

COMS 4721: Machine Learning for Data Science

Lecture 13, 3/2/2017

Prof. John Paisley

Department of Electrical Engineering
& Data Science Institute
Columbia University

BOOSTING

General method for taking any binary classifier and improving it.

BAGGING CLASSIFIERS

Algorithm: Bagging binary classifiers

Given $(x_1, y_1), \dots, (x_n, y_n), x \in \mathcal{X}, y \in \{-1, +1\}$

- ▶ For $b = 1, \dots, B$ \hookrightarrow labelled data.

\hookrightarrow same size as original dataset.

sample uniformly with replacement.

- ▶ Sample a bootstrap dataset \mathcal{B}_b of size n . For each entry in \mathcal{B}_b , select (x_i, y_i) with probability $\frac{1}{n}$. Some (x_i, y_i) will repeat and some won't appear in \mathcal{B}_b .
- ▶ Learn a classifier f_b using data in \mathcal{B}_b .

- ▶ Define the classification rule to be

(So we now have B classifiers where each one is learned on a specific bootstrapping dataset.)

$$f_{\text{bag}}(x_0) = \text{sign} \left(\sum_{b=1}^B f_b(x_0) \right)$$

\downarrow bagging of classifiers

\rightarrow evaluate x_0 on each classifier & then simply take a majority vote.

- ▶ With bagging, we observe that a *committee* of classifiers votes on a label.
- ▶ Each classifier is learned on a *bootstrap sample* from the data set.
- ▶ Learning a collection of classifiers is referred to as an *ensemble method*.

BOOSTING

(development of this idea, motivation → how can we take a mediocre classifier, and boost them so the sum of these classifiers gives a

How is it that a committee of blockheads can somehow arrive at highly reasoned decisions, very accurate prediction despite the weak judgment of the individual members?

- Schapire & Freund, "Boosting: Foundations and Algorithms" *John*

Boosting is another powerful method for ensemble learning. It is similar to bagging in that a set of classifiers are combined to make a better one.

It works for any classifier, but a "weak" one that is easy to learn is usually chosen. (weak = accuracy a little better than random guessing)

Short history

1984 : Leslie Valiant and Michael Kearns ask if "boosting" is possible.

1989 : Robert Schapire creates first boosting algorithm.

1990 : Yoav Freund creates an optimal boosting algorithm.

1995 : Freund and Schapire create AdaBoost (Adaptive Boosting), the major boosting algorithm.

So boosting at a very high level is almost the same but different in a very crucial way.

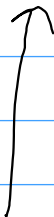
So what boosting does is at first creates a weighted sample of our data set to learn a classifier.

And then given this weighted sample it generates a new weighted sample.

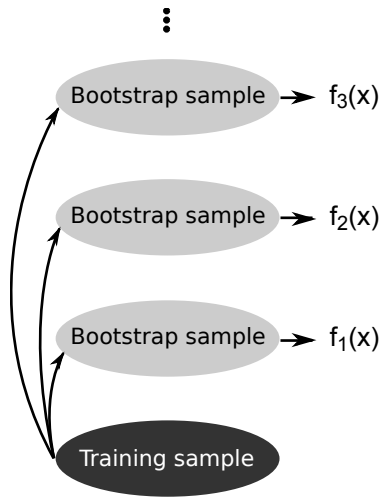
So it's not pulling from the original training sample.

It's pulling from the weighted sample of the previous round of classification.

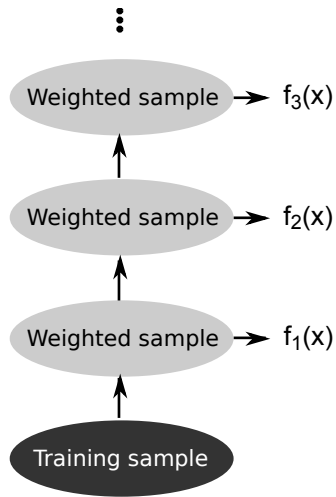
To generate a new data set for which we get a new classifier.



