

Ensemble learning: Bagging & Boosting

CE-717: Machine Learning

Sharif University of Technology

Prepared by M. Soleymani ; reviewed/Complemented by H. R. Rabiee

Fall 2018

What is ensemble learning?

- ▶ Ensembles combine multiple hypotheses to form a (hopefully) better hypothesis.
 - ▶ combining many *weak learners* in an attempt to produce a *strong learner*.
- ▶ Ensemble term is usually reserved for methods that generate multiple hypotheses using the **same base learner**.
- ▶ Multiple classifier (broader term) also covers **combination of hypotheses** that are not induced by the same base learner.

Ensemble learning

- ▶ We only talk about:
 - ▶ **Bagging: Bootstrap aggregating**
 - ▶ **Boosting**
 - ▶ One important committee method
 - ▶ The most famous boosting algorithm: AdaBoost

Bias-variance trade-off

- ▶ Weak or simple learners
 - ▶ Low variance: they don't usually overfit
 - ▶ High bias: they can't learn complex functions
- ▶ Boosting to decrease the bias
 - ▶ boost weak learners to enhance their capabilities
- ▶ Bagging can decrease the variance

Bagging algorithm (Breiman, 96)

- ▶ Each member of the ensemble is constructed from a different training dataset
 - ▶ samples training data uniformly at random with replacement
- ▶ Predictions combined either by uniform averaging or voting over class labels.
 - ▶ works best with unstable models (high variance models)
- ▶ Despite its apparent simplicity, Bagging is still not fully understood

Bootstrap Sampling

- ▶ Bootstrap sampling: Samples the given dataset N times uniformly with replacement (resulting in a set of N samples)
 - ▶ Some samples in the original set may be included several times in the bootstrap sampled data
- ▶ Bootstrap sampling: like “roll N -face dice N times”

Bagging algorithm

- ▶ **Input:** Required ensemble size M

Training set $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$

- ▶ **for** $t=1$ **to** T **do**

- ▶ Build a dataset D_t by sampling N items, randomly with replacement from D
- ▶ Train a model h_t using D_t , and add it to the ensemble.

- ▶ $H(\mathbf{x}) = \text{sign}(\sum_{t=1}^T h_t(\mathbf{x}))$

- ▶ Combine models by voting for classification and by averaging for regression

Bagging on decision trees

- ▶ Decision trees are popular classifiers:
 - ▶ interpretable
 - ▶ can handle discrete and continuous features
 - ▶ robust to outliers
 - ▶ low bias
- ▶ However, they are high variance
- ▶ Trees are perfect candidates for ensembles
 - ▶ Consider averaging many (nearly) unbiased tree estimators
 - ▶ Bias remains similar, but variance is reduced
 - ▶ This is called bagging
 - ▶ Train many trees on bootstrapped data, then average outputs

Random Forest

- ▶ Bagging on decision trees
 - ▶ Reduce correlation between trees, by introducing randomness
 - ▶ For $b = 1, \dots, B$,
 - ▶ Draw a bootstrap dataset
 - ▶ Learn a tree on this dataset
 - Select m features randomly out of d features as candidates before splitting
 - ▶ Output:
 - ▶ Regression: average of outputs
 - ▶ Classification: majority vote
- Usually: $m \leq \sqrt{d}$

Boosting idea

- ▶ We can select simple “weak” classification or regression methods and combine them into a single “strong” method
- ▶ Examples of weak classifiers: Naïve bayes, logistic regression, decision stumps or shallow decision trees
- ▶ Learn many weak classifiers that are good at different parts of the input space.
- ▶ To find the output class, find weighted vote of classifiers

Boosting

- ▶ Try to combine many simple “weak” classifiers (in sequence) to find a single “strong” classifier
 - ▶ Each component is a simple binary ± 1 classifier
 - ▶ Voted combinations of component classifiers

$$H_m(\mathbf{x}) = \alpha_1 h(\mathbf{x}; \boldsymbol{\theta}_1) + \dots + \alpha_t h(\mathbf{x}; \boldsymbol{\theta}_t)$$

- ▶ To simplify notation: $h(\mathbf{x}; \boldsymbol{\theta}_i) = h_i(\mathbf{x})$ $\alpha_i \geq 0$ are higher for more reliable classifiers

$$H_m(\mathbf{x}) = \alpha_1 h_1(\mathbf{x}) + \dots + \alpha_t h_t(\mathbf{x})$$

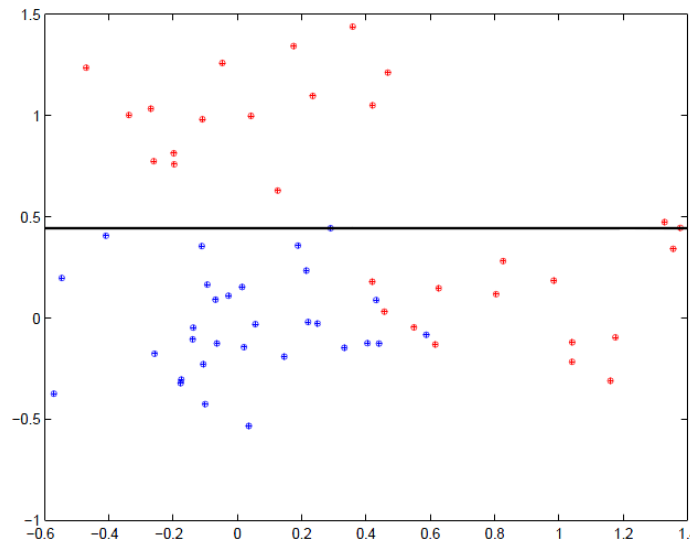
- ▶ Prediction: $\hat{y} = \text{sign}(H_t(\mathbf{x}))$

Simple component classifiers

- ▶ Simple family of component classifiers (called decision stumps):

$$h(\mathbf{x}; \boldsymbol{\theta}) = \text{sign}(w_1 x_k - w_0) \quad \boldsymbol{\theta} = \{k, w_1, w_0\}$$

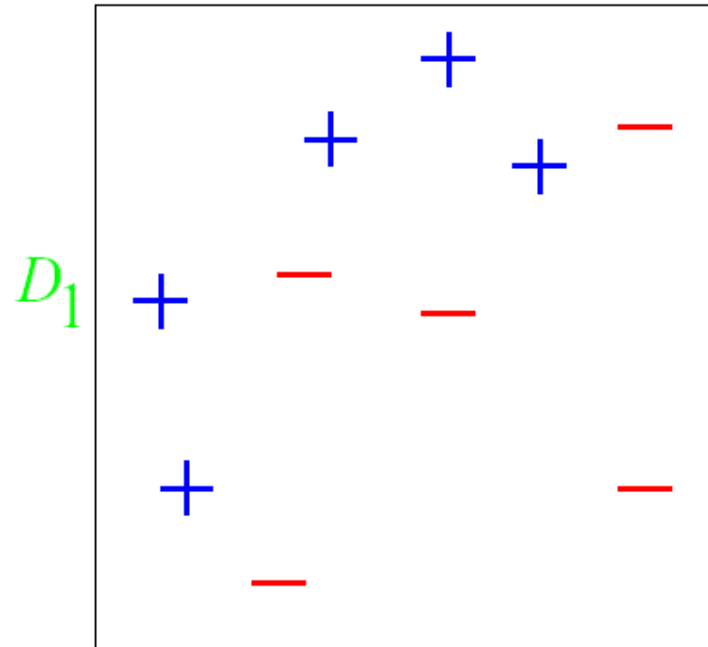
- ▶ Each classifier is based on only a single feature of \mathbf{x} (e.g., x_k): decision tree of depth one



