

Boosting = algorithm used to improve accuracy of weak learners, wrapper around other classifiers.

Recall PAC learning:

- δ = probability of failure.
- ϵ = accuracy parameter. want to make small.

PAC learning requirement:

For any choice of ϵ, δ , A should output with probability $\geq 1 - \delta$, an ϵ -accurate classifier.

A is allowed to run in time $\text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta})$,
Take # of samples $\text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta})$.

Question: what if we have an algorithm A with probability 5% outputs an ϵ -accurate classifier. How can we use A to obtain a standard PAC learner?

we want to increase that 5% probability of success to $1 - \delta$.

Solution: Run A a large # of times, say, t .

$P[A \text{ fails to output an } \epsilon\text{-accurate classifier}] \leq (0.95)^t$.
probability it fails every time for t trials. \rightarrow

we can make $(0.95)^t$ very small by choosing t to be approximately $O(\log \frac{1}{\delta})$, then we can "test" each classifier generated during these t trials to see if any of them are good classifiers.

Bottom Line: Amplifying the probability of success is not too difficult. $5\% \rightarrow 1-\delta$ by running the algorithm many times.

$$P(\text{at least 1 success}) = 1 - P(\text{all failures}) = 1 - (\text{failure prob})^{\# \text{ of trials}}$$

Trickier Question: What if ϵ is fixed to, say, 0.49?

Imagine A with probability $\geq 1-\delta$ outputs a classifier with $\epsilon = 0.49 \Rightarrow$ accuracy 51%.

Natural Question: How do we amplify/improve the accuracy parameter?

$$\boxed{A} \longrightarrow h \text{ s.t. } \text{err}(h) = 0.49$$

want $\boxed{A'} \longrightarrow h' \text{ s.t. } \text{err}(h') \leq \epsilon$
 \boxed{A}
 subroutine

Solution attempt: Try running A many times,
 h_1, \dots, h_t . Take the majority (h_1, \dots, h_t) .

Doesn't make sense, since ϵ is the same for all.

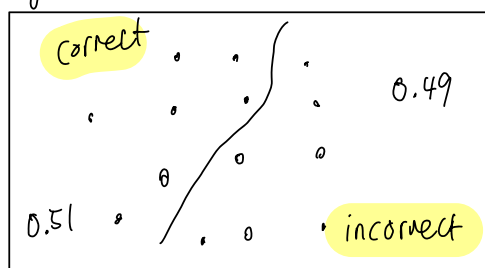
Question was posed by Valiant on PAC learning.

Solved by R. Schapire.

Solution by Freund and Schapire.

Algorithm presented: "Adaboost".

High level idea of adaboost:



← training set.

Run A on uniform distribution on points in the training set.

Assumes A outputs classifier h w/ $\epsilon = 0.49$.

Fraction of correct points is 0.51.

New distribution:

Put a little less weight on those we got right.
 " " " more " wrong.

Divide by sum of total weights \rightarrow new probability dist on points that we got wrong.
core idea:

- re-weight points we get wrong: have more weight.
- re-weight points we got right: have less weight.
- Run algorithm A again, to obtain classifier wrt the new weighting.
- At the end, take majority of classifiers generated during the process.

Adaboost (simplified) Algorithm: Ada = adapt

- training set of size m .

1) Initially D_0 = uniform distribution corresponds to $w_i = 1 \quad \forall i$: each point has weight 1.

Dist is obtained by dividing W = sum of weights.
 \uparrow capital W

2) let:

error rate $\equiv E$. At each iteration, assume: $E = \frac{1}{2} - \gamma$

accuracy $A \equiv 1 - E$.

update factor $\beta \equiv E/A = \frac{E}{1-E} = \frac{\frac{1}{2} - \gamma}{1 - E} = \frac{\frac{1}{2} - \gamma}{1 - (\frac{1}{2} - \gamma)}$

$$= \boxed{\frac{\frac{1}{2} - \gamma}{\frac{1}{2} + \gamma}} \leq 1 (?)$$

3) How to update the weights:

At iteration t , run A to obtain h_t .

- For each x_i s.t. $h_t(x_i)$ is correct:

$$\text{update } w_i^{\text{new}} = \beta \cdot w_i^{\text{old}}$$

- For each x_i s.t. $h_t(x_i)$ is incorrect: do nothing,

4) Repeat for T steps, output $\text{MAJ}(h_1, \dots, h_T)$.

Claim: After T iterations, error $h_{\text{final}} = \text{MAJ}(h_1, \dots, h_T)$

$$\leq e^{-2T\gamma^2} \Rightarrow \text{choose } T \approx \frac{1}{2\gamma^2} \log\left(\frac{1}{\epsilon}\right) \text{ then}$$

$$\text{error of } h_{\text{final}} \leq \epsilon.$$

total weight after an iteration:

- W is the weight of all points before iteration t .