BAGGING VS BOOSTING (OVERVIEW)

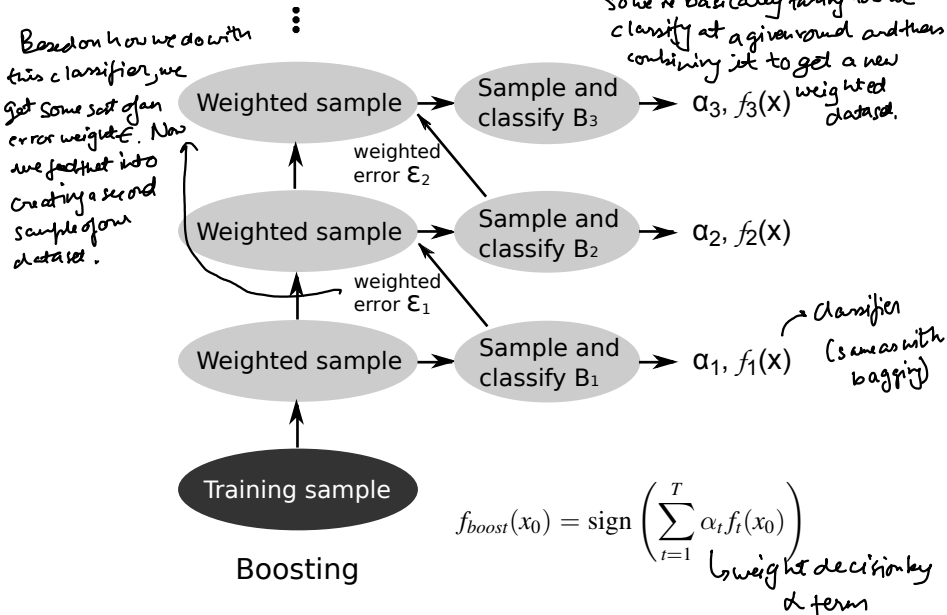


Bagging



Boosting

THE ADABOOST ALGORITHM (SAMPLING VERSION)



THE ADABOOST ALGORITHM (SAMPLING VERSION)

Algorithm: Boosting a binary classifier

Given $(x_1, y_1), \dots, (x_n, y_n)$, $x \in \mathcal{X}$, $y \in \{-1, +1\}$, set $w_1(i) = \frac{1}{n}$ for $i = 1:n$

For $t = 1, \dots, T$

1. Sample a bootstrap dataset \mathcal{B}_t of size n according to distribution w_t .
Notice we pick (x_i, y_i) with probability $w_t(i)$ and not $\frac{1}{n}$.

2. Learn a classifier f_t using data in \mathcal{B}_t .

3. Set $\epsilon_t = \sum_{i=1}^n w_t(i) \mathbb{1}\{y_i \neq f_t(x_i)\}$ and $\alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$.

4. Scale $\hat{w}_{t+1}(i) = w_t(i) e^{-\alpha_t y_i f_t(x_i)}$ and set $w_{t+1}(i) = \frac{\hat{w}_{t+1}(i)}{\sum_j \hat{w}_{t+1}(j)}$

Set the classification rule to be

$$f_{\text{boost}}(x_0) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x_0) \right).$$

Comment: Description usually simplified to “learn classifier f_t using distribution w_t .”

Calculate the prob. of error at iteration t on the distribution on the data at iteration t .

update the weights at iteration $t+1$ for each of the n data points.

data period.
iteration no.
Like a prob. distribution over the data. Initial set to be uniform. [same we used in hypothesis]

because we want to turn it into a prob. wt we simply normalize.

* 1st:

So on the first iteration we sample a bootstrap data set of size n using a uniform distribution on our data. (Idea from bagging)

t^m :

However, at iteration t this distribution on the data might be different. And so whatever the distribution on our data set is at iteration t , we sample n times from that to construct a bootstrap sample bt . Then we learn a classifier at iteration t using the data in bt .

Difference from bagging:

- ① So this is exactly the same thing as was done in bagging. The only difference is that the way that we constructed this data set was different because it's not necessarily with a uniform distribution. So that's the key difference.
- ② i) Another key difference is that after we construct this classifier, before we can go to the next iteration. We need to first decide, how do we update this weight probability distribution on our data?

ii) And also we have to get these alphas that I've discussed previously that are used in the classification.