AdaBoost: basic ideas

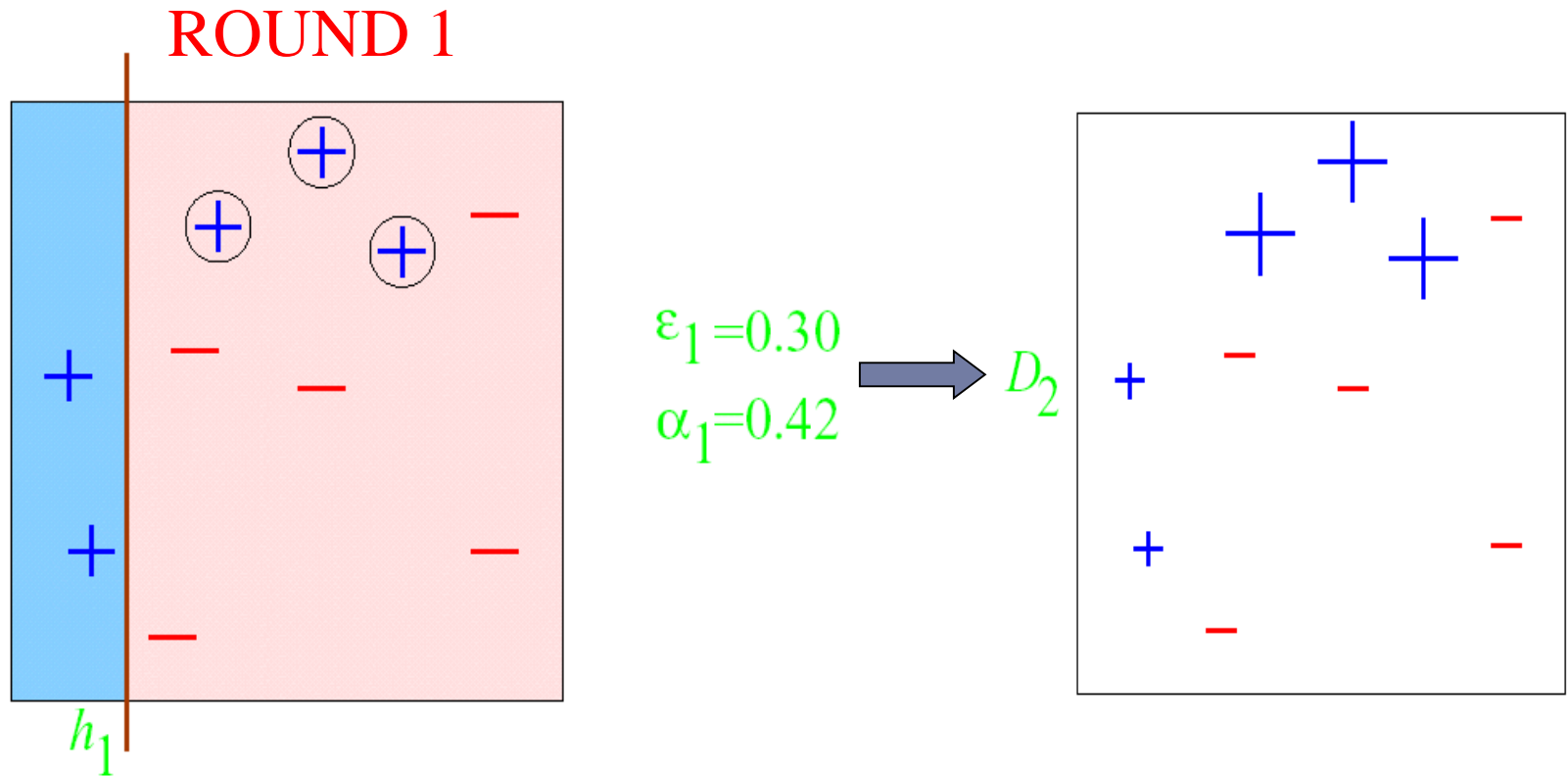
- ▶ Sequential production of classifiers
 - ▶ choose classifier whose addition will be most helpful.
- ▶ Each classifier is dependent on the previous ones
 - ▶ focuses on the previous ones' errors
- ▶ Incorrectly predicted samples in previous classifiers are weighted more heavily

