

Bagging and Boosting: Resampling for Classifier Design

Sargur Srihari

srihari@cedar.buffalo.edu

Bagging

- *Arcing*– adaptive re-weighting and combining– refers to reusing or selecting data to improve classification
- Includes both *bagging* and *boosting*
- Simplest method is *bagging*
- Name derived from *bootstrap aggregation*
- Uses multiple versions of training set

Bagging Method

- Drawing $n' < n$ samples from D with replacement
- Each bootstrap set is used to train a different component classifier
- Final classification decision based on vote of each component classifier
- Component classifiers are all of the same form
 - all HMMs, all neural networks, all decision trees
 - Different parameter values due to different training sets

Bagging Overcomes Classifier Instability

- Unstable if small changes in training data lead to significantly different classifiers or large changes in accuracy
- Decision Tree algorithms can be unstable
 - Slight change in the position of a training point can lead to a radically different tree
- Bagging improves recognition for unstable classifiers by smoothing over discontinuities

Boosting

- Powerful technique for combining multiple “base” classifiers to form a committee whose performance can be significantly better than any of the base classifiers
- *AdaBoost* (adaptive boosting) can give good results even if base classifier performance is only slightly better than random
- Hence base classifiers are known as *weak learners*
- Designed for classification, can be used for regression as well

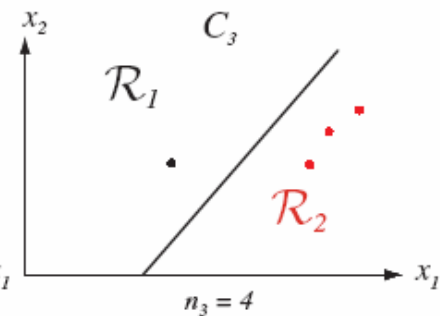
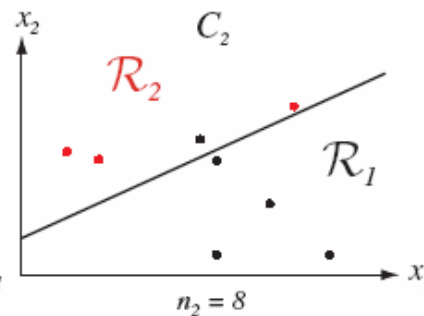
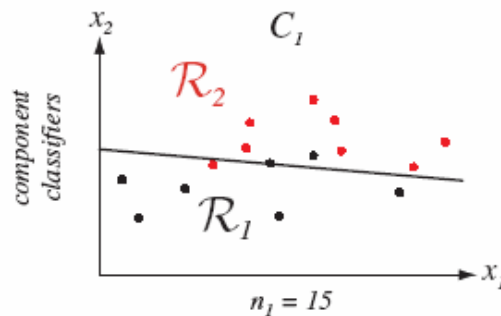
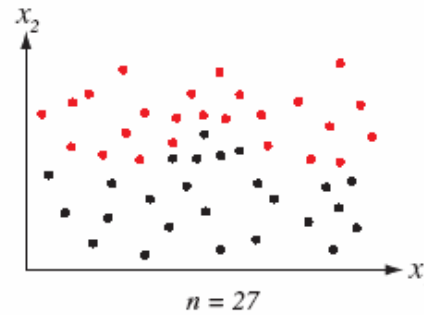
Boosting versus Bagging

- Base classifiers are trained in sequence
- Each base classifier trained using a weighted form of the dataset
 - Weighting coefficient depends on the performance of the previous classifiers
 - Points misclassified by one of the classifiers is given more weight when used to train next classifier
- Decisions are combined using a weighted majority weighting scheme

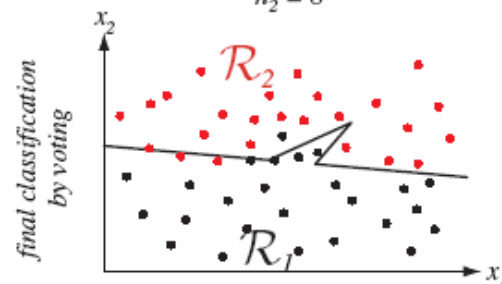
Boosting with 2D two-category problem

- Component classifiers trained using LMS
- Training sets selected by basic boosting procedure