* * We then calculate the probability of error at iteration t according to the distribution on the data at iteration t . So what that's saying is that the error at iteration t is going to be equal to the sum of all of the probability weights on the data that are misclassified. According to the classifier that we learned only at iteration t .

↳ So, we learned f_t from the data in b_t , we then evaluate that classifier at all n data points in our original data set. And if our classifier gets it wrong then we include that probability of that observation and sum up those misclassified weights.

We're summing the weights only on the misclassified data according to the classifier that we learnt at iteration t .

So we do this many times, capital T times, where T is arbitrary but large. And then we get out of this T pairs, capital T pairs of classifier and its associated weight α .

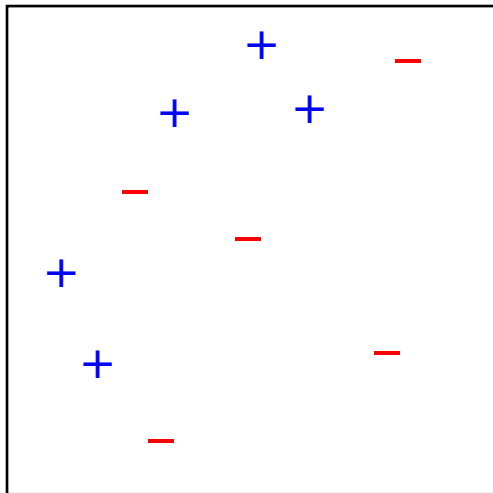
And then for a new data point, the boosted classifier is equal to the sign of each individual classifier's prediction.

Weighted by α for that particular classifier.

So we sum that up.

And then again it's like a majority vote except it's a majority weighted vote in this case. And then we take the sign, that's our final classification.

BOOSTING A DECISION STUMP (EXAMPLE 1)



Original data

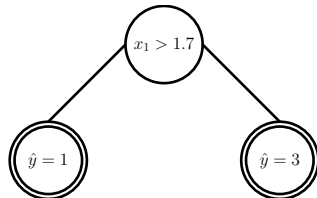
Uniform distribution, w_1

Learn *weak classifier*

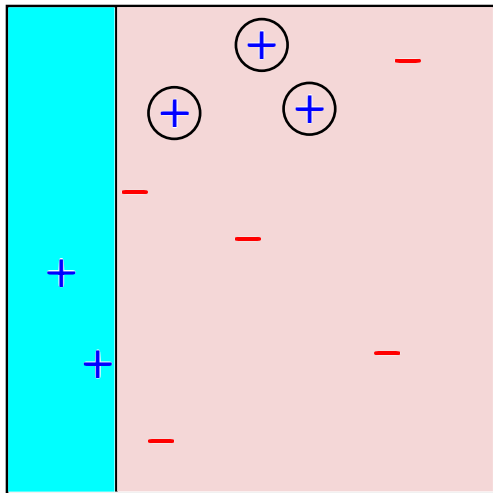
Here: Use a decision stump

Take weak
classifier's output
to make it good.

every sample has
some weight.



BOOSTING A DECISION STUMP (EXAMPLE 1)



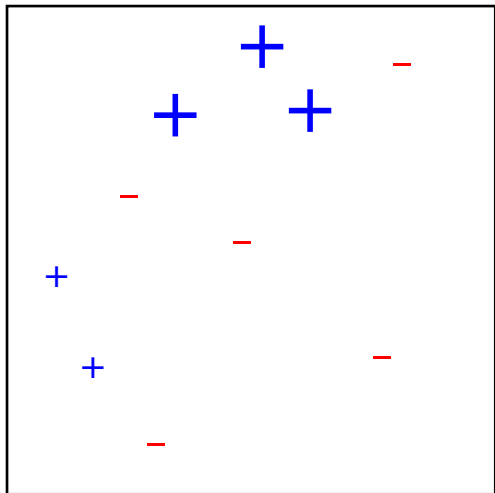
Round 1 classifier

Weighted error: $\epsilon_1 = 0.3$

Weight update: $\alpha_1 = 0.42$

$$\hookrightarrow \frac{1}{2} \ln \left(\frac{0.7}{0.3} \right)$$

BOOSTING A DECISION STUMP (EXAMPLE 1)



Now we weight data

Weighted data

After round 1

data point i
 $\hat{w}_2(i) = w_1(i) e^{-x_i y_i f_1(x_i)}$
iteration 2.

