

**CS 2750 Machine Learning  
Lecture 21**

**Ensemble methods.  
Boosting**

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

**Administrative announcements**

- **Term projects:**
  - Reports due on Wednesday, April 23, 2003 at 2pm.
  - Presentations on Friday, April 25, 2003 at 1pm.
- **Quiz**
  - Monday, April 21, 2003
  - Closed book
  - Short (~30 minutes)
  - Main ideas of methods covered after the midterm
    - EM, Dimensionality reduction, Clustering, Non-parametric methods, Decision trees, Mixtures of experts, Bagging and Boosting, Reinforcement learning.

---

CS 2750 Machine Learning

## Ensemble methods

- **Mixture of experts**
  - Different ‘base’ models (classifiers, regressors) cover different parts of the input space
- **Committee machines:**
  - Train several ‘base’ models on the complete input space, but on slightly different train sets
  - Combine their decision to produce the final result
  - **Goal:** Improve the accuracy of the ‘base’ model
  - **Methods:**
    - **Bagging**
    - **Boosting**
    - Stacking (not covered)

---

CS 2750 Machine Learning

## Bagging (Bootstrap Aggregating)

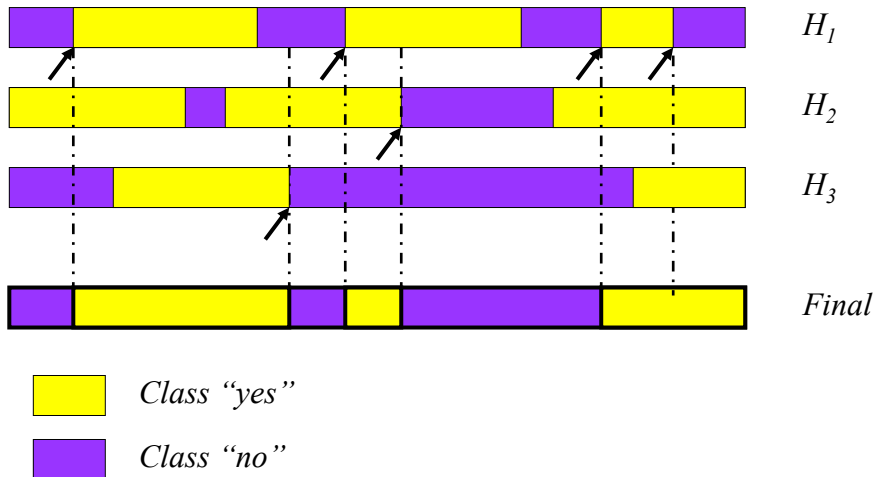
- **Given:**
  - Training set of  $N$  examples
  - A class of learning models (e.g. decision trees, neural networks, ...)
- **Method:**
  - Train multiple ( $k$ ) models on different samples (data splits) and average their predictions
  - Predict (test) by averaging the results of  $k$  models
- **Goal:**
  - Improve the accuracy of one model by using multiple copies of the model
  - Average of misclassification errors on different data splits gives a better estimate of the predictive ability of a learning method

---

CS 2750 Machine Learning

## Simple Majority Voting

Test examples



CS 2750 Machine Learning

## Analysis of Bagging

- **Expected error= Bias+Variance**

- *Expected error* is the expected discrepancy between the estimated and true function

$$E \left[ \left( \hat{f}(X) - E[f(X)] \right)^2 \right]$$

- *Bias* is squared discrepancy between *averaged* estimated and true function

$$\left( E[\hat{f}(X)] - E[f(X)] \right)^2$$

- *Variance* is expected divergence of the estimated function vs. its average value

$$E \left[ \left( \hat{f}(X) - E[\hat{f}(X)] \right)^2 \right]$$

CS 2750 Machine Learning

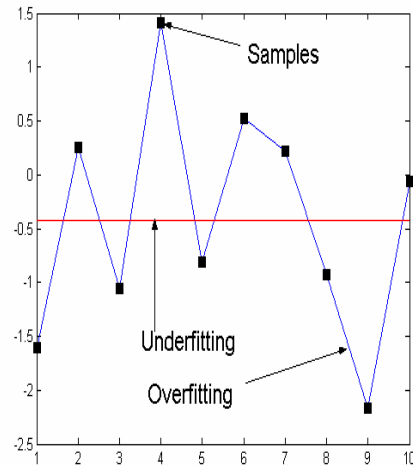
## When Bagging works? Under-fitting and over-fitting