AdaBoost: example

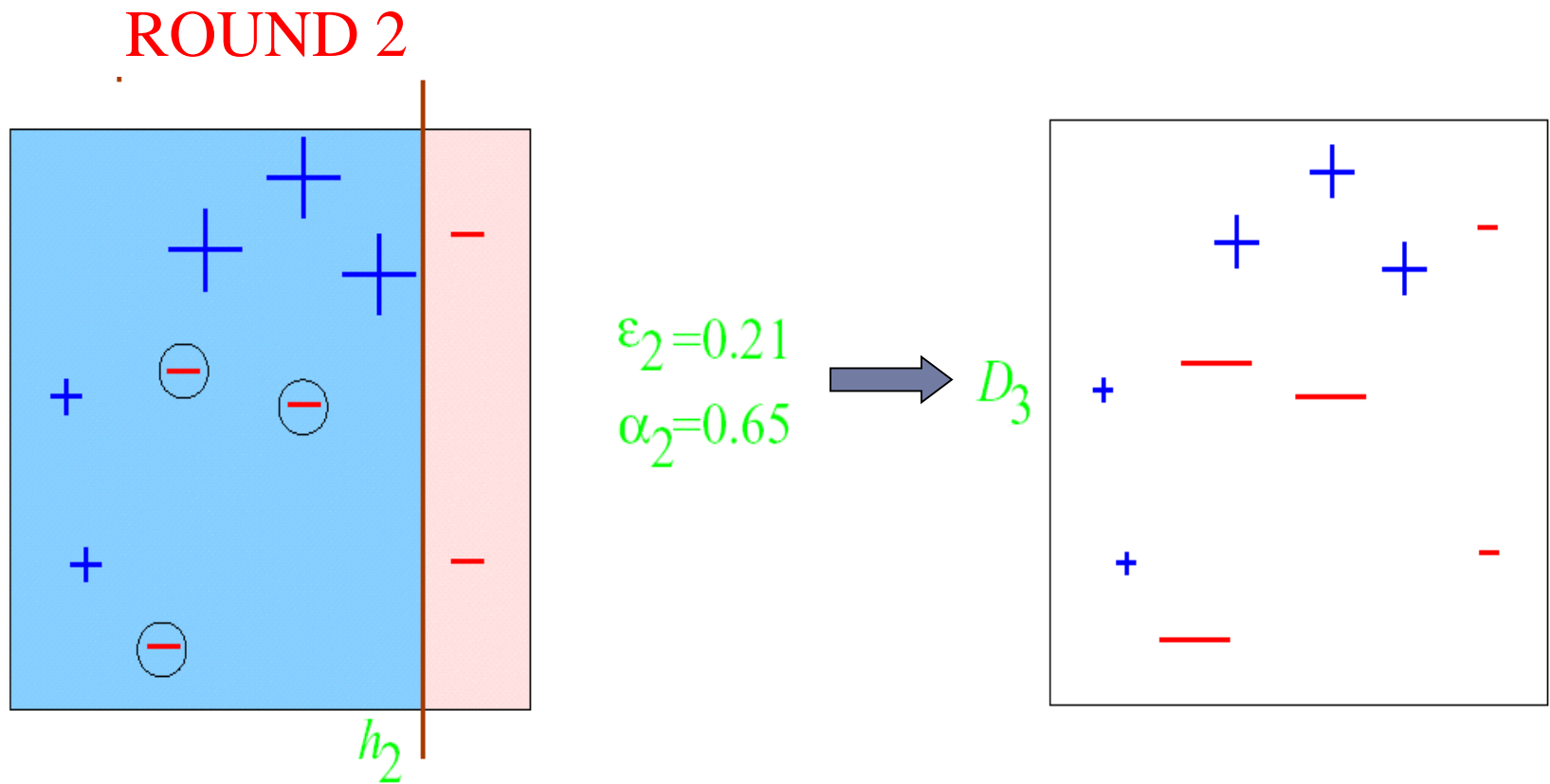


Equal Weights to all training samples

AdaBoost: example

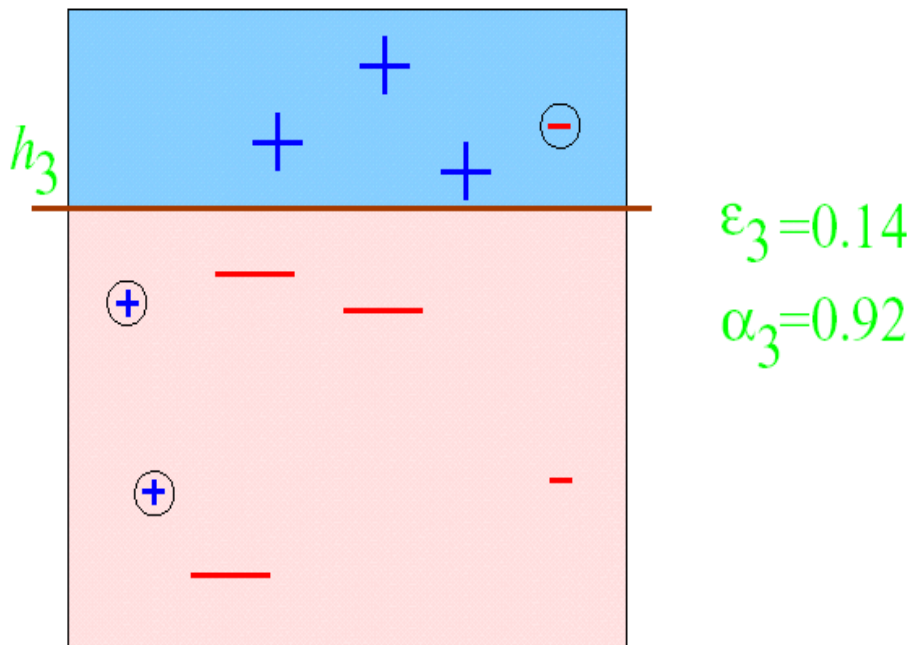


AdaBoost: example



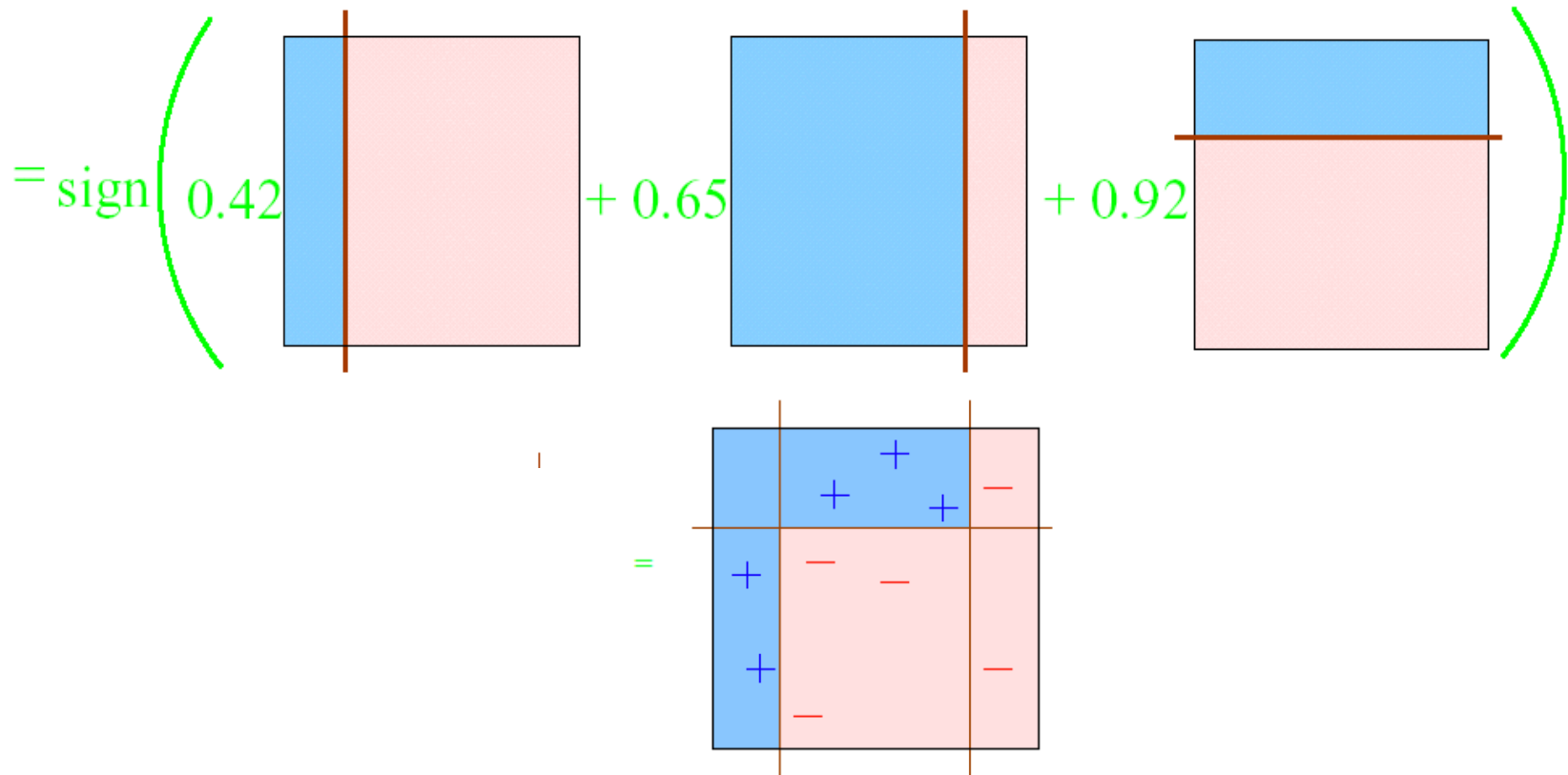
AdaBoost: example

ROUND 3

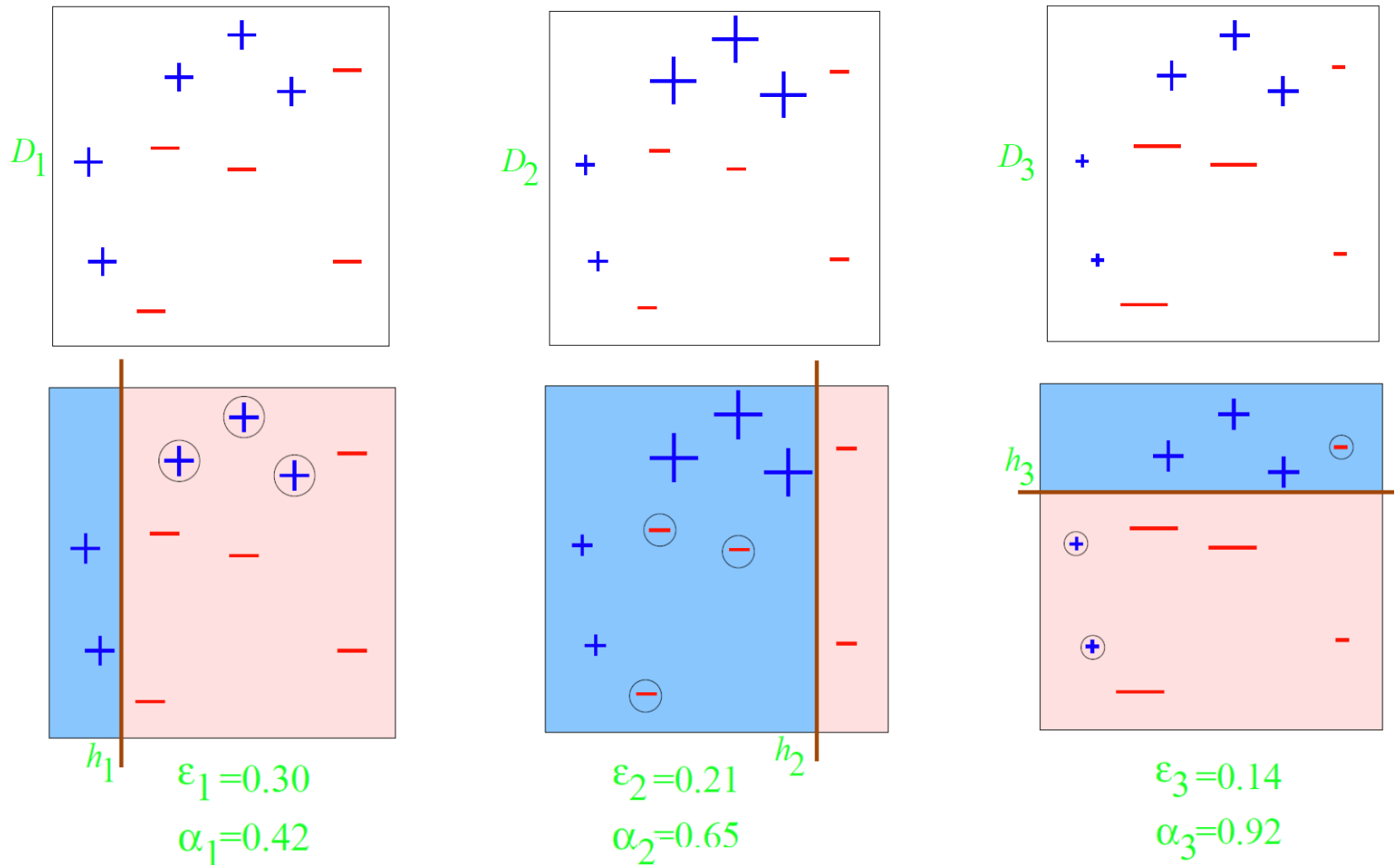


AdaBoost: example

H_{final}



AdaBoost: example



AdaBoost algorithm

1) For $i = 1$ to N Initialize the data weight $w_1^{(i)} = \frac{1}{N}$

2) For $t = 1$ to T

a) Find a classifier $h_t(\mathbf{x})$ by minimizing the weighted error function:

$$J_t = \sum_{i=1}^N w_t^{(i)} \times I(y^{(i)} \neq h_t(\mathbf{x}^{(i)}))$$

Only when $h_t(\mathbf{x})$ with $\epsilon_t < 0.5$ (better than chance) is found, boosting continues.

b) Find the weighted error of $h_t(\mathbf{x})$:

$$\epsilon_t = \frac{\sum_{i=1}^N w_t^{(i)} \times I(y^{(i)} \neq h_t(\mathbf{x}^{(i)}))}{\sum_{i=1}^N w_t^{(i)}}$$

and the new component is assigned votes based on its error:

$$\alpha_t = \ln((1 - \epsilon_t)/\epsilon_t)$$

c) The normalized weights are updated:

$$w_{t+1}^{(i)} = w_t^{(i)} e^{\alpha_t I(y^{(i)} \neq h_t(\mathbf{x}^{(i)}))}$$

3) Combined classifier $\hat{y} = \text{sign}(H_T(\mathbf{x}))$ where $H_T(\mathbf{x}) = \sum_{t=1}^M \alpha_t h_t(\mathbf{x})$



Notation explanation

- ▶ $w_t^{(i)}$: Weighting coefficient of data point i in iteration t
- ▶ α_t : weighting coefficient of t -th base classifier in the final ensemble
 - ▶ ϵ_t : weighted error rate of t -th base classifier

Boosting: main ideas

- ▶ Boosting algorithms maintain weights on training data:
 - ▶ Initially, all weights are equal, $w_1^{(i)} = 1/N$.
 - ▶ In t -th iteration, the weights are updated:
 - ▶ If $x^{(i)}$ is misclassified by h_t , $w_t^{(i)}$ goes up;
- ▶ Fitting of h_{t+1} is guided by weights of samples
 - ▶ Force h_{t+1} to focus on already misclassified examples.