Correctly classified:

-1 label

misclassified class -1

so $y_i f_1(x_i) = +1$

$\alpha_i > 0$ always, because we're assuming our classifier is better than random guessing

$$\ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \begin{matrix} \nearrow \text{getting right} \\ \searrow \text{getting wrong} \end{matrix} > 0$$

if we are assuming it is better than random guessing.

So that's gonna be a +ve no.

For correctly classified points.

$$e^{-\alpha_i(1)} \quad \text{because } \alpha_i > 0$$

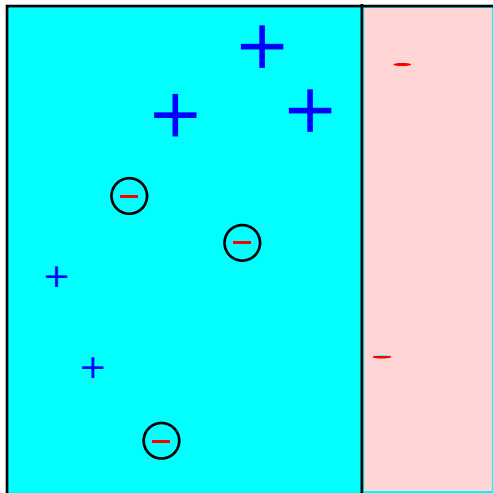
\hookrightarrow is $e^{-\alpha} < 1$

So we are downweighting all of the points that we got correct.

Whatever the weight was at previous iteration we multiply by a no. str 0 and 1, we downweight it.

Misclassified points \rightarrow upweighting
Correctly points \rightarrow

BOOSTING A DECISION STUMP (EXAMPLE 1)



Step 3: we take weights on ^{these} 3 points from the previous round.

The updated weights and sum them up.

Round 2 classifier

Weighted error: $\epsilon_2 = 0.21$

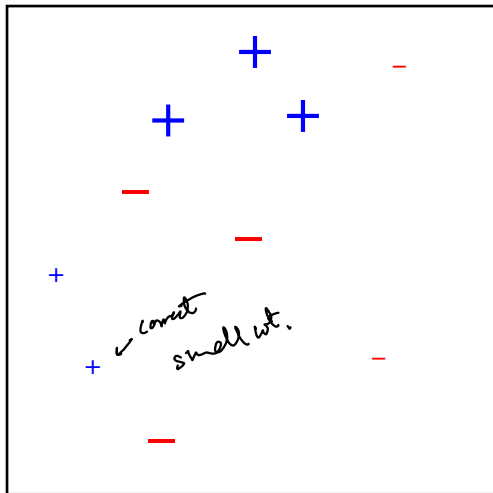
Weight update: $\alpha_2 = 0.65$

In the next bootstrap round, we have more of these misclassified points.

$\alpha_i \rightarrow$ sum of weighted errors

And update weights. All correct \rightarrow down weight
" Incorrect \rightarrow up weight

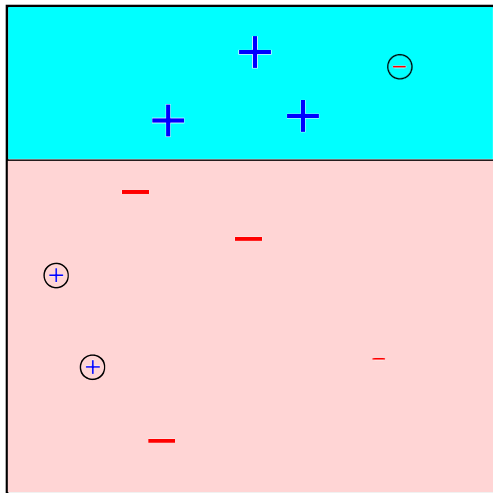
BOOSTING A DECISION STUMP (EXAMPLE 1)