- **Under-fitting:**

- High bias (models are not accurate)
- Small variance (smaller influence of examples in the training set)

- **Over-fitting:**

- Small bias (models flexible enough to fit well to training data)
- Large variance (models depend very much on the training set)



CS 2750 Machine Learning

## Averaging decreases variance

- **Example**

- Assume we measure a random variable  $x$  with a  $N(\mu, \sigma^2)$  distribution
- If only one measurement  $x_1$  is done,
  - The expected mean of the measurement is  $\mu$
  - Variance is  $\text{Var}(x_1) = \sigma^2$
- If random variable  $x$  is measured  $K$  times ( $x_1, x_2, \dots, x_k$ ) and the value is estimated as:  $(x_1 + x_2 + \dots + x_k) / K$ ,
  - Mean of the estimate is still  $\mu$
  - But, variance is smaller:
    - $[\text{Var}(x_1) + \dots + \text{Var}(x_k)] / K^2 = K\sigma^2 / K^2 = \sigma^2 / K$

- Observe: **Bagging is a kind of averaging!**

CS 2750 Machine Learning

## When Bagging works

- **Main property of Bagging** (proof omitted)
  - Bagging decreases variance of the base model without changing the bias!!!
  - Why? averaging!
- **Bagging typically helps**
  - When applied with an over-fitted base model
    - High dependency on actual training data
- **It does not help much**
  - High bias. When the base model is robust to the changes in the training data (due to sampling)

---

CS 2750 Machine Learning

## Boosting. Theoretical foundations.

- **PAC: Probably Approximately Correct framework**
  - **( $\epsilon$ - $\delta$ ) solution**
- **PAC learning:**
  - Learning with the pre-specified accuracy  $\epsilon$  and confidence  $\delta$
  - **the probability that the misclassification error is larger than  $\epsilon$  is smaller than  $\delta$**

$$P(ME(c) > \epsilon) \leq \delta$$

- **Accuracy ( $\epsilon$ ):** Percent of correctly classified samples in test
- **Confidence ( $\delta$ ):** The probability that in one experiment some accuracy will be achieved

---

CS 2750 Machine Learning

## PAC Learnability

### Strong (PAC) learnability:

- There exists a learning algorithm that **efficiently** learns the classification with a pre-specified **accuracy and confidence**

### Strong (PAC) learner:

- A learning algorithm  $P$  that given an arbitrary
  - classification error  $\epsilon$  ( $<1/2$ ), and
  - confidence  $\delta$  ( $<1/2$ )
- Outputs a classifier
  - With a classification accuracy  $> (1-\epsilon)$
  - A confidence probability  $> (1-\delta)$
  - And runs in time polynomial in  $1/\delta, 1/\epsilon$ 
    - Implies: number of samples  $N$  is polynomial in  $1/\delta, 1/\epsilon$

## Weak Learner

### Weak learner:

- A learning algorithm (learner)  $W$ 
  - Providing classification accuracy  $>1-\epsilon_0$
  - With probability  $>1-\delta_0$
- For some **fixed and uncontrollable**
  - classification error  $\epsilon_0$  ( $<1/2$ )
  - confidence  $\delta_0$  ( $<1/2$ )

**on an arbitrary problem**

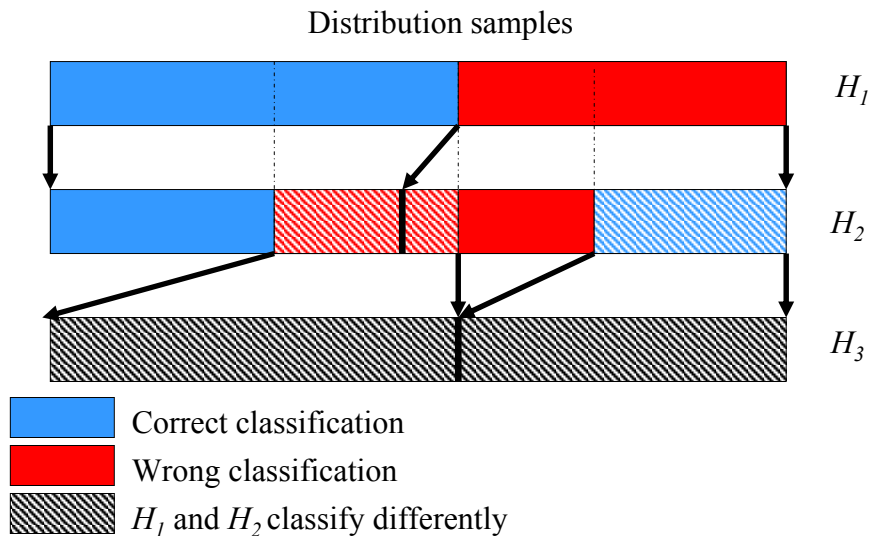
## Weak learnability=Strong (PAC) learnability