- What is the weight of correct points after iteration t ?

$$W_t^+ = A \cdot \beta \cdot W = \left(\frac{1}{2} + \gamma\right) \cdot \beta \cdot W$$

\uparrow accuracy \uparrow update factor

- What is weight of the incorrect points after iteration t ?

$$W_t^- = E \cdot W = \left(\frac{1}{2} - \gamma\right) \cdot W \text{ since we don't update their weights.}$$

\uparrow error rate

- New sum of all weights (add W_t^+ and W_t^-).

$$W_t^+ + W_t^- = W \left(\frac{1}{2} \beta + \gamma \beta + \frac{1}{2} - \gamma \right)$$

$$\Rightarrow W \left(\left(\frac{1}{2} + \gamma \right) \cdot \beta + \frac{1}{2} - \gamma \right) = W \left(\frac{\left(\frac{1}{2} + \gamma \right) \left(\frac{1}{2} - \gamma \right)}{\left(\frac{1}{2} + \gamma \right)} + \frac{1}{2} - \gamma \right)$$

$$= W(1 - 2\gamma) = W(2 \cdot (\frac{1}{2} - \gamma))$$

After i^{th} iterations, the sum of the weights =

$$W_i = W_0 (2(\frac{1}{2} - \gamma))^i$$

After T iterations, the sum of all weights

$$\leq (2(\frac{1}{2} - \gamma))^T \cdot W_0 \leftarrow \text{initial weight.}$$

Consider a point X_i that h_{final} gets wrong: that means

$$W_T(X_i) \geq \beta^{T/2} \text{ (lower bound).}$$

If h_{final} gets it wrong then it was wrong for at least $T/2$ iterations. So for at least $T/2$ iterations we didn't multiply by β . That means its weight has to be at least $\beta^{T/2}$.

\Rightarrow If h_{final} has error ϵ , then weight of points h_{final} misclassifies $\geq \epsilon m \beta^{T/2}$ (lower bound).

$$\underbrace{\epsilon m \beta^{T/2}}_{\substack{\text{sum of weights} \\ \text{of those classified} \\ \text{incorrectly. Lower} \\ \text{bound.}}} \leq \underbrace{(2 \cdot (\frac{1}{2} - \gamma))^T \cdot m}_{\substack{\text{sum of all the weights,} \\ \text{upper bound.}}} \xRightarrow[\text{cancel } m]{} \epsilon \beta^{T/2} \leq (2\epsilon)^T \xRightarrow[\text{divide by } \beta^{T/2}]{} \epsilon \leq \frac{(2\epsilon)^T}{\beta^{T/2}}$$

$$\epsilon \leq \frac{(2\epsilon)^T}{\beta^{T/2}} \Rightarrow \epsilon \leq \frac{(2\epsilon)^{2T/2}}{\beta^{T/2}} \Rightarrow \epsilon \leq \left(\frac{4\epsilon^2}{\beta} \right)^{T/2}$$

Using $\beta = \frac{\frac{1}{2} - \gamma}{\frac{1}{2} + \gamma}$

$$\epsilon \leq (1 - 4\gamma^2)^{T/2} \leq e^{-2\gamma^2 T} \quad (1 + x \approx e^x)$$

$$\Rightarrow \boxed{\epsilon \leq e^{-2T\gamma^2}}$$

The previous algorithm is simplified bc we assumed the accuracy was exactly $\frac{1}{2} + \gamma$. Can make small modifications such that if in some iterations you get a very accurate diagnosis, this algorithm will adapt and converge to a small error and quicker # of iterations.

Adaboost Modifications

In adaboost,

$$\beta_t = \frac{\bar{E}_t}{A_t}$$

$$h_t \in \{-1, +1\}$$

weighing factor of each subclassifier.

output: $\text{sign} \left(\sum_t \alpha_t h_t - \frac{1}{2} \right)$, where $\alpha_t = \frac{\log \left(\frac{1}{\beta_t} \right)}{\sum_t \log \left(\frac{1}{\beta_t} \right)}$

↑
subclassifier

How do we guarantee that h_{final} generalizes?

We need to make sure that my # of training points is sufficiently large.

If γ (accuracy) is independent from m (size of training set), then we can choose m to be sufficiently large.

Adaboost is a special case of an algorithm due to Freund and Schapire called "hedge" algorithms.

Hedge: "Best Expert" set up:

C_1, \dots, C_m at each iteration t , expert

C_i suffers a loss $l_i^t \in [0, 1]$

l^t = a vector of losses suffered by all experts at the t^{th} iteration,

Intuition: we want to have a mixed strategy of experts, a weight average of experts.

Goal: the sum of our losses after T iterations should be "close" to the best expert in hindsight.

At each iteration, we maintain a set of weights:

w_1, \dots, w_m weighted average = $\frac{w_i}{\sum_i w_i} = p_i$

$w_1, \dots, w_m \longrightarrow$ probability distribution p^t

$$\underbrace{p_i^t}_{\text{weighted average}} = \frac{w_i^t}{\sum_i w_i}$$

* loss we suffer at the t iteration is $\underbrace{p^t \odot l^t}_{\substack{\text{weighted average of loss} \\ \text{of experts.}}}$ dot product.

* Total loss we suffer after T iterations:

$$= \sum_{t=1}^T p^t \odot l^t$$

Hedge Algorithm : $W_i^{\text{new}} = W_i^{\text{old}} \cdot \beta^{l_i^t}$

\nwarrow t^{th} iteration \nwarrow raised to power l_i^t
 $\beta^{l_i^t}$

Claim: $\text{your loss} \leq \underbrace{\min_i \sum_{t=1}^T l_i^t}_{\text{loss of best expert}} + O(\sqrt{T \log n})$

After T iterations, the sum of all of those average losses, is actually going to be the loss of the best expert (since we take a min over i), plus $O(\sqrt{T \log n})$.

$$\text{Your average loss} = \frac{\text{Your Loss}}{T} \leq \frac{\min_i \sum_{t=1}^T l_i^t}{T} + \underbrace{\frac{O(\sqrt{T \log n})}{T}}_{\substack{\text{goes to } 0 \\ \text{quickly if } T \\ \text{is large.}}}$$

(total loss divided by T)

Hedge Algorithm = multiplicative weight update algorithm.

In contrast to additive update algorithms (GD or Perceptron).