Something Unexpected: AdaBoost tends not to overfit

The Bias / Variance Curve we expect:
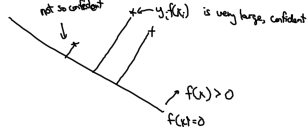


test error

underfit → just right → overfit

complexity →
(approx by # boosting rounds)

Luckily, experimentalists ignored this advice!

Drucker & Cortes '96, Quinlan '96, Breiman '98 found curves like this instead

test error

← no overfitting

# boosting rounds →

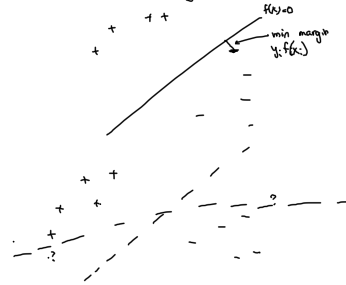Seems to contradict Occam's Razor!

Explanation is in "Boosting the Margin" (Schapire, Freund, Bartlett, Lee)

complexity ≠ # of boosting rounds
· complexity is related to the __margin__.

The __margin__ of training example $i$ is proportional to $y_i f(x_i)$.
 - confidence measure of a classifier's ability
 - "distance" from the training example to the decision bdry.

not so confident ←    $y_i f(x_i)$ is very large, confident

$f(x) > 0$
$f(x) = 0$

The __minimum margin__ (or the margin of the classifier $f$) is the min of the margins over training examples ie. the "distance" from the decision boundary to the nearest training example.

$f(x) = 0$

← min margin $y_i f(x_i)$

"Boosting the Margin" has 2 results:

1) Large Margins are good (statistical guarantee)

P Seudo theorem: w. high prob,

$$\text{misclassific. error on the test set} \le \text{misclassific. error on training set} + fctn\left( (VCdim)^{\frac{1}{2}}, \frac{1}{\sqrt{m}}, \frac{1}{margin} \right)$$

measure of complexity of a set of fctns

This is why it is believed that lg margins are the key to success.
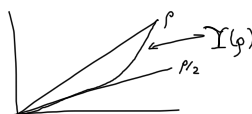
2) AdaBoost achieves large margins

P Seudo Thm:
If the max margin (for our data & weak classifiers) is $\rho$,
AdaBoost achieves a margin of at least $\rho/2$



margin achieved

$\rho$
$\rho/2$
$\rho$

No more contradiction of Occam's Razor. But theory is still incomplete

Experimental results often indicated that AdaBoost maximizes the margin... except for Breiman. Breiman claimed his alg arc-gv achieved a max. margin & that AdaBoost did not (1999).

Rätsch & Warmuth (2005) tightened the bd slightly



$\rho$
$\Upsilon(\varphi)$
$\rho/2$

If max margin is $\rho$, R & W proved AdaBoost achieves a margin of at least
$$\Upsilon(\rho) = \frac{-\ln(1-\rho^2)}{\ln\left(\frac{1+\rho}{1-\rho}\right)}.$$

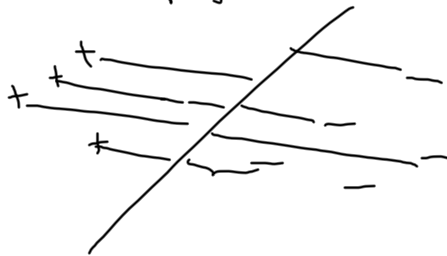It turns out that R & W were right!

Feb 21-2:37 PM

Recently, it has been shown :

    1) AdaBoost does not necessarily maximize the margin
       (RDS 04)

    2) Rätsch & Warmuth's bound is tight, i.e. that
       AdaBoost (in a special case) can acheive a
       margin of exactly $\Upsilon(\rho)$.   (RSD 06,07)

It also turns out that Breiman was right:

    ·) arc-gv does maximize the margin (Breiman, R&W)
       with a fast convergence rate (RSD 07)

    2) AdaBoost still beats arcgv experimentally.
       (Breiman, Reyzin & Schapire 06)

          exploring dist of margins



In fact, there have been a number of other algorithms
(w/ fast convergence rates) designed to maximize the margin,
none of which has been shown to beat AdaBoost experimentally.
    arc-gv (Breiman 99)
    AdaBoost$_\rho$ & AdaBoost* (Rätsch & Warmuth 05)
    Approx. Coord Descent Boosting (RSD 07)
    LP-AdaBoost (Grove & Schuurmans (98))

There are many unsolved problems - this is the current
state of theoretical results on generalization of AdaBoost!

Experimental work - there are *lots* of variations of
    AdaBoost (almost all heuristically based, all
    claim to beat AdaBoost.)

One more important fact about AdaBoost:
    It solves the bipartite ranking problem at the
    same time as it solves the problem of classification.

$$\ell(\lambda)$$