Classification Task



Voting of Components



Creating Classifier Sequence

- Randomly select $n_1 < n$ patterns form full training set D
- Train the first classifier C_1 with D_1
- C_1 only need be a weak learner, accuracy only slightly better than chance
- Seek to find second learning set D_2 that is most “informative”

Creating Second Training Set

- Informative set will be taken from the remaining samples from D
 - such that half are classified correctly and half are classified incorrectly by C_1
- Can be done by flipping a coin
 - if heads choosing a misclassified sample
 - if tails a correctly classified sample
- Now train a second component classifier C_2 with D_2

Creating third data set

- Seek data set D_3 which are not well classified by voting by C_1 and C_2
- Select samples remaining in D and include them only if C_1 and C_2 disagree
- Use of the ensemble of three classifiers:
 - If C_1 and C_2 agree then use that label
 - If they disagree use the label given by C_3
- Choice of n_1
 - Roughly equal numbers for the three classifiers

Boosting

- Boosting be recursively applied to the classifiers themselves creating more classifiers
 - With 9 or 27 classifiers very low training error can be achieved
- Continue adding weak learners until some desired low error has been achieved
- Focuses on difficult patterns by giving them more weight

Adaptive Boosting

- Continue adding weak learners until some desired low error has been achieved
- Each training pattern receives a weight that determines its probability of being selected
- Lower weight for those correctly classified
- Initially uniform weights across training set
- On iteration k draw training set according to weights and train classifier C_k
- Repeat to get C_{k+1}
- Until sufficient number k_{max} of classifiers are generated

Algorithm AdaBoost

Class label

Weights uniform

1 begin initialize $\mathcal{D} = \{\mathbf{x}^1, y_1, \dots, \mathbf{x}^n, y_n\}, k_{max}, W_1(i) = 1/n, i = 1, \dots, n$

2 $k \leftarrow 0$

Error measured w.r.t. distrib W

3 do $k \leftarrow k + 1$

4 train weak learner C_k using \mathcal{D} sampled according to $W_k(i)$

5 $E_k \leftarrow$ training error of C_k measured on \mathcal{D} using $W_k(i)$

6 $\alpha_k \leftarrow \frac{1}{2} \ln[(1 - E_k)/E_k]$ As E increases α decreases

7 $W_{k+1}(i) \leftarrow \frac{W_k(i)}{Z_k} \times \begin{cases} e^{-\alpha_k} & \text{if } h_k(\mathbf{x}^i) = y_i \text{ (correctly classified)} \\ e^{\alpha_k} & \text{if } h_k(\mathbf{x}^i) \neq y_i \text{ (incorrectly classified)} \end{cases}$

8 until $k = k_{max}$ Class decision of C_k

9 return C_k and α_k for $k = 1$ to k_{max} (ensemble of classifiers with weights)

