

ONLINE MASTERS IN DATA SCIENCE

DSC 255 - MACHINE LEARNING FUNDAMENTALS

THE ADABOOST ALGORITHM

SANJOY DASGUPTA, PROFESSOR

UC San Diego

COMPUTER SCIENCE & ENGINEERING
HALICIOĞLU DATA SCIENCE INSTITUTE

Boosting Weak Learners

A **weak classifier** just has to be marginally better than random guessing:

$$\Pr(h(X) \neq Y) \leq \frac{1}{2} - \epsilon$$

A learning algorithm that can consistently generate such classifiers is called a **weak learner**.

Boosting Weak Learners

A **weak classifier** just has to be marginally better than random guessing:

$$\Pr(h(X) \neq Y) \leq \frac{1}{2} - \epsilon$$

A learning algorithm that can consistently generate such classifiers is called a **weak learner**.

Given: data set $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$.

- Initially give all points equal weight.
- Repeat for $t = 1, 2, \dots$:
 - Feed weighted data set to the weak learner, get back a weak classifier h_t
 - Reweight data to put more emphasis on points that h_t gets wrong
- Combine all these h_t 's linearly

AdaBoost

Given: data set $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$, labels $y^{(i)} \in \{-1, +1\}$.

1 Initialize $D_1(i) = 1/n$ for all $i = 1, 2, \dots, n$

2 For $t = 1, 2, \dots, T$:

- Give D_t to weak learner, get back some $h_t: \mathcal{X} \rightarrow [-1, 1]$
- Compute h_t 's margin of correctness:

$$r_t = \sum_{i=1}^n D_t(i) y^{(i)} h_t(x^{(i)}) \in [-1, 1]$$

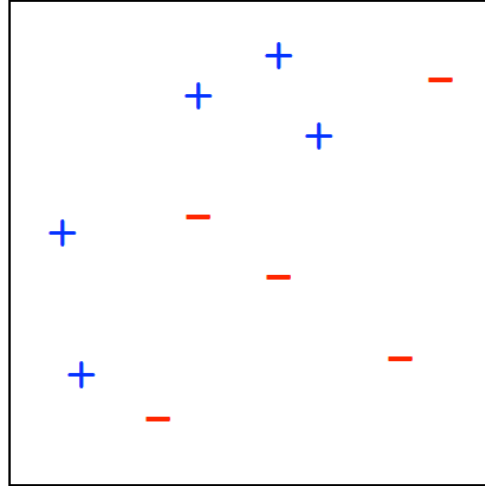
$$\alpha_t = \frac{1}{2} \ln \frac{1+r_t}{1-r_t}$$

- Update weights: $D_{t+1}(i) \propto D_t(i) \exp(-\alpha_t y^{(i)} h_t(x^{(i)}))$

3 Final classifier: $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$

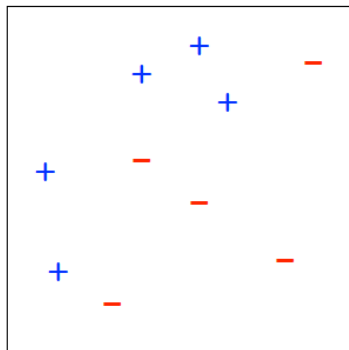
Example (Freund-Schapire)

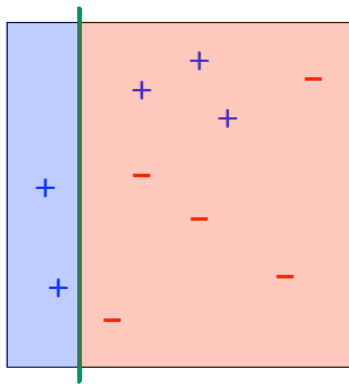
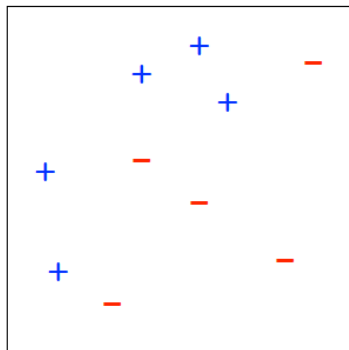
Training set:



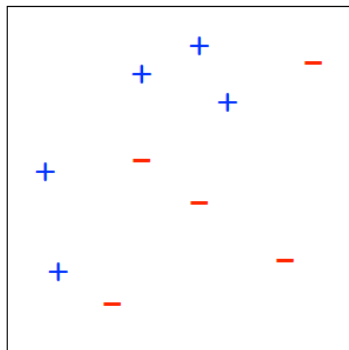
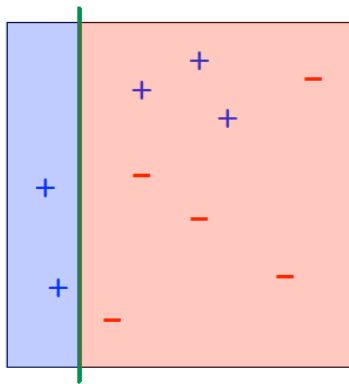
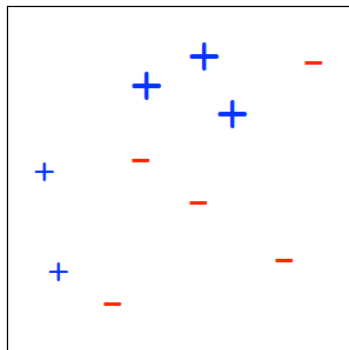
Use "decision stumps" (single-feature thresholds) as weak classifiers

D_1

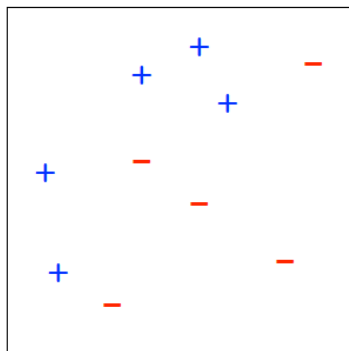
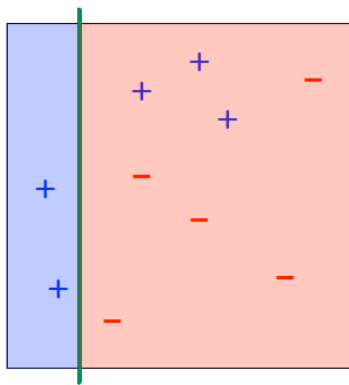
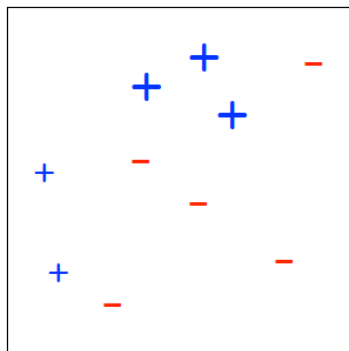


D_1  h_1

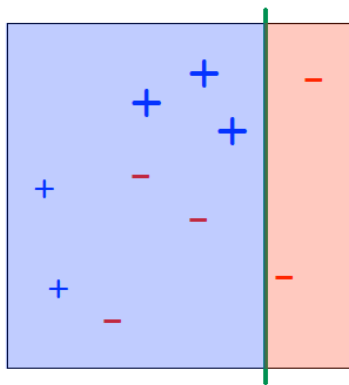
$$r_1 = 0.40, \alpha_1 = 0.42$$

D_1  D_2  h_1

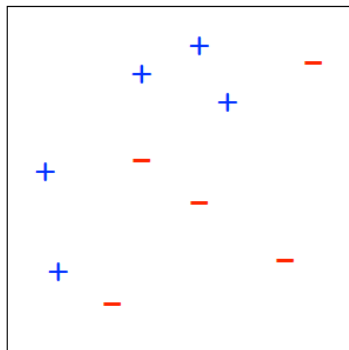
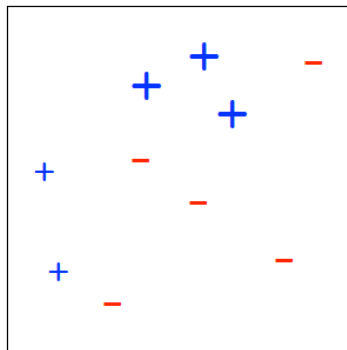
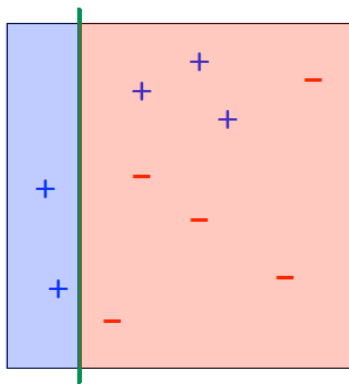
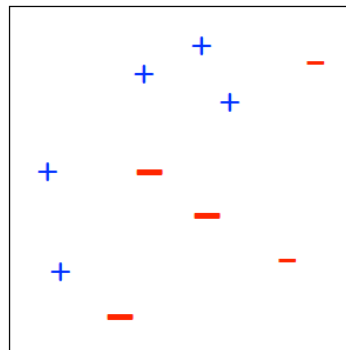
$$r_1 = 0.40, \alpha_1 = 0.42$$