Adaptive boosting (AdaBoost): intuition

- ▶ First iteration: a usual procedure for training a single (weak) classifier
- ▶ Subsequent iterations:
 - ▶ $w_t^{(i)}$ are increased for misclassified data points
 - ▶ Then, successive classifiers are forced to place greater emphasis on points misclassified by previous classifiers

Boosting: loss function

- ▶ We need a loss function for the combination
 - ▶ determine which new component $h(\mathbf{x}; \boldsymbol{\theta})$ to add
 - ▶ and how many votes it should receive

$$H_t(\mathbf{x}) = \frac{1}{2} (\alpha_1 h_1(\mathbf{x}) + \dots + \alpha_t h_t(\mathbf{x}))$$

- ▶ Many options for the loss function
 - ▶ AdaBoost is equivalent to using the following exponential loss

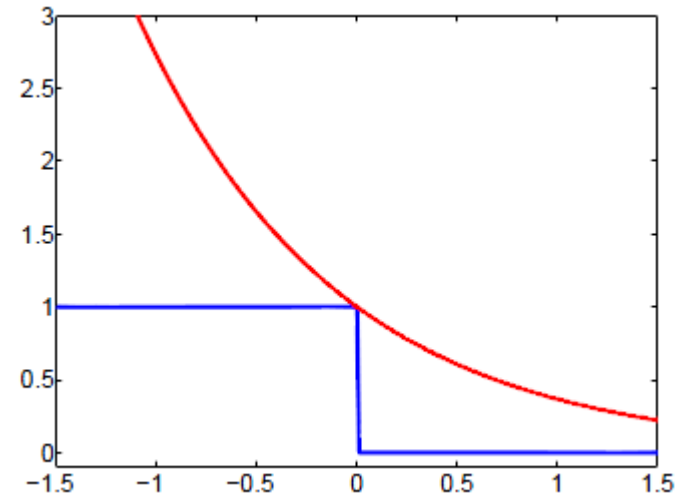
$$Loss(y, H_t(\mathbf{x})) = e^{-y \times H_t(\mathbf{x})}$$

$$\hat{y} = \text{sign}(H_t(\mathbf{x}))$$

- ▶ A simple interpretation of boosting in terms of the sequential minimization of the exponential loss function [Friedman *et al.*, 2000].

Boosting: exponential loss function

- ▶ Differentiable approximation (bound) of 0/1 loss
 - ▶ Easy to optimize
 - ▶ Optimizing an upper bound on classification error.
- ▶ Other options are possible.



$$H_t(\mathbf{x}) = \frac{1}{2} [\alpha_1 h_1(\mathbf{x}) + \dots + \alpha_t h_t(\mathbf{x})]$$

AdaBoost: loss function

- Consider adding the t -th component:

$$E = \sum_{i=1}^N e^{-y^{(i)} H_t(\mathbf{x}^{(i)})} = \sum_{i=1}^N e^{-y^{(i)} [H_{t-1}(\mathbf{x}^{(i)}) + \frac{1}{2} \alpha_t h_t(\mathbf{x}^{(i)})]}$$

$$= \sum_{i=1}^N \underbrace{e^{-y^{(i)} H_{t-1}(\mathbf{x}^{(i)})}}_{\text{Suppose it is fixed at stage } t} e^{-\frac{1}{2} \alpha_t y^{(i)} h_t(\mathbf{x}^{(i)})}$$

Suppose it is fixed at stage t

$$= \sum_{i=1}^N w_t^{(i)} \underbrace{e^{-\frac{1}{2} \alpha_t y^{(i)} h_t(\mathbf{x}^{(i)})}}_{\text{Need to be optimized at stage } t \text{ by seeking } h_t(\mathbf{x}) \text{ and } \alpha_t}$$

$$w_t^{(i)} = e^{-y^{(i)} H_{t-1}(\mathbf{x}^{(i)})}$$

Need to be optimized at stage t
by seeking $h_t(\mathbf{x})$ and α_t

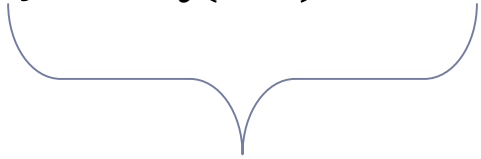
Weighted exponential loss

$$E = \sum_{i=1}^N w_t^{(i)} e^{-\alpha_t y^{(i)} h_t(x^{(i)})}$$

- ▶ sequentially adds a new component trained on reweighted training samples
- ▶ $w_t^{(i)}$: history of classification of $x^{(i)}$ by H_{t-1} .
 - ▶ \Rightarrow Loss weighted towards mistakes
- ▶ Iteration t optimization:
 - ▶ choose the new component $h_t = h(x; \theta_t)$
 - ▶ and the vote α_t that optimize the weighted exponential loss.

Minimizing loss: finding h_t

$$\begin{aligned} E &= \sum_{i=1}^N w_t^{(i)} e^{-\frac{1}{2}\alpha_t y^{(i)} h_t(x^{(i)})} \\ &= e^{\frac{-\alpha_t}{2}} \sum_{y^{(i)}=h_t(x^{(i)})} w_t^{(i)} + e^{\frac{\alpha_t}{2}} \sum_{y^{(i)} \neq h_t(x^{(i)})} w_t^{(i)} \\ &= \left(e^{\frac{\alpha_t}{2}} - e^{\frac{-\alpha_t}{2}} \right) \sum_{y^{(i)} \neq h_t(x^{(i)})} w_t^{(i)} + e^{\frac{-\alpha_t}{2}} \sum_{i=1}^N w_t^{(i)} \end{aligned}$$


$$J_t = \sum_{i=1}^N w_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))$$

Find $h_t(x)$ that
minimizes J_t

Minimizing loss: finding α_m

$$\frac{\partial E}{\partial \alpha_t} = 0$$

$$\Rightarrow \frac{1}{2} \left(e^{\frac{\alpha_t}{2}} + e^{\frac{-\alpha_t}{2}} \right) \sum_{y^{(i)} \neq h_t(\mathbf{x}^{(i)})} w_t^{(i)} - \frac{1}{2} e^{\frac{-\alpha_t}{2}} \sum_{i=1}^N w_t^{(i)} = 0$$

$$\Rightarrow \frac{e^{\frac{-\alpha_t}{2}}}{\left(e^{\frac{\alpha_t}{2}} + e^{\frac{-\alpha_t}{2}} \right)} = \frac{\sum_{y^{(i)} \neq h_t(\mathbf{x}^{(i)})} w_t^{(i)}}{\sum_{i=1}^N w_t^{(i)}}$$

$$\alpha_t = \ln((1 - \epsilon_t)/\epsilon_t)$$

$$\epsilon_t = \frac{\sum_{i=1}^N w_t^{(i)} I(y^{(i)} \neq h_t(\mathbf{x}^{(i)}))}{\sum_{i=1}^N w_t^{(i)}}$$

Updating weights

- ▶ Updating weights in AdaBoost algorithm:

$$\begin{aligned}w_i^{t+1} &= w_i^t e^{-\frac{1}{2}\alpha_t y^{(i)} h_t(x^{(i)})} \\&= w_i^t e^{-\frac{1}{2}\alpha_t} e^{\alpha_t I(y^{(i)} \neq h_t(x^{(i)}))}\end{aligned}$$

$y^{(i)} h_t(x^{(i)}) = 1 - 2I(y^{(i)} \neq h_t(x^{(i)}))$

Independent of i and
can be ignored

$$\Rightarrow w_i^{t+1} = w_i^t e^{\alpha_t I(y^{(i)} \neq h_t(x^{(i)}))}$$

Another perspective for AdaBoost

Define a uniform distribution $D_1(i)$ over elements of S .

for $t = 1$ to T **do**

Train a model h_t using distribution D_t .