10 end

Normalizing constant to ensure W is a distrib

AdaBoost Ensemble Decision

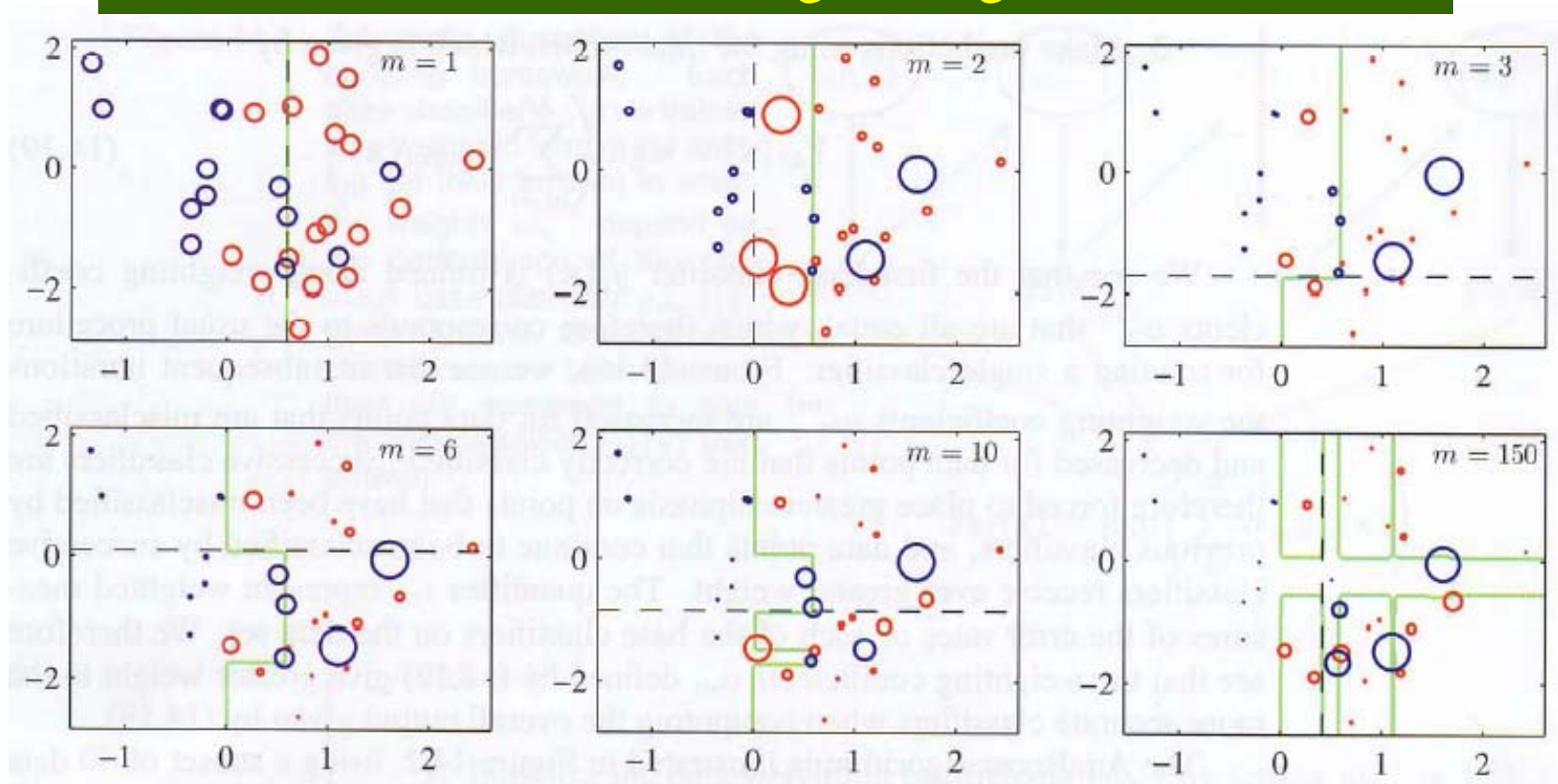
- Final classification is given by

$$g(x) = \sum_{k=1}^{k_{\max}} \alpha_k h_k(x)$$

- In two class case this is same as $\text{sgn}[g(x)]$

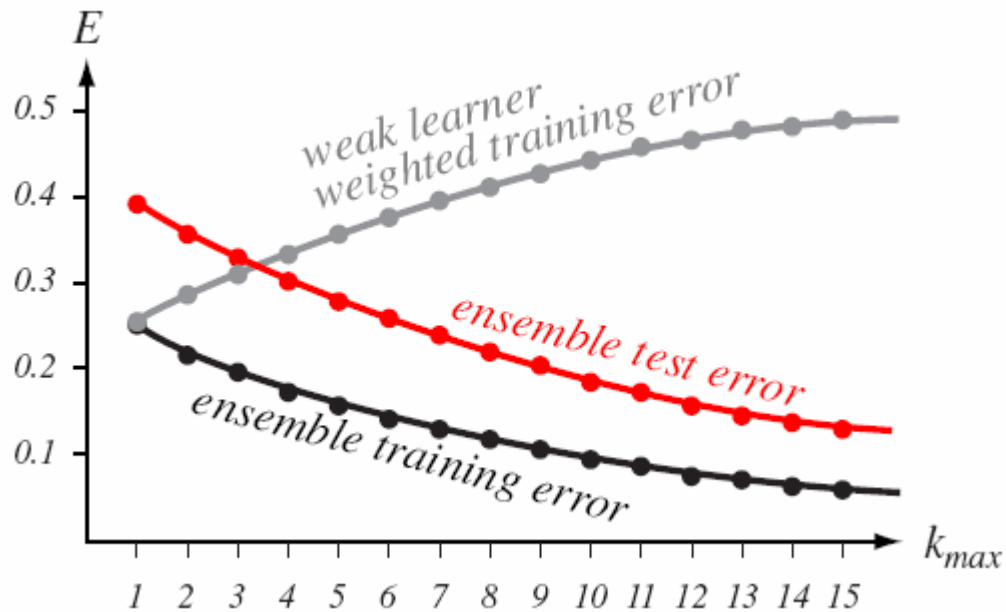
Illustration of AdaBoost

- Weak Learner: Threshold applied to one or other axis
- Ensemble decision boundary is Green
- Current base learner is dashed black line
- Size of circle indicates weight assigned



AdaBoost Error Rates

- Ensemble Error rate keeps decreasing.
- Individual Classifiers perform worse on weighted training sets



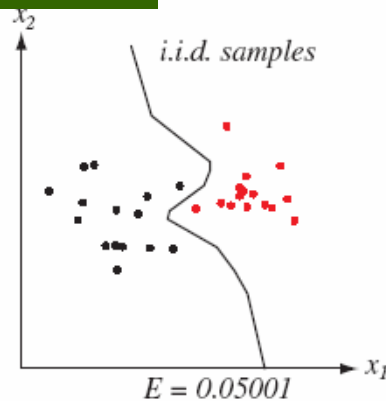