- Assume there exists a **weak learner**
  - **On any problem** it is better than a random guess (50 %) with confidence higher than 50 %
- **Question:**
  - Is problem also PAC-learnable?
  - Can we generate an algorithm  $P$  that achieves an arbitrary  $(\epsilon-\delta)$  accuracy?
- **Why is important?**
  - Usual classification methods (decision trees, neural nets), have specified, but uncontrollable performances.
  - Can we improve performance to achieve pre-specified accuracy (confidence)?

## Weak=Strong learnability!!!

- **Proof due to R. Schapire**
    - An arbitrary  $(\epsilon-\delta)$  improvement is possible
  - **Idea:** combine multiple weak learners together
    - Weak learner  $W$  with confidence  $\delta_0$  and maximal error  $\epsilon_0$
    - It is possible:
      - To improve (boost) the confidence
      - To improve (boost) the accuracy
- by training different weak learners on slightly different datasets

## Boosting accuracy Training



CS 2750 Machine Learning

## Boosting accuracy

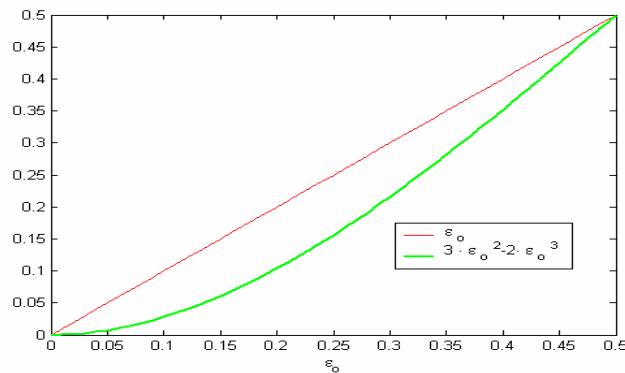
- **Training**
  - Sample randomly from the distribution of examples
  - Train hypothesis  $H_1$  on the sample
  - Evaluate accuracy of  $H_1$  on the distribution
  - Sample randomly such that for the half of samples  $H_1$  provides correct, and for another half, incorrect results; Train hypothesis  $H_2$ .
  - Train  $H_3$  on samples from the distribution where  $H_1$  and  $H_2$  classify differently
- **Test**
  - For each example, decide according to the majority vote of  $H_1$ ,  $H_2$  and  $H_3$

CS 2750 Machine Learning



## Theorem

- If each hypothesis has an error  $\epsilon_0$ , the final classifier has error  $< g(\epsilon_0) = 3\epsilon_0^2 - 2\epsilon_0^3$
- **Accuracy improved !!!!**
- **Apply recursively to get to the target accuracy !!!**



CS 2750 Machine Learning

## Theoretical Boosting algorithm

- Similarly to boosting the accuracy we can boost the confidence at some restricted accuracy cost
- **The key result:** we can improve both the accuracy and confidence
- **Problems with the theoretical algorithm**
  - A good (better than 50 %) classifier on all problems
  - We cannot properly sample from data-distribution
  - Method requires large training set
- **Solution to the sampling problem:**
  - Boosting by sampling
    - **AdaBoost** algorithm and variants

