only estimate the cluster *centers*, but **parameters** in general (e.g., the clusters' shape+prior)

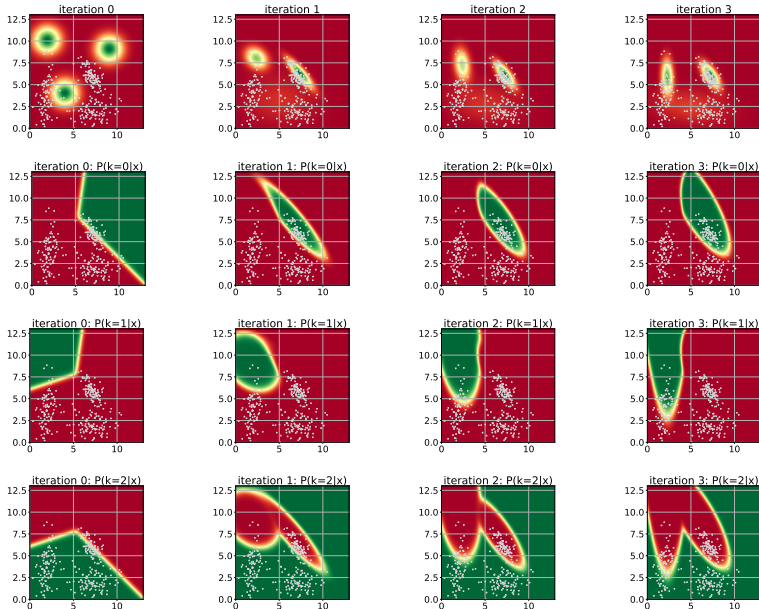
K-Means vs. Expectation Maximization (EM)



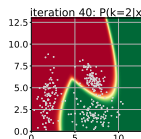
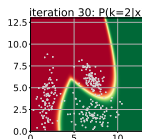
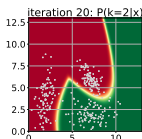
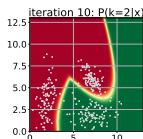
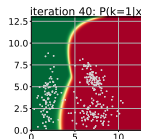
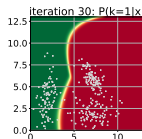
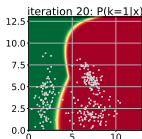
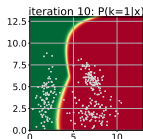
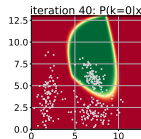
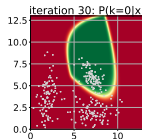
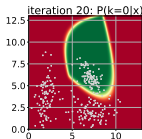
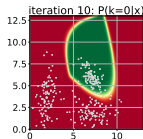
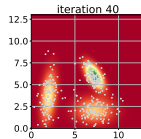
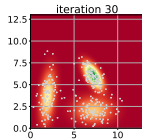
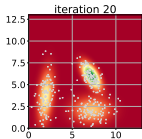
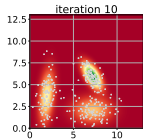
Illustration

	K-Means	EM
E-Step	$k(i) := \arg \min_k \ \mathbf{x}_i - \mu_{k(i)}\ $	$w_{ki} := P(k(i) = k \mathbf{x}_i) = \frac{p(\mathbf{x}_i; \mu_k, \Sigma_k)}{\sum_{k'} p(\mathbf{x}_i; \mu_{k'}, \Sigma_{k'})}$
M-Step	$\mu_k := \frac{\sum_{\mathbf{x} \in X_k} \mathbf{x}}{ X_k }$	$\mu_k := \frac{\sum_i w_{ki} \cdot \mathbf{x}_i}{\sum_i w_{ki}}$
	—	$\Sigma_k := \frac{\sum_i w_{ki} \cdot (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{\sum_i w_{ki}}$
	—	$P_k := \frac{\sum_i w_{ki}}{\sum_{k'} \sum_i w_{k'i}}$

EM: Example



EM: Example

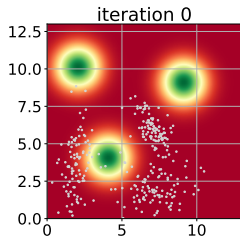


EM: Goodness-of-Fit

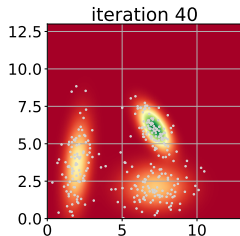


- ▶ Goal: restart EM many times, pick the 'best' model.
- ▶ Given an **EM model** $\Theta = (\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, P_1, \dots, P_K)$, we want to measure its “**goodness-of-fit**”.
- ▶ Approach: We measure the **likelihood** of the data

$$\begin{aligned}L(\mathbf{x}_1, \dots, \mathbf{x}_n; \Theta) &= \prod_i p(\mathbf{x}_i | \Theta) \\ &= \prod_i \sum_k P_k \cdot p(\mathbf{x}_i; \mu_k, \Sigma_k)\end{aligned}$$



low likelihood



high likelihood



EM as a general Learning Scheme



- ▶ EM for **Gaussian Mixture Models** is just a special case!