Weighted data

After round 2

BOOSTING A DECISION STUMP (EXAMPLE 1)

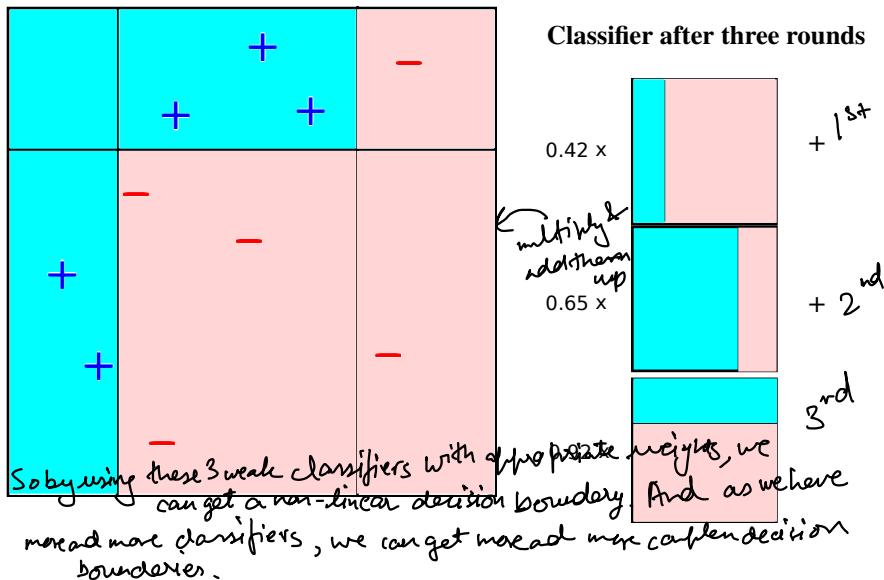


Round 2 classifier

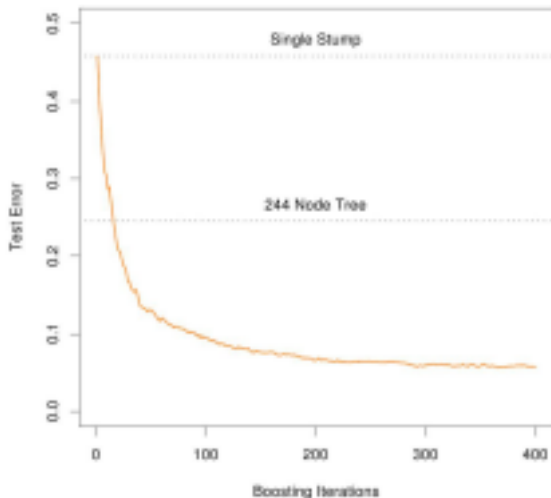
Weighted error: $\epsilon_3 = 0.14$

Weight update: $\alpha_3 = 0.92$

BOOSTING A DECISION STUMP (EXAMPLE 1)



BOOSTING A DECISION STUMP (EXAMPLE 2)



Example problem

Random guessing
50% error

Decision stump
45.8% error

Full decision tree
24.7% error

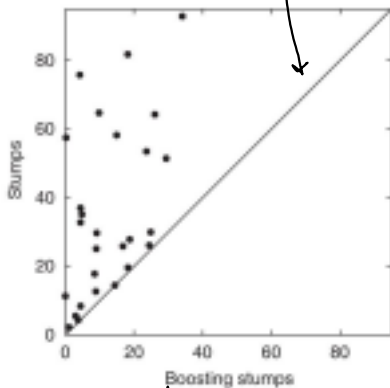
Boosted stump
5.8% error

BOOSTING

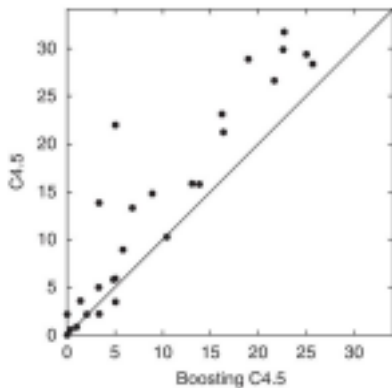
Many datasets

above the line here
we are boosting improved
results.