Calculate $\epsilon_t = P_{D_t}(h_t(x) \neq y)$

If $\epsilon_t \geq 0.5$ break

Set $\alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$

Update $D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$

where Z_t is a normalization factor so that D_{t+1} is a valid distribution.

end for

For a new testing point (x', y') ,

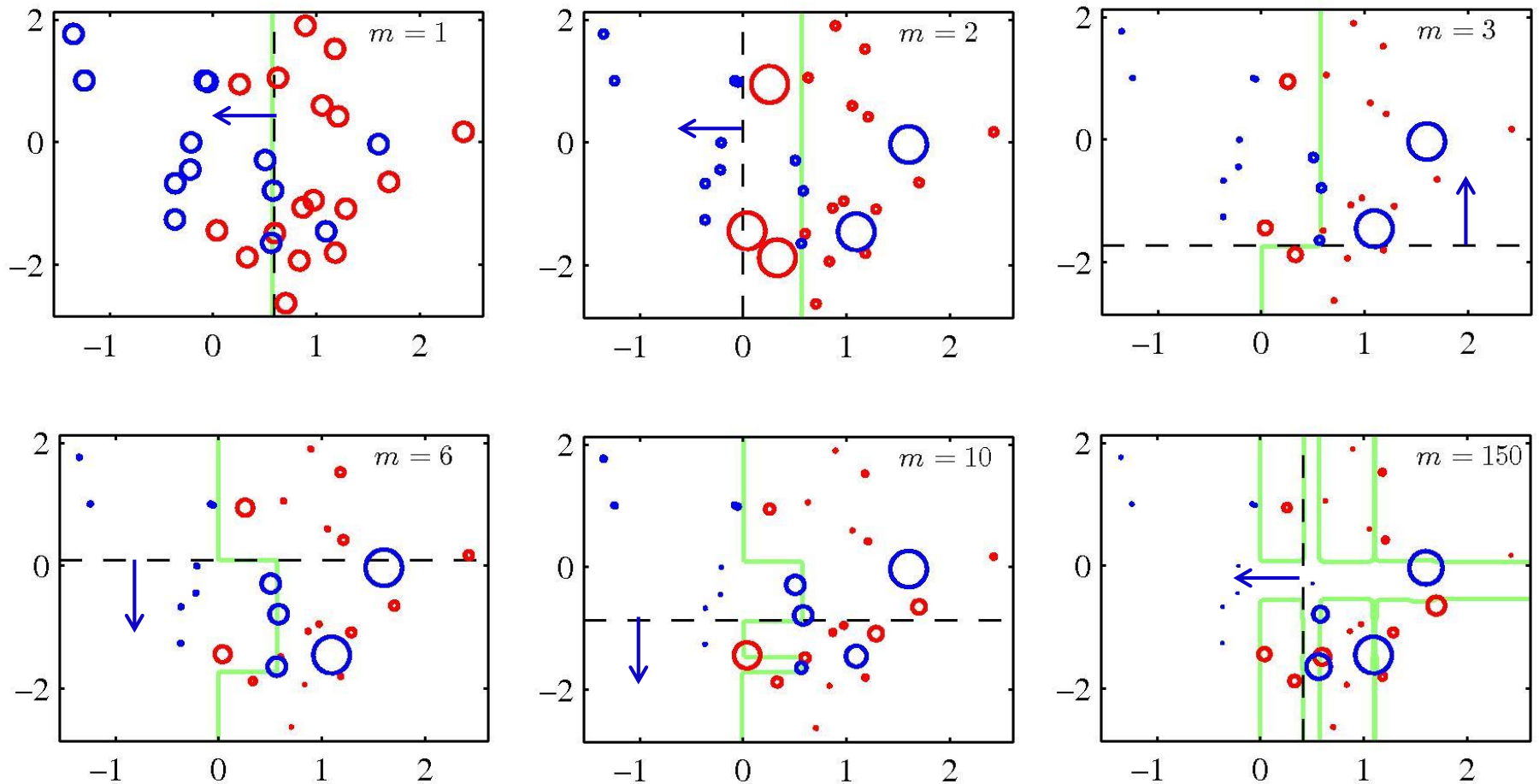
$H(x') = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x'))$

$$H_T(\mathbf{x}) = \frac{1}{2} [\alpha_1 h_1(\mathbf{x}) + \dots + \alpha_T h_T(\mathbf{x})]$$

AdaBoost algorithm: summary

- 1) For $i = 1$ to N Initialize the data weight $w_1^{(i)} = \frac{1}{N}$
- 2) For $t = 1$ to T
 - a) Find a classifier $h_t(x)$ by minimizing the weighted error function
 - b) Find the normalized weighted error of $h_t(x)$ as ϵ_t
 - c) Compute the new component weight as α_t
 - d) Update example weights for the next iteration $w_{t+1}^{(i)}$
- 3) Combined classifier $\hat{y} = \text{sign}(H_T(x))$ where $H_T(x) = \sum_{t=1}^T \alpha_t h_t(x)$

AdaBoost: example 2



How to train base learners

- ▶ Weak learners used in practice:
 - ▶ Decision stumps
 - ▶ Decision trees
 - ▶ Multi-layer neural networks
- ▶ Can base learners operate on weighted examples?
 - ▶ In many cases they can be modified to accept weights along with the examples
 - ▶ In general, we can sample the examples (with replacement) according to the distribution defined by the weights

AdaBoost: typical behavior

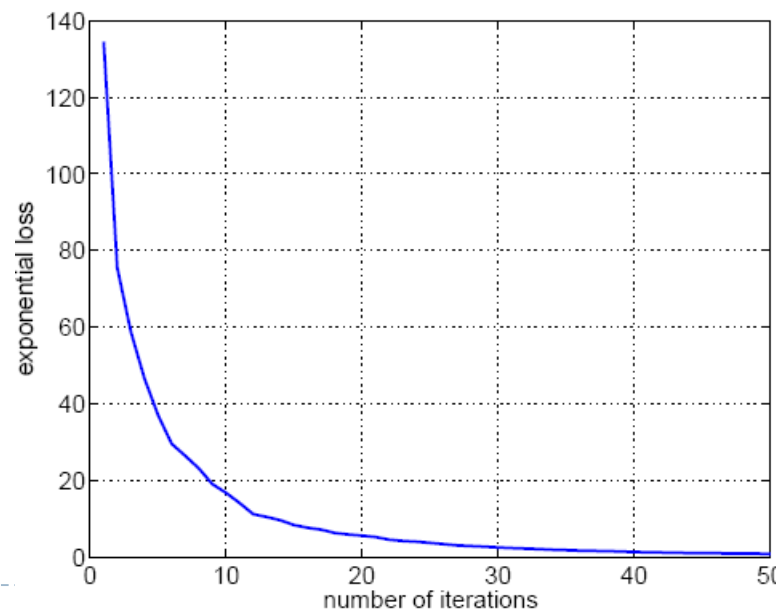
- ▶ Exponential loss goes strictly down.
- ▶ Training error of H goes down
- ▶ Weighted error ϵ_t goes up \Rightarrow votes α_t go down.

AdaBoost properties: exponential loss

- ▶ In each boosting iteration, assuming we can find $h(\mathbf{x}; \hat{\boldsymbol{\theta}}_t)$ whose weighted error is better than chance.

$$H_t(\mathbf{x}) = \frac{1}{2} [\hat{\alpha}_1 h(\mathbf{x}; \hat{\boldsymbol{\theta}}_1) + \dots + \hat{\alpha}_t h(\mathbf{x}; \hat{\boldsymbol{\theta}}_t)]$$

