

Boosting

Some slides are due to Robin Dhamankar
Vandi Verma & Sebastian Thrun

Boosting: motivation

- It is usually hard to design an accurate classifier which generalizes well
- However it is usually easy to find many “rule of thumb” *weak* classifiers
 - A classifier is *weak* if it is only slightly better than random guessing
- Can we combine several weak classifiers to produce an accurate classifier?
 - Question people have been working on since 1980’s

Ada Boost

- Let's assume we have 2-class classification problem, with $y_i \in \{-1, 1\}$
- Ada boost will produce a discriminant function:

$$g(x) = \sum_{t=1}^T \alpha_t f_t(x), \quad \alpha_t \geq 0$$

where $f_t(x)$ is the “weak” classifier

- The final classifier is sign of $g(x)$
- Given x , each weak classifier votes for a label $f_t(x)$ using α_t votes allocated to it. The ensemble then classifies the example according to which label receives the most votes.
- Note that $g(x) \in [-1, 1]$ whenever the votes are normalized to sum to one. So, $g(x) = 1$ only if all the weak classifiers agree that the label should be $y = 1$.

Idea Behind Ada Boost

- Algorithm is iterative
- Maintains distribution of weights over the training examples
- Initially distribution of weights is uniform
- At successive iterations, the weight of misclassified examples is increased, forcing the weak learner to focus on the hard examples in the training set

More Comments on Ada Boost

- Ada boost is very simple to implement, provided you have an implementation of a “weak learner”
- Will work as long as the “basic” classifier $f_t(x)$ is at least slightly better than random
- Can be applied to boost any classifier, not necessarily weak

Ada Boost *(slightly modified from the original version)*

- $d(x)$ is the distribution of weights over the N training points $\sum d(x_i)=1$
- Initially assign uniform weights $d_0(x_i) = 1/N$ for all x_i
- At each iteration t :
 - Find best weak classifier $f_t(x)$ using weights $d_t(x)$
 - Compute the error rate ϵ_t as
$$\epsilon_t = \sum_{i=1 \dots N} d_t(x_i) \cdot I[y_i \neq f_t(x_i)]$$
 - assign weight α_t the classifier f_t 's in the final hypothesis
$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$
 - For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$
 - Normalize $d_{t+1}(x_i)$ so that $\sum_{i=1} d_{t+1}(x_i) = 1$
- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$

Ada Boost

- At each iteration t :
 - Find best weak classifier $f_t(x)$ using weights $d_t(x)$
 - Compute ε_t the error rate as
$$\varepsilon_t = \sum d_t(x_i) \cdot I[y_i \neq f_t(x_i)]$$
 - assign weight α_t the classifier f_t 's in the final hypothesis
$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$
 - For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$
 - Normalize $d_{t+1}(x_i)$ so that $\sum_{t+1} d(x_i) = 1$
- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$
- If the classifier does not take weighted samples, this step can be achieved by sampling from the training samples according to the distribution $d_t(x)$

Ada Boost

- At each iteration t :