each point corresponds to a dataset

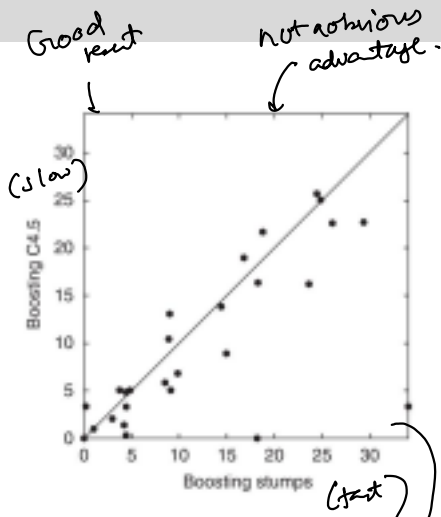
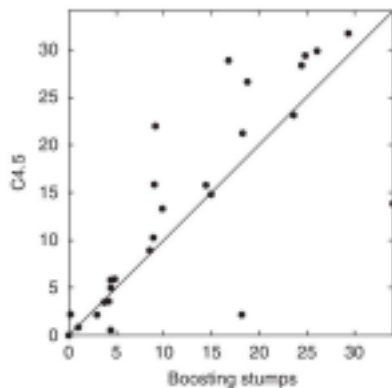


error rate.



Point = one dataset. Location = error rate w/ and w/o boosting. The boosted version of the same classifier almost always produces better results.

BOOSTING



(left) Boosting a bad classifier is often better than not boosting a good one.

(right) Boosting a good classifier is often better, but can take more time.

Overall not clear that boosting a decision stump is in general worse than boosting a very good decision tree.

BOOSTING AND FEATURE MAPS

Q: What makes boosting work so well?

A: This is a well-studied question. We will present one analysis later, but we can also give intuition by tying it in with what we've already learned.

The classification for a new x_0 from boosting is

$$\text{decision of classifier} f_{\text{boost}}(x_0) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x_0) \right).$$

Handwritten notes: An arrow points from $f_t(x_0)$ to $+1 \text{ or } -1$. Another arrow points from $\sum_{t=1}^T$ to 1000 or more .

Define $\phi(x) = [f_1(x), \dots, f_T(x)]^\top$, where each $f_t(x) \in \{-1, +1\}$.

- ▶ We can think of $\phi(x)$ as a high dimensional feature map of x .
- ▶ The vector $\alpha = [\alpha_1, \dots, \alpha_T]^\top$ corresponds to a hyperplane.
- ▶ So the classifier can be written $f_{\text{boost}}(x_0) = \text{sign}(\phi(x_0)^\top \alpha)$. *dot product*
- ▶ Boosting learns the feature mapping and hyperplane simultaneously.

Boosting is trying to learn in an online way (adaptive & online method) learning a high dimensional feature mapping along with the decision. The linear decision boundary defined by α such that we can classify all of the data in our training set.

↳ Something prove: That's going to learn a linear classifier that perfectly classifies our training set.

APPLICATION: FACE DETECTION

FACE DETECTION (VIOLA & JONES, 2001)

Boosting allows us to define a classifier that does ok and then boost it so that we get excellent performance.

Problem: Locate the faces in an image or video.

Processing: Divide image into patches of different scales, e.g., 24×24 , 48×48 , etc. Extract *features* from each patch.

Classify each patch as face or no face using a *boosted decision stump*. This can be done in real-time, for example by your digital camera (at 15 fps).



- ▶ One patch from a larger image. Mask it with many “feature extractors.”
- ▶ Each pattern gives one number, which is the sum of all pixels in black region minus sum of pixels in white region (total of 45,000+ features).

FACE DETECTION (EXAMPLE RESULTS)



ANALYSIS OF BOOSTING

ANALYSIS OF BOOSTING

Training error theorem

We can use *analysis* to make a statement about the accuracy of boosting on the training data.