D_1  D_2  h_1

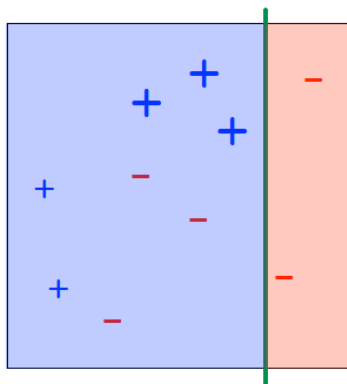
$$r_1 = 0.40, \alpha_1 = 0.42$$

 h_2

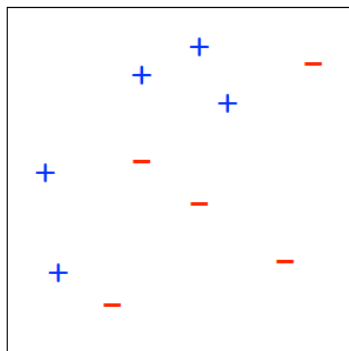
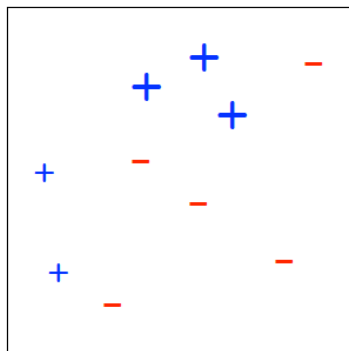
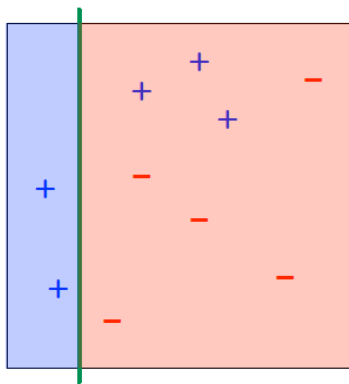
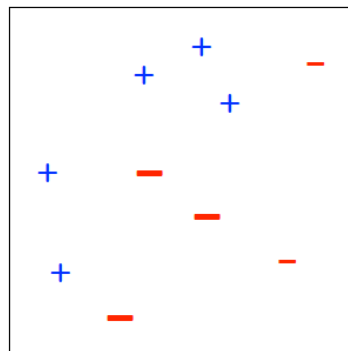
$$r_2 = 0.58, \alpha_2 = 0.65$$