- Find best weak classifier $f_t(x)$ using weights $d_t(x)$

- Compute ε_t the error rate as

$$\varepsilon_t = \sum d_t(x_i) \cdot \mathbb{I}[y_i \neq f_t(x_i)]$$

- assign weight α_t the classifier f_t 's in the final hypothesis

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

- For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$

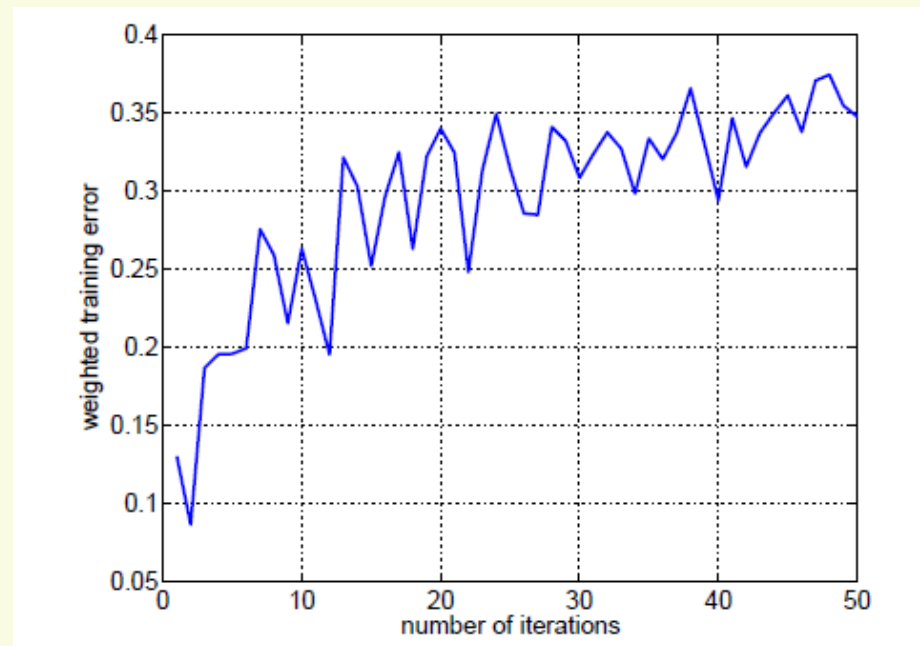
- Normalize $d_{t+1}(x_i)$ so that $\sum d_{t+1}(x_i) = 1$

- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$

- Since the weak classifier is better than random, we expect $\varepsilon_t < 1/2$

Weighted error (ε_t)

- The weighted error achieved by a new simple classifier $f_t(x)$ relative to weights $d_t(x)$ tends to increase with t , i.e., with each boosting iteration (though not monotonically).
- The reason for this is that since the weights concentrate on examples that are difficult to classify correctly, subsequent base learners face harder classification tasks.



Weighted error (ε_t)

- It can be shown that the weighted error of the simple classifier $f_t(x)$ relative to updated weights $d_{t+1}(x)$ is exactly 0.5.
- This means that the simple classifier introduced at the t -th boosting iteration will be useless (at chance level) for the next boosting iteration. So the boosting algorithm would never introduce the same simple classifier twice in a row.
- It could, however, reappear later on (relative to a different set of weights)

Ada Boost

- At each iteration t :

- Find best weak classifier $f_t(x)$ using weights $d_t(x)$

- Compute ε_t the error rate as

$$\varepsilon_t = \sum d(x_i) \cdot I(y_i \neq f_t(x_i))$$

- assign weight α_t the classifier f_t 's in the final hypothesis

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

- For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$

- Normalize $d_{t+1}(x_i)$ so that $\sum d_{t+1}(x_i) = 1$

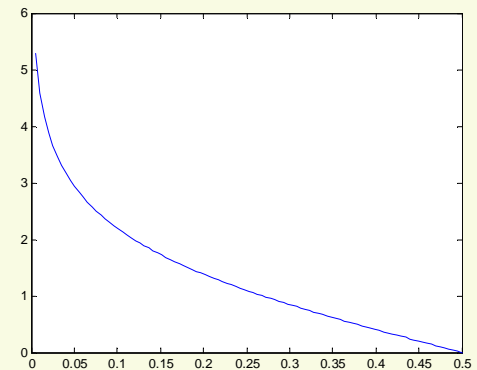
- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$

- Recall that $\varepsilon_t < \frac{1}{2}$

- Thus $(1 - \varepsilon_t) / \varepsilon_t > 1 \Rightarrow \alpha_t > 0$

- The smaller is ε_t , the larger is α_t , and thus the more importance (weight) classifier $f_t(x)$ gets in the final classifier

$$f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$$



Ada Boost

- At each iteration t :
 - Find best weak classifier $f_t(x)$ using weights $d_t(x)$
 - Compute ε_t the error rate as
$$\varepsilon_t = \sum d_t(x_i) \cdot I(y_i \neq f_t(x_i))$$
 - assign weight α_t the classifier f_t 's in the final hypothesis
$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$
 - For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$
 - Normalize $d_{t+1}(x_i)$ so that $\sum d_{t+1}(x_i) = 1$
- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$
- Weight of misclassified examples is increased and the new $d_{t+1}(x_i)$'s are normalized to be a distribution again