Theorem: Under the AdaBoost framework, if ϵ_t is the weighted error of classifier f_t , then for the classifier $f_{\text{boost}}(x_0) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x_0))$,

$$\text{training error} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i \neq f_{\text{boost}}(x_i)\} \leq \exp\left(-2 \sum_{t=1}^T (\frac{1}{2} - \epsilon_t)^2\right).$$

Annotations:

- $\frac{1}{n}$: data set size
- $\sum_{i=1}^n$: misclassified point
- $\mathbb{1}\{y_i \neq f_{\text{boost}}(x_i)\}$: goes to 0 (adding small no.)
- $\sum_{t=1}^T$: big no. \downarrow small no.
- $(\frac{1}{2} - \epsilon_t)^2$: sum of each round of boosting.

Even if each ϵ_t is only a little better than random guessing, the sum over T classifiers can lead to a large negative value in the exponent when T is large.

For example, if we set:

$$\epsilon_t = 0.45, T = 1000 \rightarrow \text{training error} \leq 0.0067.$$

So this theorem shows that a bunch of these weak classifiers that don't do very well on their own when combined pushes the training error down to 0 as the no. of classifiers we learn increases to ∞ .

PROOF OF THEOREM

Setup

We break the proof into three steps. It is an application of the fact that

$$\text{if } \underbrace{a < b}_{\text{Step 2}} \quad \text{and} \quad \underbrace{b < c}_{\text{Step 3}} \quad \text{then} \quad \underbrace{a < c}_{\text{conclusion}}$$

- ▶ Step 1 calculates the value of b .
- ▶ Steps 2 and 3 prove the two inequalities.

Also recall the following step from AdaBoost:

- ▶ Update $\hat{w}_{t+1}(i) = w_t(i)e^{-\alpha_t y_i f_t(x_i)}$.
- ▶ Normalize $w_{t+1}(i) = \frac{\hat{w}_{t+1}(i)}{\sum_j \hat{w}_{t+1}(j)} \longrightarrow$ Define $Z_t = \sum_j \hat{w}_{t+1}(j)$.

PROOF OF THEOREM ($a \leq \mathbf{b} \leq c$)

Step 1

We first want to expand the equation of the weights to show that

$$w_{T+1}(i) = \frac{1}{n} \frac{e^{-y_i \sum_{t=1}^T \alpha_t f_t(x_i)}}{\prod_{t=1}^T Z_t} = \frac{1}{n} \frac{e^{-y_i f_{boost}^{(T)}(x_i)}}{\prod_{t=1}^T Z_t} \quad (f_{boost}^{(T)} \text{ is up to step } T)$$

Derivation of Step 1:

Notice the update rule: $w_{t+1}(i) = \frac{1}{Z_t} w_t(i) e^{-\alpha_t y_i f_t(x_i)}$

Do the same expansion for $w_t(i)$ and continue until reaching $w_1(i) = \frac{1}{n}$,

$$w_{T+1}(i) = w_1(i) \frac{e^{-\alpha_1 y_i f_1(x_i)}}{Z_1} \times \cdots \times \frac{e^{-\alpha_T y_i f_T(x_i)}}{Z_T}$$

The product $\prod_{t=1}^T Z_t$ is “ \mathbf{b} ” above. We use this form of $w_{T+1}(i)$ in Step 2.

PROOF OF THEOREM ($\mathbf{a} \leq \mathbf{b} \leq \mathbf{c}$)

Step 2

Next show that the training error of $f_{boost}^{(T)}$ after T steps is $\leq \prod_{t=1}^T Z_t$.

From Step 1: $w_{T+1}(i) = \frac{1}{n} \frac{e^{-y_i f_{boost}^{(T)}(x_i)}}{\prod_{t=1}^T Z_t} \implies w_{T+1}(i) \prod_{t=1}^T Z_t = \frac{1}{n} e^{-y_i f_{boost}^{(T)}(x_i)}$

Derivation of Step 2:

(Observe that $0 < e^{z_1}$ and $1 < e^{z_2}$ for any $z_1 < 0 < z_2$.)

