

Boosting and margins

Maria-Florina Balcan

10/11/2011

Outline

- AdaBoost
 - Algorithm
 - AdaBoost Behavior in Experiments

- Generalization error as a function of Margin Distributions
 - Classification Margin
 - Finite base-classifier spaces

- The effect of Boosting on Margin Distributions

AdaBoost recap

AdaBoost combines weak learners in a weighted majority voting scheme

- given a training set $(x_1, y_1), \dots, (x_m, y_m)$
- $y_i \in \{-1, 1\}$ correct label of instance $x_i \in X$
- for $t = 1, \dots, T$
 - construct a distribution D_t on $\{1, 2, \dots, m\}$
 - find a weak hypothesis ("rule of thumb") $h_t : X \rightarrow \{-1, 1\}$ with small error ϵ_t on D_t , $\epsilon_t = \Pr_{D_t}[h_t(x_i) \neq y_i]$
- output final hypothesis H_{final}
- constructing D_t :
 - $D_1(i) = \frac{1}{m}$
 - given D_t and h_t
$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \cdot e^{-\alpha_t} \text{ if } y_i = h_t(x_i)$$
$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \cdot e^{\alpha_t} \text{ if } y_i \neq h_t(x_i) \text{ where}$$
$$\alpha_t = \frac{1}{2} \ln \left[\frac{1-\epsilon_t}{\epsilon_t} \right]$$
- final hypothesis: $H_{final}(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$

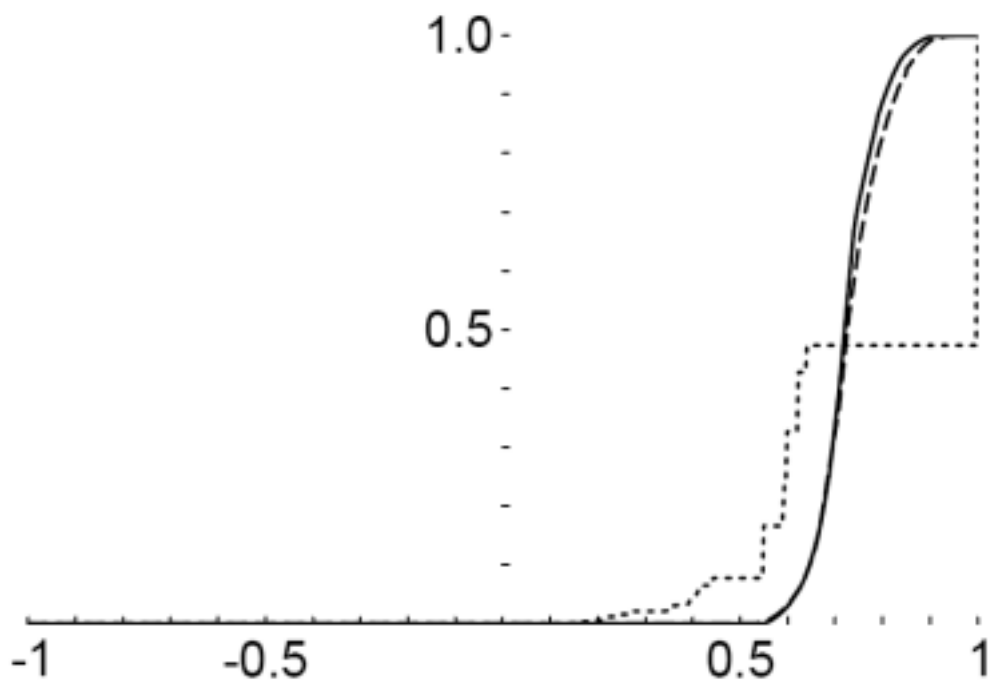
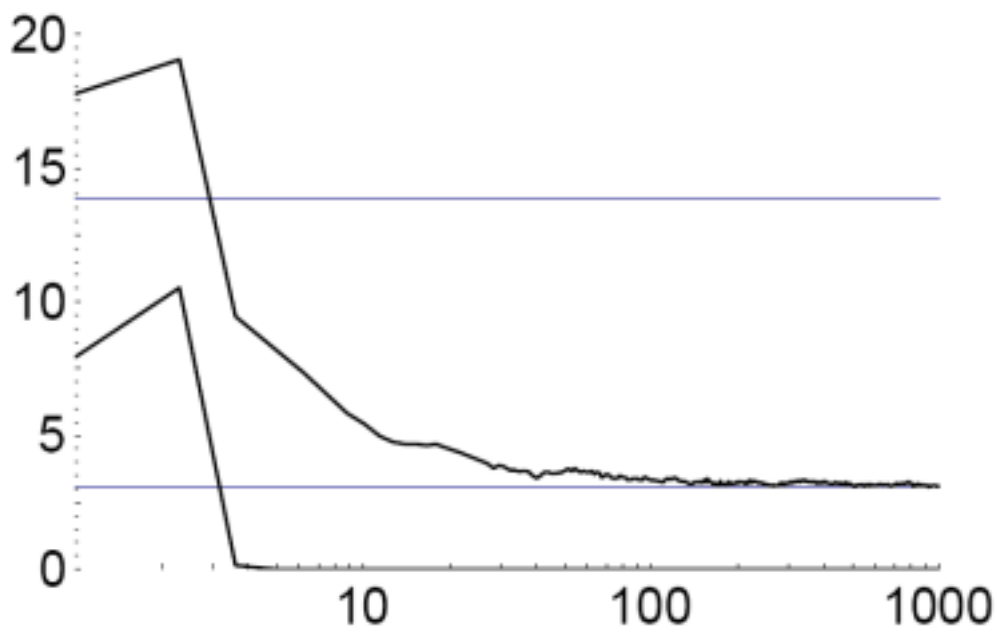
AdaBoost Behavior in Experiments

