

a category cannot be determined with certainty





Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 6





Other applications



Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008

Machine Learning (ML) & TC

- In the '80s, the typical approach used for the construction of TC system involved hand-crafting an expert system consisting of a set of rules, one per category, of the form
 - If <DNF formula> then <category> else ¬<category>
 - Where <DNF formula> is a disjunction of conjunctive clauses
- The drawback of this "manual" approach is the knowledge acquisition bottleneck: since rules must be manually defined, building a classifier is expensive, and if the set of categories is updated or the classifier is ported to a different domain, other manual work has to be done.

11



If	((wheat & farm)	or	
	(wheat & commodity)	or	
	(bushels & export)	or	
	(wheat & tonnes)	or	
	(wheat & winter & :soft))	then	WHEAT else : WHEAT

Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008

Induction

- Since the early '90s, the machine learning approach to the construction of TC systems has become dominant.
- A general inductive process automatically builds a classifier for a category c by 'observing' the characteristics of a set of documents previously classified belonging (or not) to c_i, by a domain expert.
- This is an instance of Supervised Learning (SL).

13

A	dvantages of the SL approach
	The engineering effort goes towards the construction, not of a classifier, but of an automatic builder of classifiers (learner)
	If the set of categories is updated, or if the system is ported to a different domain, all that is need is a different set of manually classified documents
	Domain expertise (for labeling), and not knowledge engineering expertise, is needed. Easier to characterize a concept extensionally than intentionally.
	Sometimes the preclassified documents are already available
	The effectiveness achievable nowadays by these classifiers rivals that of hand-crafted classifiers and that of human classifiers
Dip. Pura	di Matematica F. Aiolli - Sistemi Informativi 15 ed Applicata 2007/2008
Т	raining Set and Test Set
	raining Set and Test Set The ML approach relies on the application of a train-and-test approach to a labeled corpus Tr={d1,,d15}, i.e. a set of documents previously classified under C={c1,,cm}.
T	raining Set and Test Set The ML approach relies on the application of a train-and-test approach to a labeled corpus $Tr=\{d_1,,d_{ S }\}$, i.e. a set of documents previously classified under $C=\{c_1,,c_m\}$. The value of the function $\Phi: D \times C \rightarrow \{-1,+1\}$ are known for every pair $\langle d_j, c_i \rangle$. Tr then constitutes a 'glimpse' of the ground truth.
	raining Set and Test Set The ML approach relies on the application of a train-and-test approach to a labeled corpus $\text{Tr}=\{d_1,,d_{ S }\}$, i.e. a set of documents previously classified under $C=\{c_1,,c_m\}$. The value of the function $\Phi: D \times C \rightarrow \{-1,+1\}$ are known for every pair $\langle d_j, c_i \rangle$. Tr then constitutes a 'glimpse' of the ground truth. We assume that pair $\langle d_j, c_i \rangle$ are extracted according to a probability distribution $P(d_j, c_i)$

For evaluation purposes, a new set Te is usually provide which has elements extracted from the same pair distribution P(d_j,c_i) used for elements in Tr.

/ •	
	Most of the time, the learner is parametric. These parameters should be optimized by testing which values of the parameters yield the best effectiveness.
	Hold-out procedure A small subset of Tr, called the validation set (or hold-out set), denoted Va, is identified
	 A classifier is learnt using examples in Tr-Va. Step 2 is performed with different values of the parameters, and tested against the hold-out sample
	In an operational setting, after parameter optimization, one typically re- trains the classifier on the entire training corpus, in order to boost effectiveness (debatable step!)
	It is possible to show that the evaluation performed in Step 2 gives an unbiased estimate of the error performed by a classifier learnt with the same parameters and with training set of cardinality Tr - Va < Tr
Dip	. di Matematica F. Aiolli - Sistemi Informativi 17
K	-fold Cross Validation
K	-fold Cross Validation An alternative approach to model selection (and avaluation) is the K fold space validation mathed
K	-fold Cross Validation An alternative approach to model selection (and evaluation) is the K-fold cross-validation method K-fold CV procedure
K	 -fold Cross Validation An alternative approach to model selection (and evaluation) is the K-fold cross-validation method K-fold CV procedure K different classifiers h₁,h₂,,h_k are built by partitioning the initial corpus Tr into k disjoint sets Va₁,,Va_k and then iteratively applying the Hold-out approach on the k-pairs <tr, =="" li="" tr-va,="" va.<=""> </tr,>
	 -fold Cross Validation An alternative approach to model selection (and evaluation) is the K-fold cross-validation method K-fold CV procedure K different classifiers h₁,h₂,,h_k are built by partitioning the initial corpus Tr into k disjoint sets Va₁,,Va_k and then iteratively applying the Hold-out approach on the k-pairs <tr<sub>i = Tr-Va_i, Va_i></tr<sub> Effectiveness is obtained by individually computing the effectiveness of h₁,,h_k, and then averaging the individual results