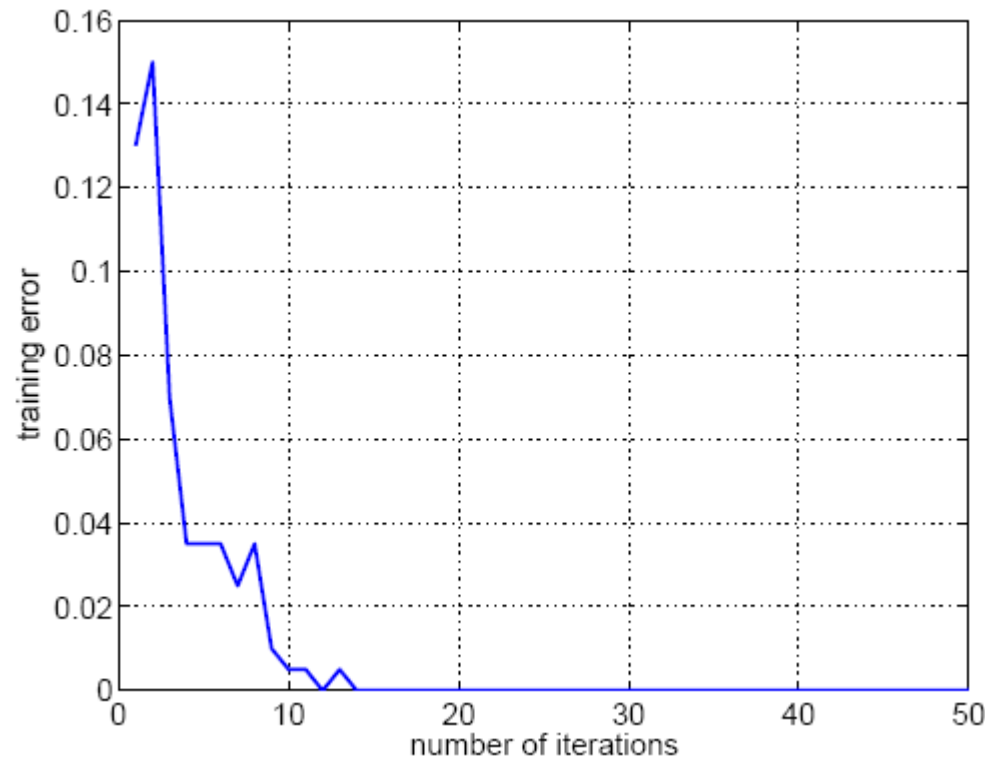
- ▶ Thus, lower exponential loss over training data is guaranteed.



$$E = \sum_{i=1}^N e^{-y^{(i)} H_t(\mathbf{x}^{(i)})}$$

AdaBoost properties: training error

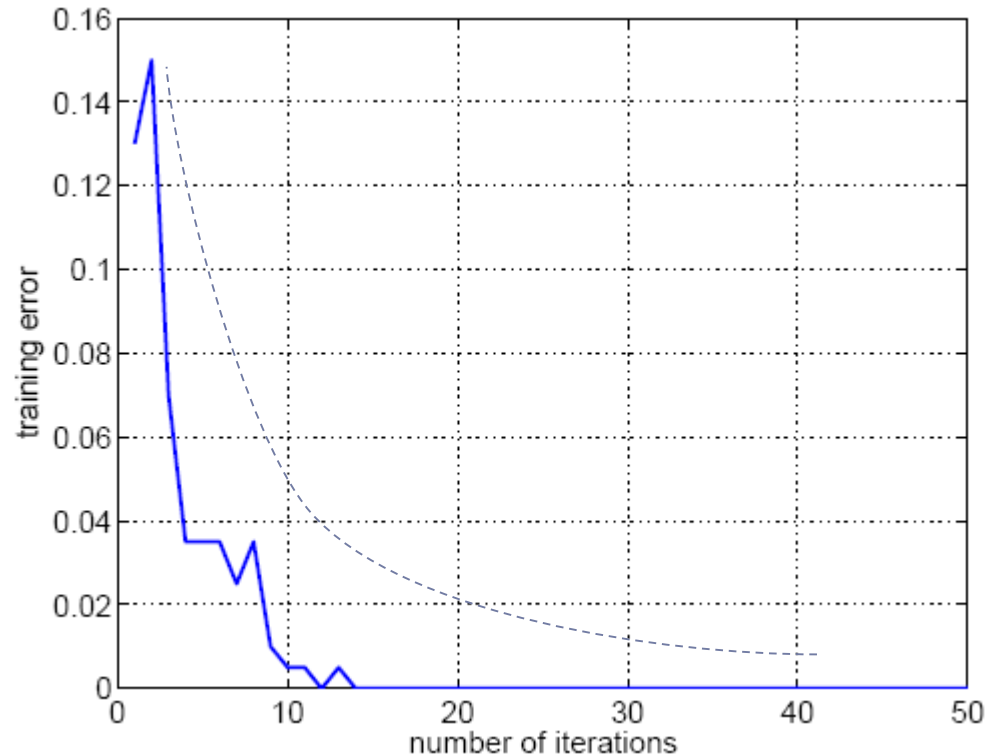
- ▶ Boosting iterations typically decrease the classification error of $H_t(x)$ over training examples.



AdaBoost properties: training error

- ▶ Training error has to go down exponentially fast if the weighted error of each h_t is strictly better than chance (i.e., $\epsilon_t < 0.5$)

Training error in t -th iteration is bounded by an exponential function a^t ($0 < a < 1$)

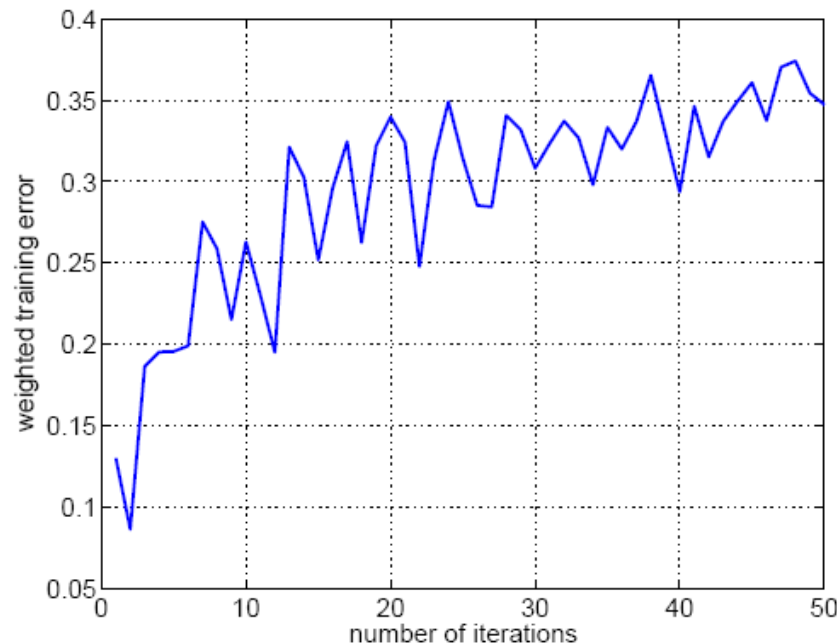


AdaBoost properties: weighted error

- ▶ Weighted error of each new component classifier

$$\epsilon_m = \frac{\sum_{i=1}^n w_m^{(i)} I(y^{(i)} \neq h_m(x^{(i)}))}{\sum_{i=1}^n w_m^{(i)}}$$

tends to increase as a function of boosting iterations.



Updates & Normalization

- Claim: D_{t+1} puts half of the weight on samples on which h_t was incorrect and other half on samples on which h_t was correct

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{-\frac{1}{2}\alpha_t y^{(i)} h_t(x^{(i)})}$$

$$\Pr_{D_{t+1}}[y^{(i)} \neq h_t(x^{(i)})] = \sum_{i: y^{(i)} \neq h_t(x^{(i)})} \frac{D_t(i)}{Z_t} e^{\frac{1}{2}\alpha_t} = \frac{\epsilon_t}{Z_t} e^{\frac{1}{2}\alpha_t} = \frac{\epsilon_t}{Z_t} \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} = \frac{1}{Z_t} \sqrt{\epsilon_t(1 - \epsilon_t)}$$

$$\Pr_{D_{t+1}}[y^{(i)} = h_t(x^{(i)})] = \sum_{i: y^{(i)} = h_t(x^{(i)})} \frac{D_t(i)}{Z_t} e^{-\frac{1}{2}\alpha_t} = \frac{1 - \epsilon_t}{Z_t} e^{-\frac{1}{2}\alpha_t} = \frac{1 - \epsilon_t}{Z_t} \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = \frac{1}{Z_t} \sqrt{\epsilon_t(1 - \epsilon_t)}$$

$$\begin{aligned} Z_t &= \sum_i D_t(i) e^{-\frac{1}{2}\alpha_t y^{(i)} h_t(x^{(i)})} = \sum_{i: y^{(i)} \neq h_t(x^{(i)})} D_t(i) e^{\frac{1}{2}\alpha_t} + \sum_{i: y^{(i)} = h_t(x^{(i)})} D_t(i) e^{-\frac{1}{2}\alpha_t} \\ &= (1 - \epsilon_t) e^{-\frac{1}{2}\alpha_t} + \epsilon_t e^{\frac{1}{2}\alpha_t} = 2\sqrt{\epsilon_t(1 - \epsilon_t)} \end{aligned}$$

Bound on training error:

Theorem: Error of h_t over D_t : $\epsilon_t = \frac{1}{2} - \gamma_t$

$$E_{train}(H_T) \leq e^{-2 \sum_{t=1}^T \gamma_t^2}$$

Thus, if $\forall t, \gamma_t \geq \gamma > 0$ then $E_{train}(H_T) \leq e^{-2\gamma^2 T}$