Learning with Queries

- Data Set is unlabeled
- There exists a process for labeling that is costly
 - Handwritten numeral classification
- Present a pattern to an oracle (teacher) who will label pattern without error

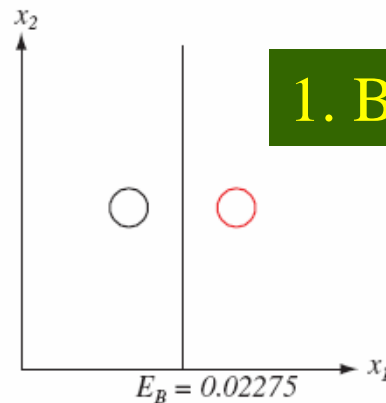
Learning with Queries

Active learning can result in more accurate classifiers than with i.i.d. sampling

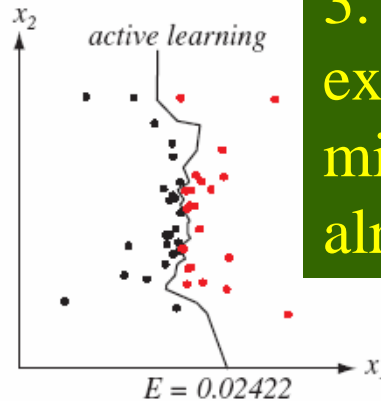
2. Nearest_neighbor with 30 samples and i.i.d. sampling



1. Bayes Classifier for Two Gaussians



3. Queries posed for the extremes of feature space midway between two points already selected



References

- Duda, Hart and Stork 2001
- Bishop, 2006