

Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2006/2007 103

Combination rules



Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2006/2007 104

Boosting

Boosting is a CC method whereby the classifiers ('weak hypothesis') are trained
sequentially by the same learner (weak
learner'), and are combined into a CC ('final
hypothesys')

The training of h_t is done in such a way to try to make the classifier to perform well on examples in which h₁,..,h_{t-1} have performed worst

AdaBoost is a popular Boosting algorithm

Freund & Schapire's AdaBoost

At iteration s:

- 1. Passes a distribution D_s of weights to the weak learner, where $D_s(d_j)$ measures how effective $h_{1,...,h_{s-1}}$ have been in classifying d_i
- 2. The weak learner returns a new weak hypothesis h_s that concentrates on documents with the highest D_s values
- 3. Runs h_s on Tr and uses the results to produce an updated distribution D_{s+1} where
 - Correctly classified documents have their weights decreased
 - □ Misclassified documents have their weights increased

Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2006/2007 106

Evaluating TC systems

- Similarly to IR systems, the evaluation of TC systems is to be conducted experimentally, rather than analytically
- Several criteria of quality:
 - Training-Time efficiency
 - Classification-Time efficiency
 - Effectiveness
- In operational situations, all three criteria must be considered, and the right tradeoff between them depends on the application



Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2006/2007



Summarizing

- All the problem setting above can be seen as homogeneous linear problems in an opportune augmented space
- Any algorithm for linear optimization (e.g. perceptron, SVMs, or a linear programming package) can be used to solve them