CS 2750 Machine Learning

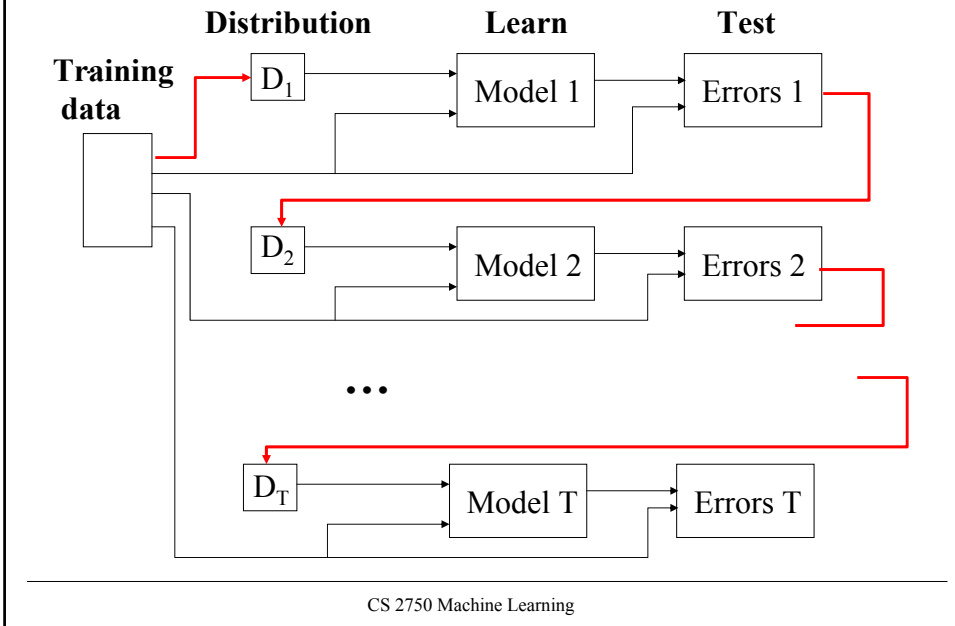
## AdaBoost

- **AdaBoost: boosting by sampling**
- **Classification** (Freund, Schapire; 1996)
  - AdaBoost.M1 (two-class problem)
  - AdaBoost.M2 (multiple-class problem)
- **Regression** (Drucker; 1997)
  - AdaBoostR

## AdaBoost

- **Given:**
  - A training set of  $N$  examples (attributes + class label pairs)
  - A “base” learning model (e.g. a decision tree, a neural network)
- **Training stage:**
  - Train a sequence of  $T$  “base” models on  $T$  different sampling distributions defined upon the training set ( $D$ )
  - A sample distribution  $D_t$  for building the model  $t$  is constructed by modifying the sampling distribution  $D_{t-1}$  from the  $(t-1)$ th step.
    - Examples classified incorrectly in the previous step receive higher weights in the new data (attempts to cover misclassified samples)
- **Application (classification) stage:**
  - **Classify according to the weighted majority** of classifiers

## AdaBoost training



## AdaBoost algorithm

### Training (step t)

- **Sampling Distribution**  $D_t$

$D_t(i)$  - a probability that example  $i$  from the original training dataset is selected

$D_1(i) = 1 / N$  for the first step ( $t=1$ )

- Take  $K$  samples from the training set according to  $D_t$
- Train a classifier  $h_t$  on the samples

- Calculate the error  $\epsilon_t$  of  $h_t$ : 
$$\epsilon_t = \sum_{i: h_t(x_i) \neq y_i} D_t(i)$$

- Classifier weight:  $\beta_t = \epsilon_t / (1 - \epsilon_t)$

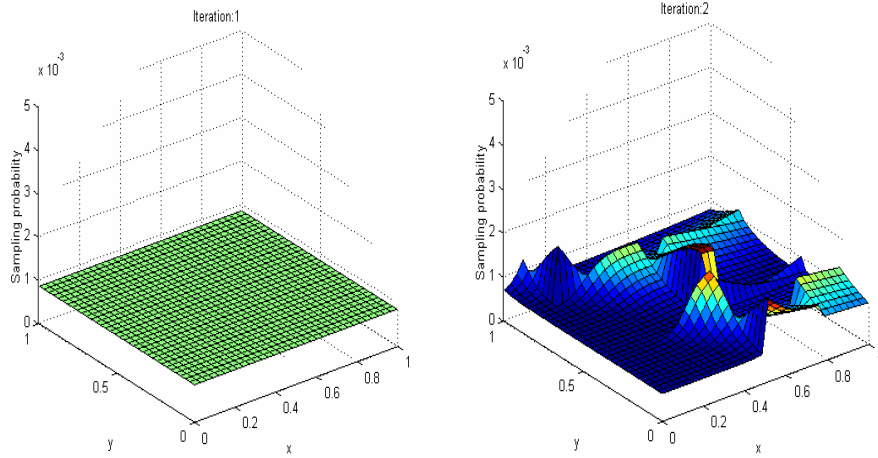
- New sampling distribution

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} \beta_t & h_t(x_i) = y_i \\ 1 & \text{otherwise} \end{cases}$$