- ▶ Training error decreases exponentially in T
- ▶ To reach $E_{train}(H_T) \leq \epsilon$ we need $T = O\left(\frac{1}{\gamma^2} \log\left(\frac{1}{\epsilon}\right)\right)$

Bound on Training error

► **Step 1:** $D_{T+1}(i) = \frac{1}{N} \left(\frac{e^{-y^{(i)} H_T(x^{(i)})}}{\prod_{t=1}^T Z_t} \right)$

► **Step 2:** $E_{train}(H_T) \leq \prod_{t=1}^T Z_t$

► **Step 3:** $\prod_{t=1}^T Z_t = \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)} = \prod_{t=1}^T \sqrt{1-4\gamma_t^2}$
 $\leq e^{-2 \sum_{t=1}^T \gamma_t^2}$


 $1 - x \leq e^{-x}$
For $0 \leq x \leq 1$

Bound on Training error

► Step 1: $D_{T+1}(i) = \frac{1}{N} \left(\frac{e^{-y^{(i)} H_T(x^{(i)})}}{\prod_{t=1}^T Z_t} \right)$ $H_t(x) = \frac{1}{2} [\alpha_1 h_1(x) + \dots + \alpha_t h_t(x)]$

► Proof:

$$\begin{aligned}
 D_1(i) &= \frac{1}{N} \\
 D_{t+1}(i) &= \frac{D_t(i)}{Z_t} e^{-\frac{1}{2} \alpha_t y^{(i)} h_t(x^{(i)})} \\
 D_{T+1}(i) &= \frac{e^{-\frac{1}{2} \alpha_T y^{(i)} h_T(x^{(i)})}}{Z_T} D_T(i) \\
 &= \frac{e^{-\frac{1}{2} \alpha_T y^{(i)} h_T(x^{(i)})}}{Z_T} \times \frac{e^{-\frac{1}{2} \alpha_{T-1} y^{(i)} h_{T-1}(x^{(i)})}}{Z_{T-1}} D_{T-1}(i) \\
 &= \frac{e^{-\frac{1}{2} \alpha_T y^{(i)} h_T(x^{(i)})}}{Z_T} \times \dots \times \frac{e^{-\frac{1}{2} \alpha_1 y^{(i)} h_1(x^{(i)})}}{Z_1} D_1(i) \\
 &= \frac{1}{N} \left(\frac{e^{-\frac{1}{2} (\alpha_1 y^{(i)} h_1(x^{(i)}) + \dots + \alpha_T y^{(i)} h_T(x^{(i)}))}}{\prod_{t=1}^T Z_t} \right) = \frac{1}{N} \left(\frac{e^{-y^{(i)} H_T(x^{(i)})}}{\prod_{t=1}^T Z_t} \right)
 \end{aligned}$$

Bound on Training error

► **Step 1:** $D_{T+1}(i) = \frac{1}{N} \left(\frac{e^{-y^{(i)} H_T(x^{(i)})}}{\prod_{t=1}^T Z_t} \right)$

► **Step 2:** $E_{train}(H_T) \leq \prod_{t=1}^T Z_t$

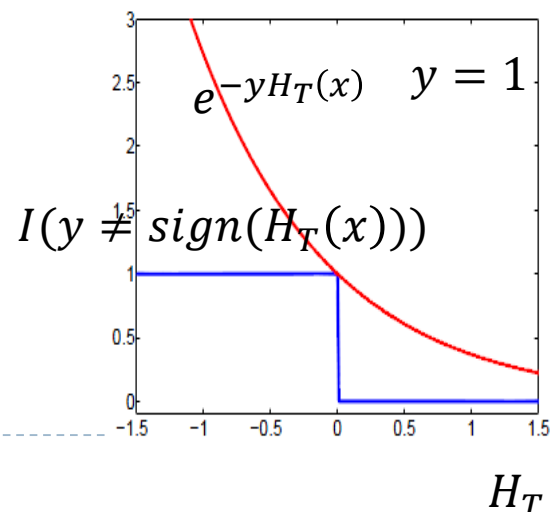
Proof:

$$E_{train}(H_T) = \frac{1}{N} \sum_i I(y^{(i)} \neq \text{sign}(H_T(x^{(i)})))$$

$$= \frac{1}{N} \sum_i I(y^{(i)} H_T(x^{(i)}) \leq 0)$$

$$\leq \frac{1}{N} \sum_i e^{-y^{(i)} H_T(x^{(i)})}$$

$$= \sum_i D_{T+1}(x^{(i)}) \prod_{t=1}^T Z_t$$
$$= \prod_{t=1}^T Z_t$$



Bound on Training error

- ▶ **Step 1:** $D_{T+1}(i) = \frac{1}{N} \left(\frac{e^{-y^{(i)} H_T(x^{(i)})}}{\prod_{t=1}^T Z_t} \right)$
- ▶ **Step 2:** $E_{train}(H_T) \leq \prod_{t=1}^T Z_t$
- ▶ **Step 3:** $\prod_{t=1}^T Z_t = \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)} = \prod_{t=1}^T \sqrt{1-4\gamma_t^2}$
 $\leq e^{-2 \sum_{t=1}^T \gamma_t^2}$

Error of h_t over D_t
 $\epsilon_t = \frac{1}{2} - \gamma_t$

$$\text{Recall: } Z_t = (1 - \epsilon_t)e^{-\frac{1}{2}\alpha_t} + \epsilon_t e^{\frac{1}{2}\alpha_t} = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$$

Bound on training error:

Theorem: Error of h_t over D_t : $\epsilon_t = \frac{1}{2} - \gamma_t$

$$E_{train}(H_T) \leq e^{-2 \sum_{t=1}^T \gamma_t^2}$$

Thus, if $\forall t, \gamma_t \geq \gamma > 0$ then $E_{train}(H_T) \leq e^{-2\gamma^2 T}$

- ▶ Training error decreases exponentially in T
- ▶ To reach $E_{train}(H_T) \leq \epsilon$ we need $T = O\left(\frac{1}{\gamma^2} \log\left(\frac{1}{\epsilon}\right)\right)$

Generalization Error Bounds

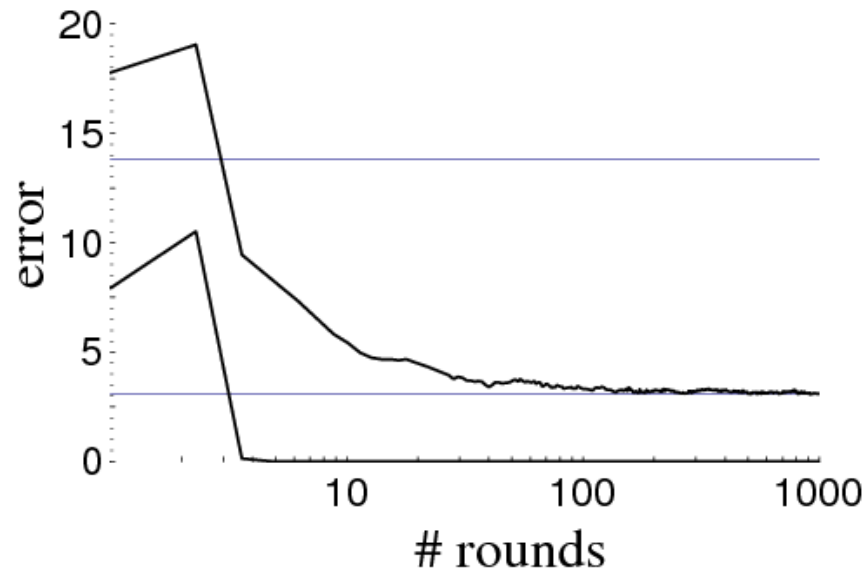
- ▶ $E_{true}(H) \leq E_{train}(H) + \tilde{O}\left(\sqrt{\frac{Td}{N}}\right)$ With high probability
[Freund & Schapire'95]
T: number of boosting rounds
d: VC dimension of weak learner, measures complexity of classifier
N: number of training examples
- ▶ According to this bound, boosting can overfit if T is large

Boosting and overfitting

- ▶ Boosting is **often robust to overfitting**
 - ▶ But not always
 - ▶ may easily overfit in the presence of labeling noise or overlap of classes
- ▶ Test set error decreases even after training error is zero

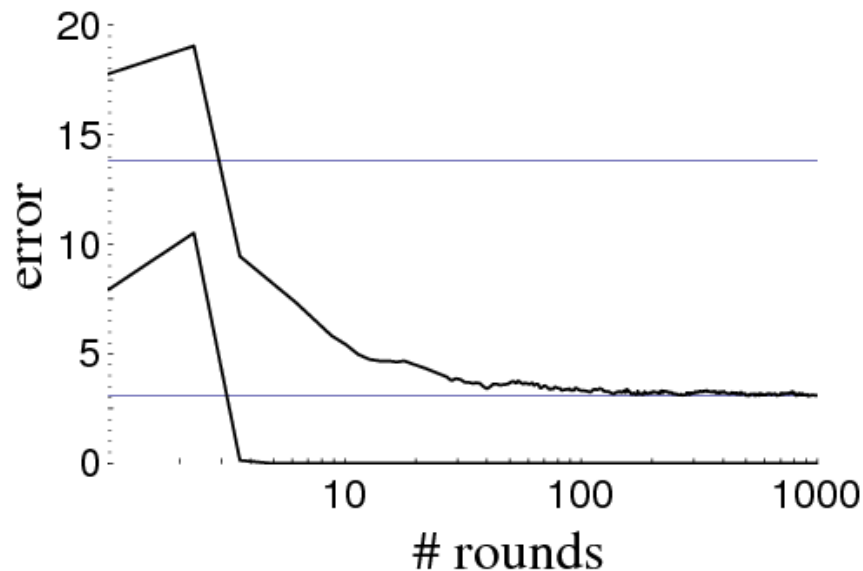
Training and test error

- ▶ Test error usually does not increase as the number of base classifiers becomes very large.