Experiments with boosting showed that the test error of the generated classifier usually **does not increase** as its size becomes very large.

Experiments with boosting showed also that continuing to add new weak learners after **correct** classification of the training set had been achieved could further **improve** test set performance!!!

These results seem to contradict Occam's razor: in order achieve good test error the classifier should be as **simple** as possible!

Error Curve, Margin Distr. Graph - Plots from [SFBL98]



Analyzing Generalization Error

Remember, usual sample complexity statements:

Theorem 1 *If H is a finite hypotheses space, then with probab. $1 - \delta$, $\forall h \in H$ we have $|err(h) - \widehat{err}(h)| < \epsilon$ given that we see*

$$m \geq O\left(\frac{1}{\epsilon^2} \left[\ln |H| + \ln \frac{1}{\delta}\right]\right)$$

labeled examples.

Or, another way to state it: with probab. $1 - \delta$, $\forall h \in H$

$$err(h) \leq \widehat{err}(h) + O\left(\sqrt{\frac{\ln |H| + \ln\left(\frac{1}{\delta}\right)}{m}}\right)$$

given that we see m labeled examples.

In general, with probab. $1 - \delta$, $\forall h \in H$,

$$err(h) \leq \widehat{err}(h) + O\left(\sqrt{\frac{\ln(C[2m]) + \ln\left(\frac{1}{\delta}\right)}{m}}\right)$$

How can we explain the experiments?

R. Schapire, Y. Freund, P. Bartlett, W. S. Lee.
present in "Boosting the margin: A new explanation for the effectiveness of voting methods"
a nice theoretical explanation.

Main Idea:

Training error does not tell the whole story.

We need also to consider the classification confidence!!

Classification Margin

Consider H to be the space of weak hypotheses. Define the **convex hull** of H to be

$$co(H) = \left\{ f = \sum_{t=1}^T a_t h_t, a_t \geq 0, \sum_{t=1}^T a_t = 1, h_t \in H \right\}$$

Let $f \in co(H)$, $f = \sum_{t=1}^T a_t h_t, a_t \geq 0, \sum_{t=1}^T a_t = 1$.

The majority vote rule H_f associated with f (given by $H_f(x) = \text{sign}(f(x))$) gives a wrong prediction on the example (x, y) iff $yf(x) \leq 0$.

Define the **margin** of H_f (or of f) on example (x, y) to be $yf(x)$.

$$\text{Note that } yf(x) = y \sum_{t=1}^T [a_t h_t(x)] = \sum_{t=1}^T [y a_t h_t(x)] = \sum_{t:y=h_t(x)} a_t - \sum_{t:y \neq h_t(x)} a_t.$$

The margin is positive iff $y = H_f(x)$.