Norm. constant

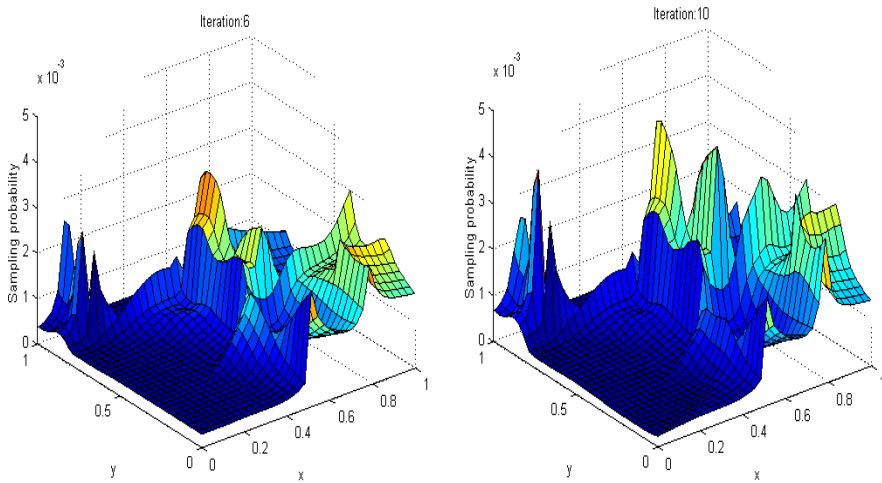
## AdaBoost. Sampling Probabilities

Example: - Nonlinearly separable binary classification  
- NN as weak learners



CS 2750 Machine Learning

## AdaBoost: Sampling Probabilities



CS 2750 Machine Learning

## AdaBoost classification

- We have  $T$  different classifiers  $h_t$ 
  - weight  $w_t$  of the classifier is proportional to its accuracy on the training set

$$w_t = \log(1 / \beta_t) = \log((1 - \epsilon_t) / \epsilon_t)$$

$$\beta_t = \epsilon_t / (1 - \epsilon_t)$$

- **Classification:**

For every class  $j=0,1$

- Compute the sum of weights  $w$  corresponding to ALL classifiers that predict class  $j$ ;
- Output class that correspond to the maximal sum of weights (weighted majority)

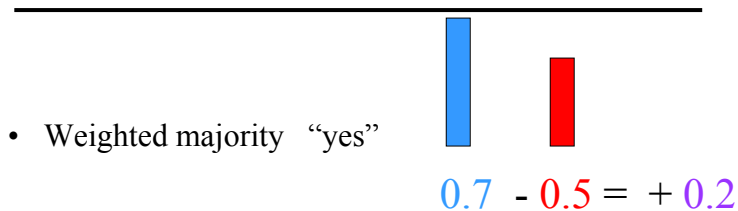
$$h_{final}(\mathbf{x}) = \arg \max_j \sum_{t: h_t(\mathbf{x})=j} w_t$$

---

CS 2750 Machine Learning

## Two-Class example. Classification.

- Classifier 1      “yes”      0.7
- Classifier 2      “no”            0.3
- Classifier 3      “no”            0.2



- The final choose is “yes” + 1

---

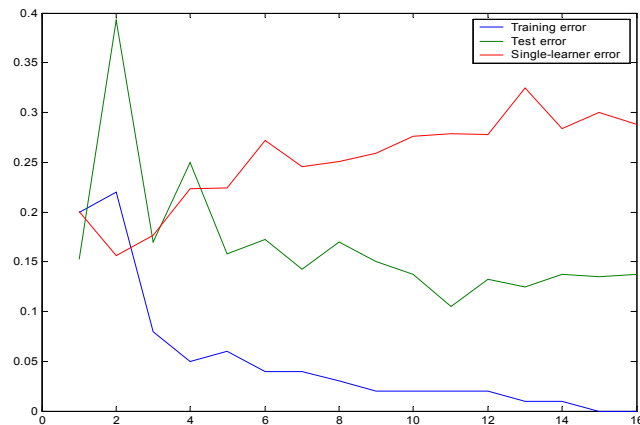
CS 2750 Machine Learning

## What is boosting doing?

- Each classifier specializes on a particular subset of examples
- Algorithm is concentrating on “more and more difficult” examples
- **Boosting can:**
  - Reduce variance (the same as Bagging)
  - But also to eliminate the effect of high bias of the weak learner (unlike Bagging)
- **Train versus test errors performance:**
  - Train errors can be driven close to 0
  - But test errors do not show overfitting
- Proofs and theoretical explanations in **readings**

CS 2750 Machine Learning

## Boosting. Error performances



CS 2750 Machine Learning