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i \neq f_{boost}^{(T)}(x_i)\} &\leq \frac{1}{n} \sum_{i=1}^n e^{-y_i f_{boost}^{(T)}(x_i)} \\ &= \sum_{i=1}^n w_{T+1}(i) \prod_{t=1}^T Z_t \\ &= \prod_{t=1}^T Z_t \end{aligned}$$

“a” is the training error – the quantity we care about.

PROOF OF THEOREM ($a \leq \mathbf{b} \leq \mathbf{c}$)

Step 3

The final step is to calculate an upper bound on Z_t , and by extension $\prod_{t=1}^T Z_t$.

Derivation of Step 3:

This step is slightly more involved. It also shows why $\alpha_t := \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$.

$$\begin{aligned} Z_t &= \sum_{i=1}^n w_t(i) e^{-\alpha_t y_i f_t(x_i)} \\ &= \sum_{i: y_i = f_t(x_i)} e^{-\alpha_t} w_t(i) + \sum_{i: y_i \neq f_t(x_i)} e^{\alpha_t} w_t(i) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \end{aligned}$$

Remember we defined $\epsilon_t = \sum_{i: y_i \neq f_t(x_i)} w_t(i)$, the probability of error for w_t .

PROOF OF THEOREM ($a \leq \mathbf{b} \leq \mathbf{c}$)

Derivation of Step 3 (continued):

Remember from Step 2 that

$$\text{training error} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i \neq f_{\text{boost}}(x_i)\} \leq \prod_{t=1}^T Z_t.$$

and we just showed that $Z_t = e^{-\alpha_t}(1 - \epsilon_t) + e^{\alpha_t}\epsilon_t$.

We want the training error to be small, so we pick α_t to *minimize* Z_t . Minimizing, we get the value of α_t used by AdaBoost:

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right).$$

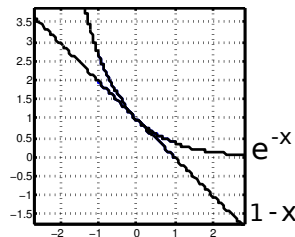
Plugging this value back in gives $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$.

PROOF OF THEOREM ($a \leq b \leq c$)

Derivation of Step 3 (continued):

Next, re-write Z_t as

$$\begin{aligned} Z_t &= 2\sqrt{\epsilon_t(1-\epsilon_t)} \\ &= \sqrt{1-4\left(\frac{1}{2}-\epsilon_t\right)^2} \end{aligned}$$



Then, use the inequality $1 - x \leq e^{-x}$ to conclude that

$$Z_t = \left(1 - 4\left(\frac{1}{2} - \epsilon_t\right)^2\right)^{\frac{1}{2}} \leq \left(e^{-4\left(\frac{1}{2} - \epsilon_t\right)^2}\right)^{\frac{1}{2}} = e^{-2\left(\frac{1}{2} - \epsilon_t\right)^2}.$$

PROOF OF THEOREM

Concluding the right inequality ($a \leq \mathbf{b} \leq c$)

Because both sides of $Z_t \leq e^{-2(\frac{1}{2}-\epsilon_t)^2}$ are positive, we can say that

$$\prod_{t=1}^T Z_t \leq \prod_{t=1}^T e^{-2(\frac{1}{2}-\epsilon_t)^2} = e^{-2 \sum_{t=1}^T (\frac{1}{2}-\epsilon_t)^2}.$$

This concludes the “ $b \leq c$ ” portion of the proof.

Combining everything

$$\text{training error} = \overbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i \neq f_{\text{boost}}(x_i)\}}^a \leq \overbrace{\prod_{t=1}^T Z_t}^b \leq \overbrace{e^{-2 \sum_{t=1}^T (\frac{1}{2}-\epsilon_t)^2}}^c.$$

We set out to prove “ $a < c$ ” and we did so by using “ b ” as a stepping-stone.

TRAINING VS TESTING ERROR

Q: Driving the training error to zero leads one to ask, does boosting overfit?

A: Sometimes, but very often it doesn't!