See $|yf(x)| = |f(x)|$ as the strength or the confidence of the vote.

Gen. error as a function of Margin Distributions

Assume that the examples are generated i.i.d. according to some distr. D over $X \times \{-1, 1\}$; denote by $\Pr_D[\cdot]$ the probability when (x, y) is chosen from D .

If S is a training set (a sample of size m , $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$), then we denote by $\Pr_S[\cdot]$ the probability when (x, y) is chosen uniformly at random from S .

Theorem 2 *If H finite, then with probability at least $1 - \delta$, $\forall f \in co(H)$, $\forall \theta > 0$,*

$$\Pr_D [yf(x) \leq 0] \leq \Pr_S [yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{\ln m \ln |H|}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

Theorem 3 *If H has VCdimension d then with probability at least $1 - \delta$, $\forall f \in co(H)$, $\forall \theta > 0$,*

$$\Pr_D [yf(x) \leq 0] \leq \Pr_S [yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{d \ln^2 \frac{m}{d}}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

Note: no dependence on number of weak hypotheses !!!

A First Lemma

- $N > 0$, C_N - the set of unweighted averages over N elements from H , i.e.

$$C_N = \left\{ g \mid g(x) = \frac{1}{N} \sum_{j=1}^N h_j(x), h_j \in H \right\}$$

- **Lemma 4** *With probability at least $1 - \delta_N$ (over the random choice of the training set), $\forall g \in C_N, \forall \theta > 0$,*

$$\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] \leq \Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] + \epsilon_N$$

where

$$\epsilon_N = \sqrt{\frac{1}{2m} \ln \left[\frac{(N+1)|H|^N}{\delta_N} \right]}$$

A First Lemma - Proof

Proof: For θ and g fixed

$$\Pr_{\text{sample}} \left[\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] > \Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] + \epsilon_N \right] \leq \exp \left[-2m\epsilon_N^2 \right]$$

By union bound, the probability (taken over a random choice of S) that $\exists g \in C_N$ such that $\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] > \Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] + \epsilon_N$ is at most $\leq |H|^N \exp \left[-2m\epsilon_N^2 \right]$

Since $yg(x)$ is always a multiple of $\frac{1}{N}$, we finally get that the probability (taken over a random choice of S) that $\exists \theta > 0, \exists g \in C_N$ such that $\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] > \Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] + \epsilon_N$ is at most $\leq (N + 1)|H|^N \exp \left[-2m\epsilon_N^2 \right]$.

We finally set $(N + 1)|H|^N \exp \left[-2m\epsilon_N^2 \right] = \delta_N$ and get the desired result. ■

A Second Lemma

Lemma 5 *With probability at least $1 - \delta$ (over the random choice of the training set), $\forall \theta > 0$, $\forall N > 0$, $\forall g \in C_N$,*

$$\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] \leq \Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] + \epsilon_N,$$

where

$$\epsilon_N = \sqrt{\frac{1}{2m} \ln \left[\frac{N(N+1)^2 |H|^N}{\delta} \right]}.$$

Proof: Just use lemma 1 and plug in $\delta_N = \frac{\delta}{N(N+1)}$. ■

Main Result

If H finite, then with probability at least $1 - \delta$, $\forall f \in \text{co}(H)$, $\forall \theta > 0$, we get

$$\Pr_D [yf(x) \leq 0] \leq \Pr_S [yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \sqrt{\frac{\ln m \ln |H|}{\theta^2} + \ln \frac{1}{\delta}}\right)$$

Proof

Consider $f \in \text{co}(H)$, $f = \sum_{t=1}^T a_t h_t$; then f can be associated with a distr. D_f over H as defined by the coefficients a_t .

Moreover we can map f to a distribution Q_f over C_N ; a function $g \in C_N$ distributed according to Q_f is generated by choosing g_1, \dots, g_N ind. at random according to D_f and then defin-

ing $g(x) = \frac{1}{N} \sum_{j=1}^N g_j(x)$.

Main Result, Proof

Note: If we fix x then $\mathbf{E}_{D_f}[g_j(x)] = \sum_{t=1}^T a_t h_t(x) = f(x)$ and $\mathbf{E}_{g \sim Q_f}[g(x)] = f(x)$.