Ada Boost

- At each iteration t :
 - Find best weak classifier $f_t(x)$ using weights $d_t(x)$
 - Compute ε_t the error rate as
$$\varepsilon_t = \sum d_t(x_i) \cdot I(y_i \neq f_t(x_i))$$
 - assign weight α_t the classifier f_t 's in the final hypothesis
$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$
 - For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp(-\alpha_t y_i f_t(x_i))$
 - Normalize $d_{t+1}(x_i)$ so that $\sum d_{t+1}(x_i) = 1$
- $f_{FINAL}(x) = \text{sign} \left[\sum \alpha_t f_t(x) \right]$

Ensemble training error

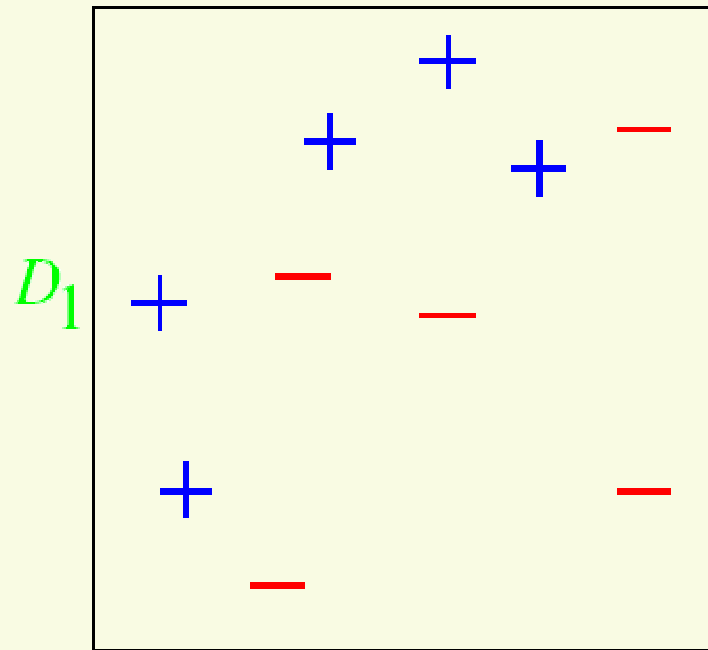
- It can be shown that the training error drops exponentially fast, if each weak classifier is slightly better than random

$$Err_{train} \leq \exp\left(-2\sum_t \gamma_t^2\right)$$

- Here $\gamma_t = \varepsilon_t - 1/2$, where ε_t is classification error at round t (weak classifier f_t)

AdaBoost Example

from "A Tutorial on Boosting" by Yoav Freund and Rob Schapire

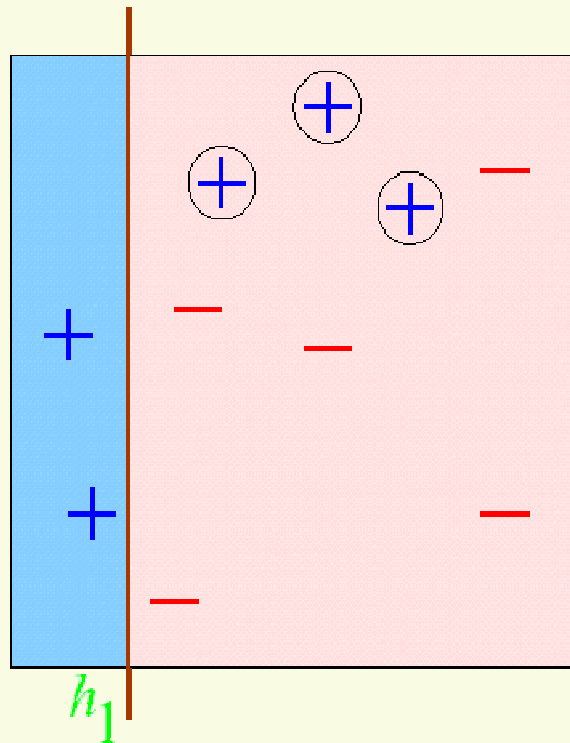


Original Training set : equal weights to all training samples

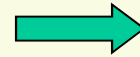
Note: in the following slides, $h_t(x)$ is used instead of $f_t(x)$, and D instead of d