AdaBoost: test error

- ▶ Continuing to add new weak learners after achieving zero training error could even decrease test error!



Generalization Error Bounds

▶ $E_{true}(H) \leq E_{train}(H) + \tilde{O}\left(\sqrt{\frac{Td}{N}}\right)$ With high probability

[Freund & Schapire'95]

T: number of boosting rounds

d: VC dimension of weak learner, measures complexity of classifier

N: number of training examples

- ▶ Is not consistent with experimental results
- ▶ The bound is too loose
- ▶ Margin-based bounds as better analysis

AdaBoost and margin

- ▶ Combined classifier in a more useful form:

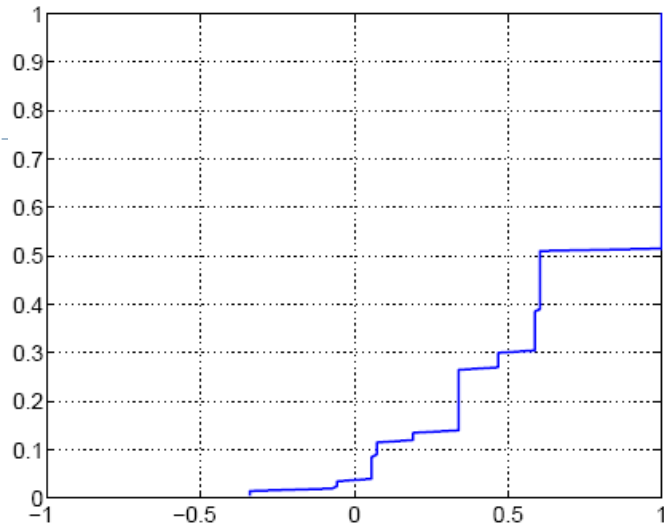
$$H_t(\mathbf{x}) = \frac{\alpha_1 h_1(\mathbf{x}) + \dots + \alpha_t h_t(\mathbf{x})}{\alpha_1 + \dots + \alpha_t}$$

- ▶ This allows us to define a margin:

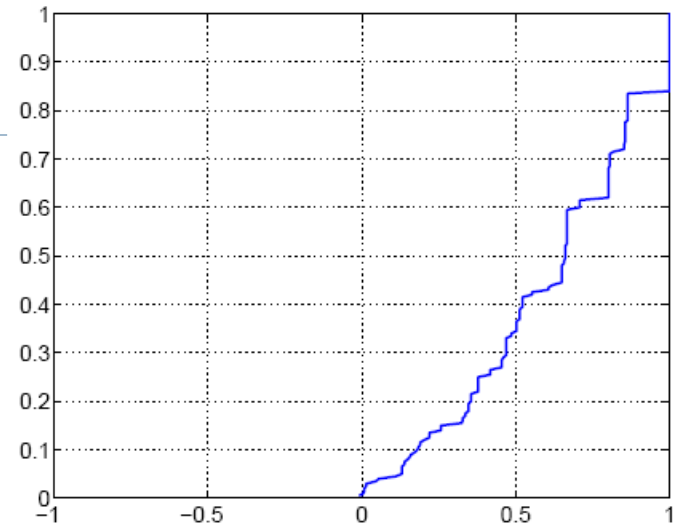
$$\text{margin}(\mathbf{x}_i) = y^{(i)} H_t(\mathbf{x}^{(i)})$$

- ▶ margin lies in $[-1, 1]$ and is negative for misclassified examples.
 - ▶ a measure of confidence in the correct decision
- ▶ Margin of training examples is increased during iterations
 - ▶ Even for correct classification can further improve confidence.

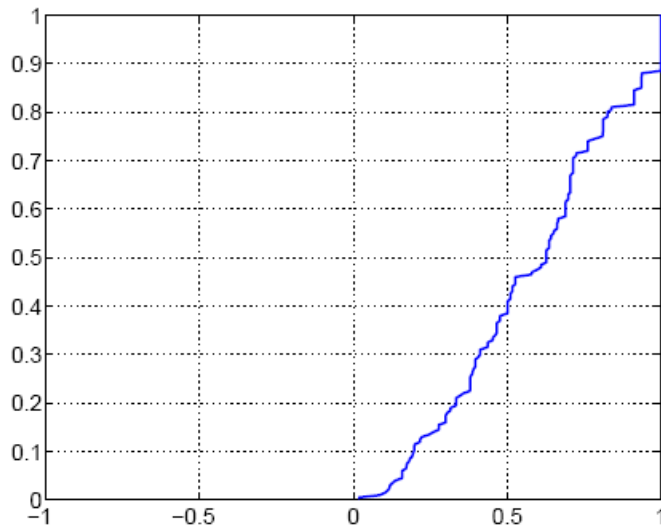
Cumulative distributions of margin values



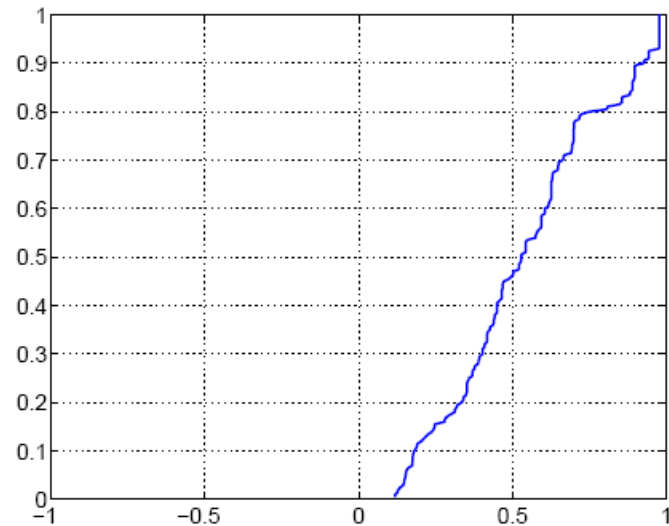
4 iterations



10 iterations



20 iterations



50 iterations

Adaboost and margin

- ▶ When a combined classifier is used, the more classifier agreeing, the more confident you are in your prediction.
- ▶ Successive boosting iterations can improve the majority vote or margin for the training examples

A Margin Bound

- ▶ For any γ , the generalization error is less than:

$$\begin{aligned} P(\text{margin}_h(\mathbf{x}, y) \leq 0) \\ \leq P_{\text{train}}(\text{margin}_h(\mathbf{x}, y) \leq \gamma) + O\left(\sqrt{\frac{d}{N\gamma^2}}\right) \end{aligned}$$

- ▶ It does not depend on T .

$$\begin{aligned} \text{margin}_h(\mathbf{x}, y) &= yH_T(\mathbf{x}) \\ H_t(\mathbf{x}) &= \frac{\alpha_1 h_1(\mathbf{x}) + \dots + \alpha_t h_t(\mathbf{x})}{\alpha_1 + \dots + \alpha_t} \end{aligned}$$

Bagging & Boosting: Summary

▶ Bagging

- ▶ Uses bootstrap sampling to construct several training sets from the original training set and then aggregate the learners trained on these datasets
- ▶ Bagging reduces the variance of high variance learners (e.g. decision tree)

▶ Boosting

- ▶ Combine many (“weak”) classifiers in sequence to find a single “strong” classifier
 - ▶ In each iteration, change the distribution of data to emphasis the samples that have been misclassified by the previous learner

Resources

- ▶ C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 14.2-14.3.
- ▶ Robert E. Schapire, The Boosting Approach to Machine Learning, 2001.
- ▶ Robert E. Schapire et. al, Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26(5):1651-1686, 1998.