Therefore

$$\Pr_{g \sim Q_f} \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right] \leq \exp \left[-N\theta^2/8 \right]$$

and so

$$\mathbf{E}_D \left[\Pr_{g \sim Q_f} \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right] \right] \leq \exp \left[-N\theta^2/8 \right]$$

or

$$\Pr_{D, g \sim Q_f} \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right] \leq \exp \left[-N\theta^2/8 \right].$$

Similarly

$$\Pr_{S, g \sim Q_f} \left[yg(x) \leq \frac{\theta}{2}, yf(x) > \theta \right] \leq \exp \left[-N\theta^2/8 \right].$$

Main Result, Proof - cont

Consider $f \in co(H)$. For any $g \in C_N$, for any $\theta > 0$ we have:

$$\Pr_D [yf(x) \leq 0] \leq \Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] + \Pr_D \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right]$$

Therefore

$$\mathbf{E}_{g \sim Q_f} [\Pr_D [yf(x) \leq 0]] \leq \mathbf{E}_{g \sim Q_f} \left[\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] \right] + \mathbf{E}_{g \sim Q_f} \left[\Pr_D \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right] \right]$$

and so

$$\Pr_D [yf(x) \leq 0] \leq \mathbf{E}_{g \sim Q_f} \left[\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] \right] + \mathbf{E}_D \left[\Pr_{g \sim Q_f} \left[yg(x) > \frac{\theta}{2}, yf(x) \leq 0 \right] \right]$$

and therefore

$$\Pr_D [yf(x) \leq 0] \leq \mathbf{E}_{g \sim Q_f} \left[\Pr_D \left[yg(x) \leq \frac{\theta}{2} \right] \right] + \exp \left[-N\theta^2/8 \right]$$

Main Result, Proof - finish

Therefore, by lemma 5 we know that with probability $1 - \delta$ (over the random choice of the training set) we have

$$\begin{aligned} \Pr_D [yf(x) \leq 0] &\leq \exp \left[-N\theta^2/8 \right] + \\ &\mathbf{E}_{g \sim Q_f} \left[\Pr_S \left[yg(x) \leq \frac{\theta}{2} \right] \right] + \\ &\sqrt{\frac{1}{2m} \ln \left[\frac{N(N+1)^2 |H|^N}{\delta} \right]} \end{aligned}$$

and so

$$\begin{aligned} \Pr_D [yf(x) \leq 0] &\leq 2 \exp \left[-N\theta^2/8 \right] + \\ \Pr_S [yf(x) \leq \theta] &+ \sqrt{\frac{1}{2m} \ln \left[\frac{N(N+1)^2 |H|^N}{\delta} \right]} \end{aligned}$$

Choosing $N = \frac{4}{\theta^2} \ln \left[\frac{m}{\ln |H|} \right]$ we get the desired result.

Boosting increases the margin

Theorem 6 *Suppose the base learning algorithm, when called by AdaBoost, generates classifiers with weighted errors $\epsilon_1, \dots, \epsilon_T$. Then for any θ we have*

$$\Pr_S [yf(x) \leq \theta] \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{(1-\theta)} (1 - \epsilon_t)^{(1+\theta)}}$$

Interpretation: if $\forall t, \epsilon_t < \frac{1}{2} - \gamma$ and if $\theta < \gamma$, then $\Pr_S [yf(x) \leq \theta]$ goes to 0 as $T \rightarrow \infty$.

(If θ is not too large, then the fraction of the training examples for which $yf(x) \leq \theta$ decreases exponentially to 0 exponentially fast with the number of base classifiers.)

References

[SFBL98] R. Schapire, Y. Freund, P. Bartlett and W. S. Lee. Boosting the margin: A new explanation for the effectiveness of voting methods, *The Annals of Statistics*, 26(5), 1651–1686, 1998

[DGL] L. Devroye, L. Györfi and G. Lugosi. Probabilistic Theory of Pattern Recognition, Springer-Verlag, 1996

[SBWA98] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson and M. Anthony. Structural Risk Minimization over Data-Dependent Hierarchies, *IEEE Transactions on Information Theory*, 44(5), 1926-1940, 1998.