AdaBoost Example

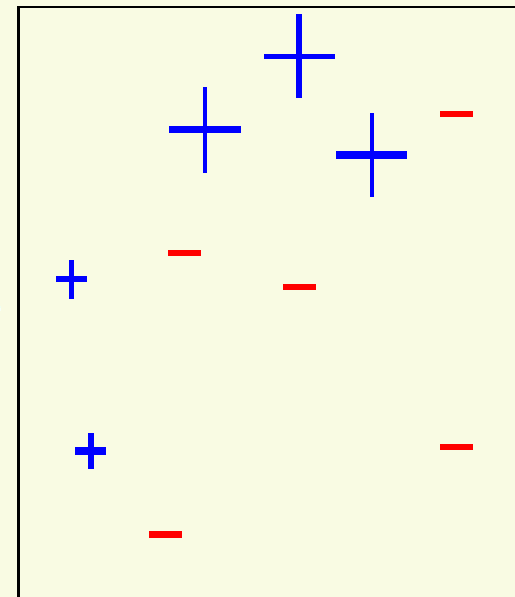
ROUND 1



$$\begin{aligned}\epsilon_1 &= 0.30 \\ \alpha_1 &= 0.42\end{aligned}$$

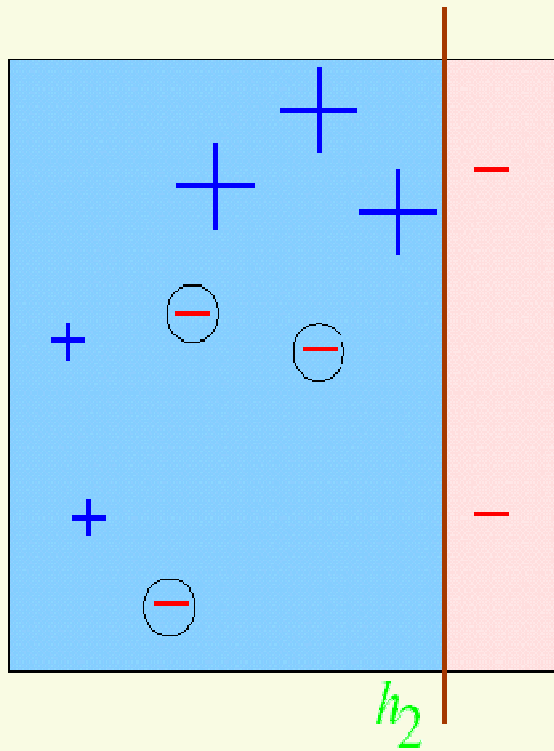


D_2



AdaBoost Example

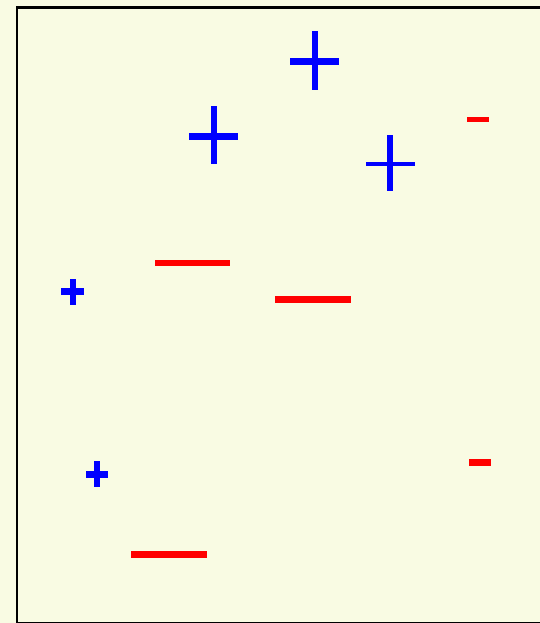
ROUND 2



$\epsilon_2 = 0.21$
 $\alpha_2 = 0.65$

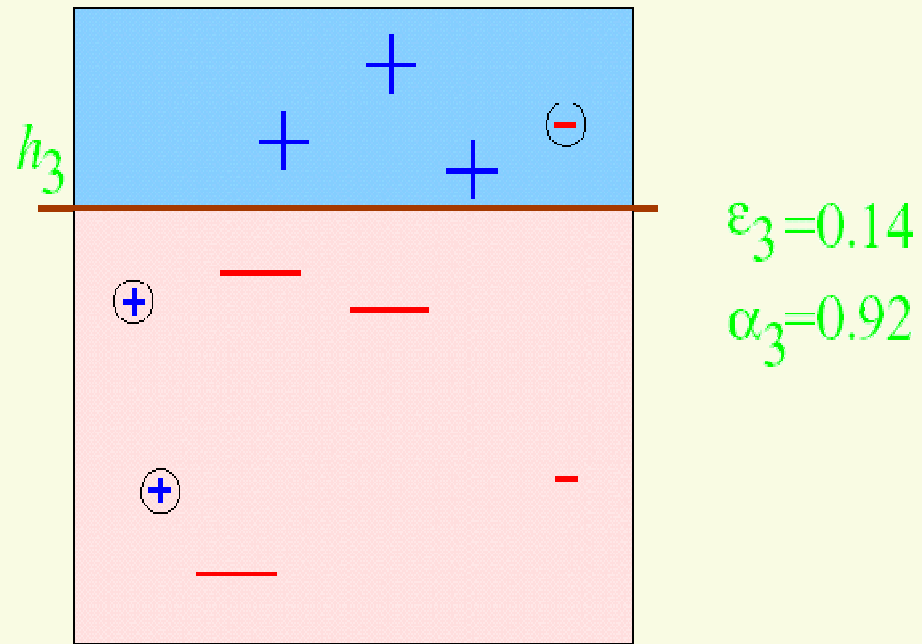


D_3



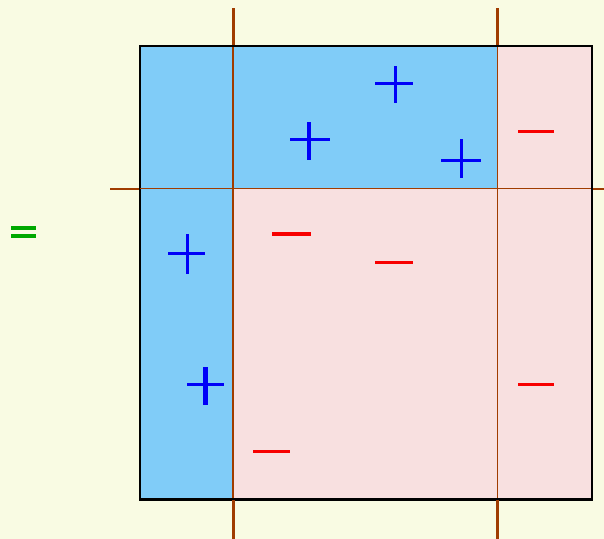
AdaBoost Example

ROUND 3



AdaBoost Example

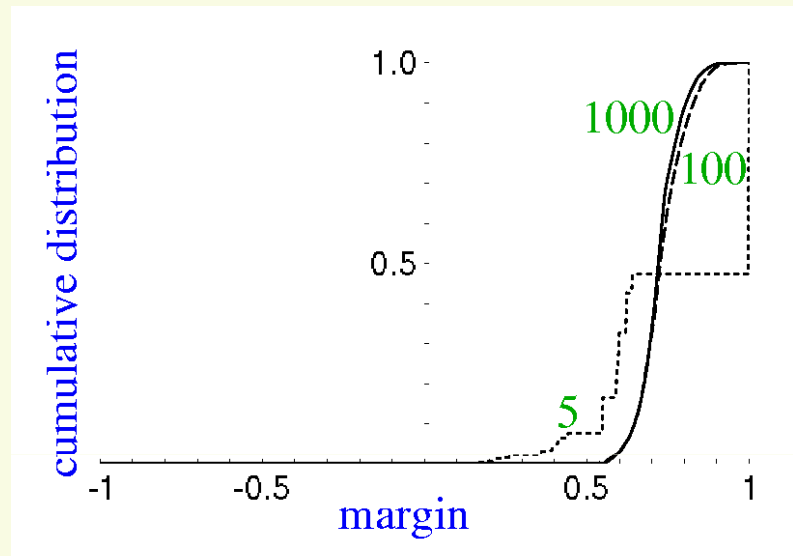
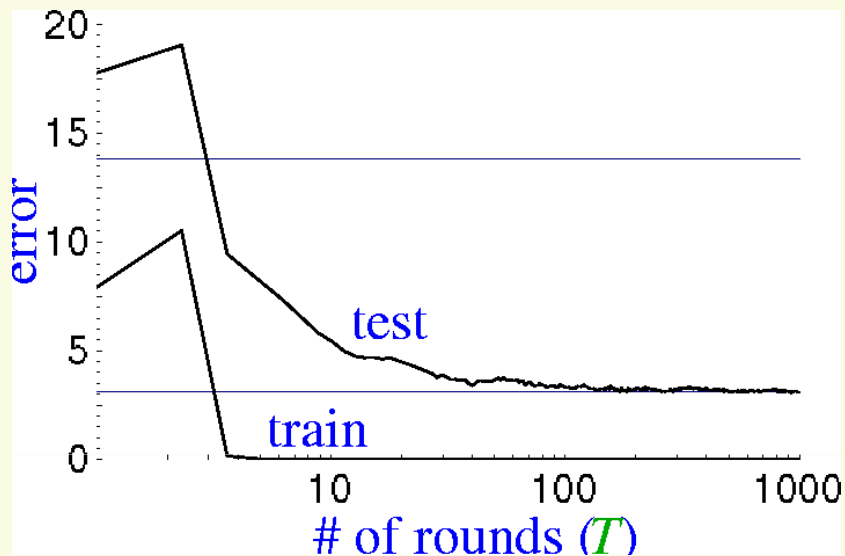
$$f_{\text{FINAL}}(X) = \text{sign} \left(0.42 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right)$$



AdaBoost Comments

- But we are really interested in the generalization properties of $f_{\text{FINAL}}(x)$, not the training error
- AdaBoost was shown to have excellent generalization properties in practice.
- It can be shown that boosting “aggressively” increases the margins of training examples, as iterations proceed
 - margins continue to increase even when training error reaches zero
 - Helps to explain empirically observed phenomena: test error continues to drop even after training error reaches zero

The Margin Distribution



epoch	5	100	1000
training error	0.0	0.0	0.0
test error	8.4	3.3	3.1
%margins ≤ 0.5	7.7	0.0	0.0
Minimum margin	0.14	0.52	0.55

Practical Advantages of AdaBoost

- fast
- simple
- Has only one parameter to tune (T)
- flexible: can be combined with any classifier
- provably effective (assuming weak learner)
 - shift in mind set: goal now is merely to find hypotheses that are better than random guessing
- finds outliers
 - The hardest examples are frequently the “outliers”

Caveats

- performance depends on data & weak learner
- AdaBoost can fail if
 - weak hypothesis too complex (overfitting)
 - weak hypothesis too weak ($\gamma_t \rightarrow 0$ too quickly),
 - underfitting
 - Low margins \rightarrow overfitting
- empirically, AdaBoost seems especially susceptible to noise