D_1  D_2  D_3  h_1

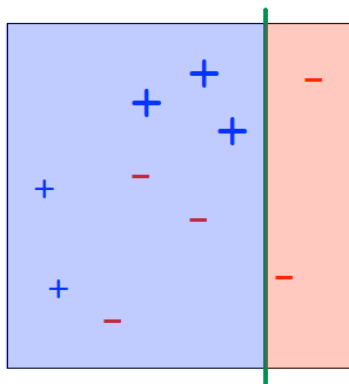
$$r_1 = 0.40, \alpha_1 = 0.42$$

 h_2

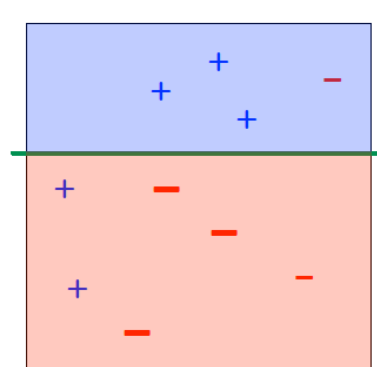
$$r_2 = 0.58, \alpha_2 = 0.65$$

D_1  D_2  D_3  h_1

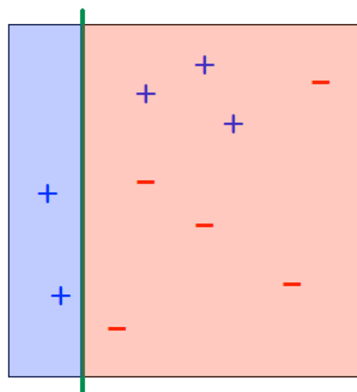
$$r_1 = 0.40, \alpha_1 = 0.42$$

 h_2

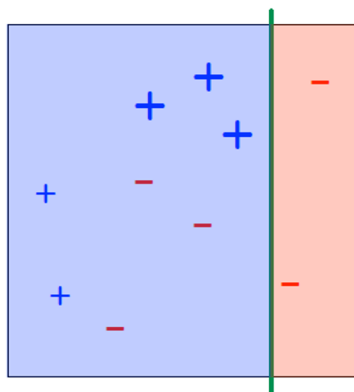
$$r_2 = 0.58, \alpha_2 = 0.65$$

 h_3

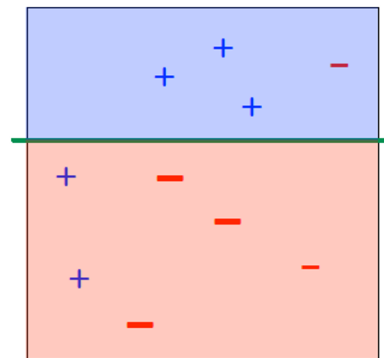
$$r_3 = 0.72, \alpha_3 = 0.92$$



h_1
 $\alpha_1 = 0.42$



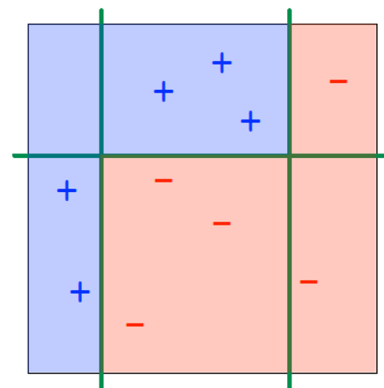
h_2
 $\alpha_2 = 0.65$



h_3
 $\alpha_3 = 0.92$

Final classifier:

$$\text{sign}(0.42h_1(x) + 0.65h_2(x) + 0.92h_3(x))$$



The Surprising Power of Weak Learning

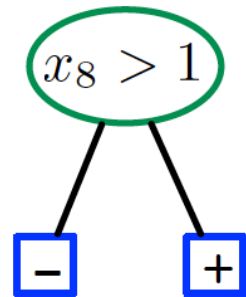
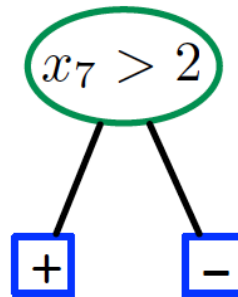
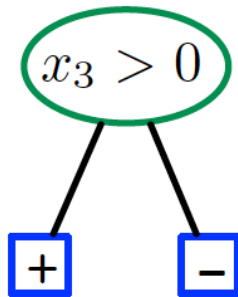
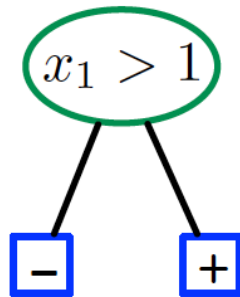
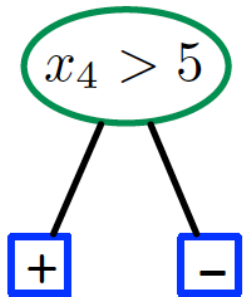
Suppose that on each round t , the weak learner returns a rule h_t whose error on the time- t weighted data distribution is $\leq 1/2 - \gamma$.

Then, after T rounds, the training error of the combined rule

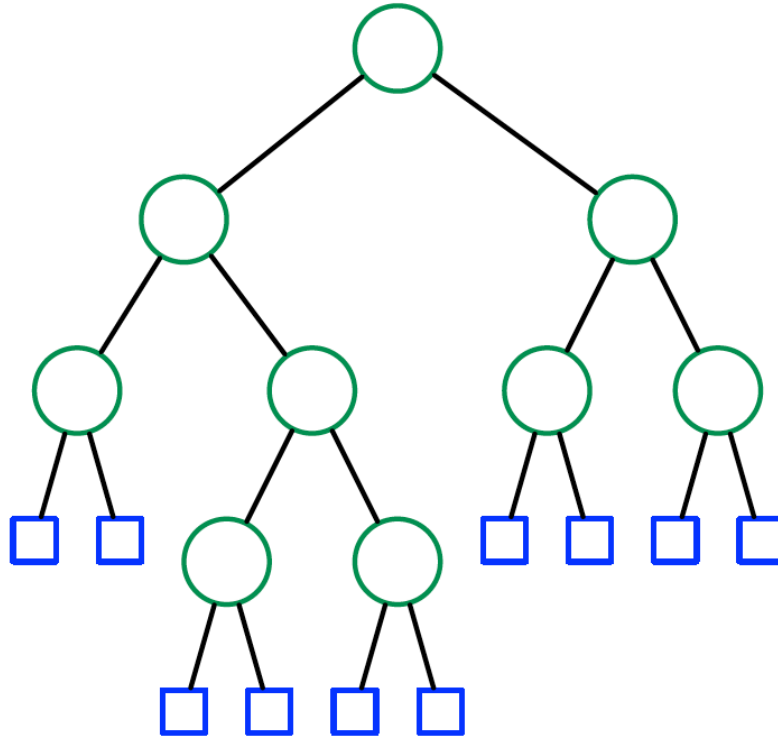
$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

Is at most $e^{-\gamma^2 T/2}$.

Boosting Decision Stumps



Boosting Decision Trees



Boosting Decision Trees