symbol	general EM	Gaussian Mixture Models
X	(known) input data	the features $\mathbf{x}_1, \dots, \mathbf{x}_n$
Θ	parameters	means μ_1, \dots, μ_K , shapes $\Sigma_1, \dots, \Sigma_K$, priors P_1, \dots, P_K
U	unknown data	the mapping from \mathbf{x}_i to clusters k

EM: General Learning Scheme

```
1 function EM( $X$ )
2   initialize  $\Theta$  randomly
3   repeat
4     compute  $P(U|X, \Theta)$  // E-step
5     optimize parameters [6], obtaining a new  $\Theta$  // M-step
6   until convergence
7   return  $\Theta$ 
8
```



1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
4. Expectation Maximization
- 5. Document Clustering**
6. Agglomerative Clustering

Document Clustering



We can also **cluster text** using EM. The resulting method is called **Probabilistic Latent Semantic Analysis (PLSA)**

- ▶ Given: a set of documents with their **bag-of-words** features
- ▶ PLSA divides the set of documents into clusters of **semantically similar** documents
- ▶ The cluster centers correspond to prototypical word distributions (or **topics**)
- ▶ We also call PLSA a **topic model**

Remarks

- ▶ This is clustering, not classification! Categories/topics are not pre-defined, but PLSA discovers them by itself!

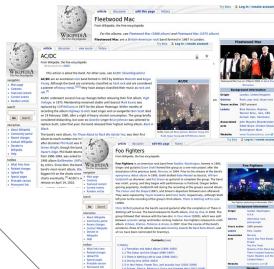
PLSA: Illustration



index
information
search
retrieval

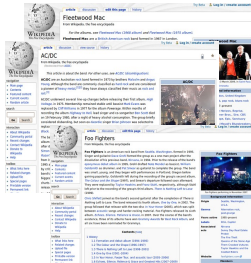
team
soccer
Bundesliga
champion

rock
concert
band
album



PLSA: Notation

- ▶ **Input:** a collection of documents d_1, \dots, d_n and a vocabulary of terms w_1, \dots, w_m
- ▶ Each document d is represented by its **bag-of-words** feature. This gives us a **probability distribution** of words $P(w|d)$.
- ▶ We assume the document collection to consist of K **topics** z_1, \dots, z_K
- ▶ Each topic z has a **word distribution** $P(w|z)$ (just like a document)
- ▶ A document d can be seen a **mixture of topics**, $P(z|d)$



rock concert
band album

PLSA: Sampling Process



Words are **sampled** from a document d in **two steps**

1. Choose a random topic z' from $P(z_1|d), \dots, P(z_K|d)$
2. Given z' , pick a word from $P(w_1|z'), \dots, P(w_m|z')$

Document d



1. choose topic $P(z|d)$

$$P(z_1|d) = 0.01$$

$$P(z_2|d) = 0.6$$

$$P(z_3|d) = 0.39$$

Topics z_1, z_2, z_3



2. choose term $P(w|z)$

soccer
team
text music
information

soccer band
cup player
retrieval
team

soccer
sound band
music live
information

PLSA: Derivation



PLSA Clustering estimates two **probability distributions**:

1. $P(z|d)$

- ▶ $P(z|d)$ tells us **which topics appear in a document** (or which topics (clusters) a document *belongs to*)
- ▶ $P(z|d)$ is a $K \times n$ probability table (covering *all topic-document combinations*)

2. $P(w|z)$

- ▶ This distributions tells us **which words appear in a topic**
- ▶ $P(w|z)$ is an $m \times K$ probability table (covering *all word-topic combinations*)

PLSA Approach

To estimate the above distributions, PLSA uses the EM Algorithm

- ▶ **E-Step** (*assign samples to clusters*)
→ assign **words to topics** (= compute $P(z|w, d)$)
- ▶ **M-Step** (*estimate cluster parameters*)
→ estimate **topics and mixtures** (= $P(w|z), P(z|d)$)

PLSA: Algorithm

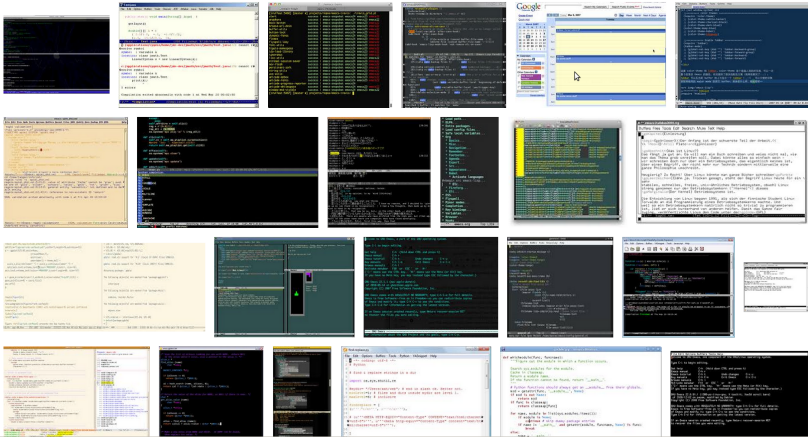


- ▶ **Given:** documents d_1, \dots, d_n , terms w_1, \dots, w_m
- ▶ **Given:** bag-of-words $P(w|d)$, number of topics K
- ▶ Initialize $P(w|z)$, $P(z|d)$ randomly
- ▶ Repeat until convergence:





PLSA: Code Sample







1. Clustering: Basics
2. K-Means
3. Model Selection: Selecting K
4. Expectation Maximization
5. Document Clustering
6. Agglomerative Clustering

Agglomerative Clustering

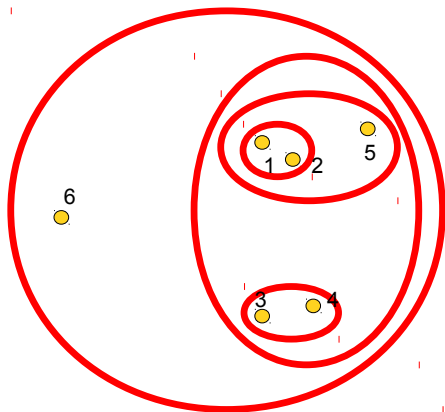


- ▶ We call K-Means/EM **divisive** clustering techniques, because they divide the dataset top-down

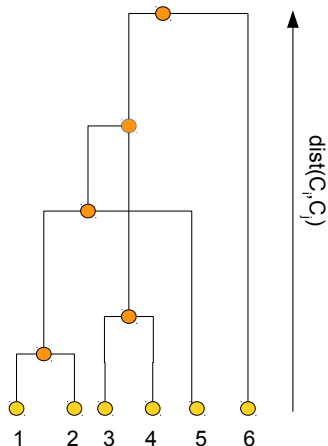
Agglomerative clustering

- ▶ initially, each sample belongs to its own cluster (*singleton*)
- ▶ iteratively, we merge the two “most similar” clusters and obtain a new, bigger cluster
- ▶ The result can be illustrated in form of a tree-like graph, a so-called **dendrogram**

Agglomerative Clustering: Illustration



Dendrogramm



Agglomerative Clustering: Further Issues



Stopping Criterion

- ▶ heuristics (see *model selection*)

Distance Measure

- ▶ We need to define a **distance between clusters** to pick the “most similar” clusters to fuse.
- ▶ The three common alternatives (let X, Y be clusters):

single linkage	$dist(X, Y) := \min_{x \in X, y \in Y} \ x - y\ ^2$
complete linkage	$dist(X, Y) := \max_{x \in X, y \in Y} \ x - y\ ^2$
average linkage	$dist(X, Y) := \frac{1}{ X \cdot Y } \sum_{x \in X, y \in Y} \ x - y\ ^2$

Agglomerative Clustering: Discussion




Agglomerative Clustering: Application Example [5]



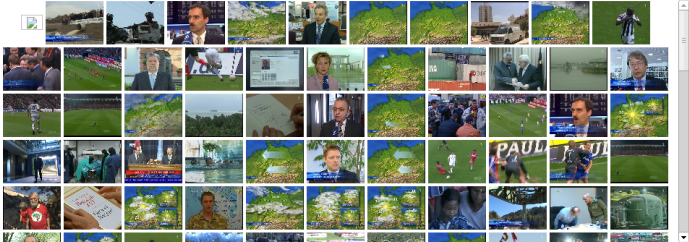
navidator^{.beta}

Start Browsing ::

Click History:



Database: Level 7 | 222 Items | Resample



References I



- [1] M. Dukia: DoS(Denial of Service) Attacks.
<http://www.securitykiller.org/2015/12/dosdenial-of-service-attacks.html> (retrieved: Nov 2016).
- [2] P. Ipeirotis: A Plea to Amazon: Fix Mechanical Turk! .
<http://www.behind-the-enemy-lines.com/2010/10/plea-to-amazon-fix-mechanical-turk.html>
(retrieved: Nov 2016).
- [3] The Value of a Professional Network?
<https://www.linkedin.com/pulse/value-professional-network-daniel-tunkelang> (retrieved: Oct 2016).
- [4] University of Vermont: Complex Networks (Course Page).
<http://www.uvm.edu/pdodds/teaching/courses/2010-01UVM-303/content/pictures.html> (retrieved: Nov 2016).
- [5] D. Borth, C. Schulze, A. Ulges, and T. Breuel.
Navidgator - Similarity Based Browsing for Image & Video Databases.
In Proc. KI 2008, pages 22–29, September 2008.
- [6] F. Dellaert.
The Expectation Maximization Algorithm.
Technical Report GIT-GVU-02-20, Georgia Institute of Technology, 2002.
- [7] S. Owen, R. Anil, T. Dunning, and E. Friedman.
Mahout in Action.
Manning Publications, 2012.
- [8] A. Ulges, M. Koch, and D. Borth.
Linking Visual Concept Detection with Viewer Demographics.
In ACM International Conference on Multimedia Retrieval, 6 2012.