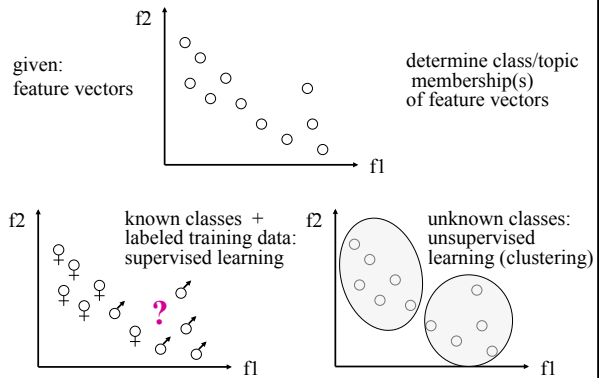


Classification & Clustering

Sergej Sizov
Information Retrieval
Summer term 2013



Classification Problem (Categorization)



Assessment of Classification Quality

empirical by automatic classification of documents that do not belong to the training data (but in benchmarks class labels of test data are usually known)

For **binary classification** with regard to class C:

- a = #docs that are classified into C and do belong to C
- b = #docs that are classified into C but do not belong to C
- c = #docs that are not classified into C but do belong to C
- d = #docs that are not classified into C and do not belong to C

$$\text{Accuracy (Genauigkeit)} = \frac{a + d}{a + b + c + d} \quad \text{Error (Fehler)} = 1 - \text{accuracy}$$

$$\text{Precision (Präzision)} = \frac{a}{a + b} \quad \text{Recall (Ausbeute)} = \frac{a}{a + c}$$

$$\text{F1 (harmonic mean of precision and recall)} = \left(\frac{1}{\text{precision}} + \frac{1}{\text{recall}} \right)^{-1}$$

For **manyway classification** with regard to classes C_1, \dots, C_k :

- macro average over k classes or
- micro average over k classes

Estimation of Classifier Quality

use benchmark collection of completely labeled documents (e.g., Reuters newswire data from TREC benchmark)

cross-validation (with held-out training data):

- partition training data into k equally sized (randomized) parts,
- for every possible choice of k-1 partitions
 - train with k-1 partitions and apply classifier to kth partition
 - determine precision, recall, etc.
- compute micro-averaged quality measures

leave-one-out validation/estimation:

variant of cross-validation with two partitions of unequal size: use n-1 documents for training and classify the nth document

Distance-based Classifiers: k-Nearest-Neighbor Method (kNN)

Step 1:
find among the training documents of all classes the k (e.g. 10-100) most similar documents (e.g., based on cosine similarity):
the k nearest neighbors of \vec{d}

Step 2:
Assign \vec{d} to class C_j for which the function value

$$f(\vec{d}, C_j) = \sum_{\vec{v} \in kNN(\vec{d})} sim(\vec{d}, \vec{v}) * \begin{cases} 1 & \text{if } \vec{v} \in C_j \\ 0 & \text{otherwise} \end{cases}$$

is maximized

With binary classification assign \vec{d} to class C if $f(\vec{d}, C)$ is above some threshold δ ($\delta > 0.5$)

Automatic clustering

Clustering: Classification based on Unsupervised Learning

given:

n **m-dimensional data records** $d_j \in D \subseteq \text{dom}(A_1) \times \dots \times \text{dom}(A_m)$
with attributes A_i (e.g. term frequency vectors $\subseteq N_0 \times \dots \times N_0$)
or n **data points** with pair-wise **distances (similarities)** in a **metric space**

wanted:

k **clusters** c_1, \dots, c_k and an assignment $D \rightarrow \{c_1, \dots, c_k\}$ such that the

average **intra-cluster similarity** $\frac{1}{k} \sum_k \left(\frac{1}{|c_k|} \sum_{d \in c_k} \text{sim}(\vec{d}, \vec{c}_k) \right)$
is high and

the average **inter-cluster similarity** $\frac{1}{k(k-1)} \sum_{\substack{i,j \\ i \neq j}} \text{sim}(\vec{c}_i, \vec{c}_j)$
is low,

where the **centroid** \vec{c}_k of c_k is: $\vec{c}_k = \frac{1}{|c_k|} \sum_{d \in c_k} \vec{d}$

Hierarchical vs. Flat Clustering

Hierarchical Clustering:

- detailed and insightful
- hierarchy built in natural manner from fairly simple algorithms
- relatively expensive
- no prevalent algorithm

Flat Clustering:

- data overview & coarse analysis
- level of detail depends on the choice of the number of clusters
- relatively efficient
- K-Means and EM are simple standard algorithms

Hierarchical Clustering: Agglomerative Bottom-up Clustering (HAC)

Principle:

- start with each d_i forming its own singleton cluster c_i
 - in each iteration combine the most similar clusters c_i, c_j into a new, single cluster
- for $i:=1$ to n do $c_i := \{d_i\}$ od;
 $C := \{c_1, \dots, c_n\}$; /* set of clusters */
while $|C| > 1$ do
 determine $c_i, c_j \in C$ with maximal inter-cluster similarity;
 $C := C - \{c_i, c_j\} \cup \{c_i \cup c_j\}$;
od;

Alternative Similarity Metrics for Clusters

given: similarity on data records - $\text{sim}: D \times D \rightarrow \mathbb{R}$ oder $[0,1]$
define: similarity between clusters - $\text{sim}: 2^D \times 2^D \rightarrow \mathbb{R}$ or $[0,1]$

Alternatives:

- **Centroid method:** $\text{sim}(c, c') = \text{sim}(d, d')$ with centroid d of c and centroid d' of c'
- **Single-Link method:** $\text{sim}(c, c') = \text{sim}(d, d')$ with $d \in c, d' \in c'$, such that d and d' have the highest similarity
- **Complete-Link method:** $\text{sim}(c, c') = \text{sim}(d, d')$ with $d \in c, d' \in c'$, such that d and d' have the lowest similarity
- **Group-Average method:** $\frac{1}{|c| \cdot |c'|} \sum_{d \in c, d' \in c'} \text{sim}(d, d')$

For hierarchical clustering the following axiom must hold:
 $\max \{ \text{sim}(c, c'), \text{sim}(c, c'') \} \geq \text{sim}(c, c' \cup c'')$ for all $c, c', c'' \in 2^D$

Cluster Quality Measures (1)

With regard to **ground truth**:

known class labels L_1, \dots, L_g for data points d_1, \dots, d_n :
 $L(d_i) = L_j \in \{L_1, \dots, L_g\}$

With cluster assignment $\Gamma(d_1), \dots, \Gamma(d_n) \in c_1, \dots, c_k$
cluster c_j has **purity** $\max_{v=1..g} |\{d \in c_j \mid L(d) = L_v\}| / |c_j|$

Complete clustering has purity $\sum_{j=1..k} \text{purity}(c_j) / k$

Alternatives:

- **Entropy** within cluster $\sum_{v=1..g} \frac{|c_j \cap L_v|}{|c_j|} \log_2 \frac{|c_j|}{|c_j \cap L_v|}$

• **MI** between cluster and classes

$$\sum_{c \in \{c_j, \bar{c}_j\}, L \in \{L_1, \dots, L_g\}} \frac{|c \cap L| / n}{|c| \cdot |L| / n} \log_2 \frac{|c| \cdot |L| / n}{|c \cap L| / n}$$

Cluster Quality Measures (2)

Without any ground truth:

ratio of intra-cluster to inter-cluster similarities

$$\frac{1}{k} \sum_k \left(\frac{1}{|c_k|} \sum_{d \in c_k} \text{sim}(\vec{d}, \vec{c}_k) \right) / \left(\frac{1}{k(k-1)} \sum_{\substack{i,j \\ i \neq j}} \text{sim}(\vec{c}_i, \vec{c}_j) \right)$$

or other **cluster validity measures** of this kind
(e.g. considering variance of intra- and inter-cluster distances)

Flat Clustering: Simple Single-Pass Method

given: data records d_1, \dots, d_n
wanted: (up to) k clusters $C := \{c_1, \dots, c_k\}$

```

C := {{d1}}; /* random choice for the first cluster */
for i:=2 to n do
  determine cluster  $c_j \in C$  with the largest value of
   $\text{sim}(d_i, c_j)$  (e.g.  $\text{sim}(d_i, \bar{c}_j)$  with centroid  $\bar{c}_j$ );
  if  $\text{sim}(d_i, c_j) \geq \text{threshold}$ 
    then assign  $d_i$  to cluster  $c_j$ 
  else if  $|C| < k$ 
    then  $C := C \cup \{d_i\}$ ; /* create new cluster */
    else assign  $d_i$  to cluster  $c_j$ 
  fi
fi
od
  
```

K-Means Method for Flat Clustering (1)

Idea:

- determine k **prototype vectors**, one for each cluster
- **assign each data record to the most similar prototype vector** and compute new prototype vector (e.g. by averaging over the vectors assigned to a prototype)
- **iterate** until clusters are sufficiently stable

```

randomly choose  $k$  prototype vectors  $\bar{c}_1, \dots, \bar{c}_k$ 
while not yet sufficiently stable do
  for i:=1 to n do
    assign  $d_i$  to cluster  $c_j$  for which  $\text{sim}(d_i, \bar{c}_j)$  is minimal
  od;
  for j:=1 to k do  $\bar{c}_j := \frac{1}{|c_j|} \sum_{d \in c_j} \bar{d}$  od;
od;
  
```

Example for K-Means Clustering

