# BOOSTING

Boosting is an _____ method.

Unlike bagging or random forests, boosting determines a _____ majority vote.

In particular, if the class labels are $y = +1, -1$, then boosting determines

$$f_1, \ldots, f_T \qquad \text{_____} \text{ classifiers}$$

$$\alpha_1, \ldots, \alpha_T > 0 \qquad \text{weights}$$

and returns the ensemble rule

$$h_T(x) \stackrel{\text{\tiny ?}}{=}$$

Intuitively, $\alpha_t$ reflects the _____ in $f_t$.

Let training data $(x_1, y_1), \ldots, (x_n, y_n)$ be fixed.
Let $\mathcal{F}$ be a fixed set of classifiers

(B) called the _____ class.

Definition | A _____ _____ for $\mathcal{F}$ is

a rule that takes as input a set of

weights $w_1, \ldots, w_n$ $\left( w_i \geq 0, \sum_{i=1}^{n} w_i = 1 \right)$

and returns a classifier $f \in \mathcal{F}$ such

that

is minimized (or at least small)

Notation |

A set of weights will be expressed

as a vector:

$$w = (w_1, \ldots, w_n)$$

# The Boosting Principle

Choose an initial weight vector $w^1$.

Fix $T$ and a base class $\mathcal{F}$.

For $t = 1:T$

- Given the weight

  © the _____ _____ to generate

  a classifier $f_t$

- Determine a confidence $\alpha_t > 0$

  in $f_t$

- If $f_t(x_i) \neq y_i$, then $\quad w_i^t$

  If $f_t(x_i) = y_i$, then $\quad w_i^t$

End

Output

$$h_T(x) =$$

## Examples of base classes

- Decision trees

$$+$$

$$-$$

- Decision _____ ( trees of depth ___ )

$$f(x) = \pm \ \text{sign} \left\{ x^{(j)} - c \right\}$$

- Radial basis functions

$$f(x) = \pm \ \text{sign} \left\{ k_\sigma (x - x_i) - b \right\}$$

where $k_\sigma$ is a radially symmetric _____ .


Recall the advantages of ensemble methods :

- increased stability ( decision trees )

- combine simple classifiers ( stumps, RBFs )


## Adaboost

→ the first successful boosting algorithm, introduced by Yoav Freund + Robert Schapire.

## Adaboost

Given $(x_1, y_1), \ldots, (x_n, y_n)$, $\quad y_i \in \{-1, +1\}$

Initialize $w_i^1 = \frac{1}{n}$.

For $t = 1, \ldots, T$

- Apply base learner with weights $w^t$ to produce classifier $f_t$

- Set
$$r_t = \sum_{i=1}^{n} w_i^t \, \mathbb{1}\{f_t(x_i) \neq y_i\}$$

- Set
$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 - r_t}{r_t}\right)$$

- Update
$$w_i^{t+1} = \frac{w_i^t \cdot \exp\{-\alpha_t \, y_i \, f_t(x_i)\}}{Z_t}$$

where $Z_t$ is a normalization constant

End

Output
$$h_T(x) = \text{sign}\left\{\sum_{t=1}^{T} \alpha_t \, f_t(x)\right\}$$

## Weak learning

The success of Adaboost is reflected in the following result.

Denote $\gamma_t = \frac{1}{2} - r_t$. We may assume $\gamma_t \geq 0$ (why?)

__Theorem__  The training error of Adaboost satisfies

$$\frac{1}{n} \sum_{i=1}^{h} 1_{\{h_T(x_i) \neq y_i\}} \leq \exp\left(-2 \sum_{t=1}^{T} \gamma_t^2\right)$$

In particular, if $\gamma_t \geq \gamma > 0$ for all $t$, then

Ⓔ

$$\frac{1}{n} \sum_{i=1}^{h} 1_{\{h_T(x_i) \neq y_i\}} \leq$$

The assumption $\gamma_t \geq \gamma > 0 \quad \forall t$ is
called the ___ ___ ___

and in this case the base learner is
called a ___ learner.


In words, the theorem tells us that
if our base learner does slightly better than
___ ___, the final ensemble

classifier can separate the training data
perfectly for $T$ large enough. In fact
the training error goes to zero

___ ___ .

## Remarks

Ⓔ  • If $r_t = 0$, then $\alpha_t = $

  Does this make sense?

  • Setting $T$ is a ＿＿＿＿ ＿＿＿＿ problem. If $T$ is too large we may experience ＿＿＿＿＿＿＿. Cross-validation is a common approach.

  • Empirical results suggest that Adaboost with decision trees is one of the best "off-the-shelf" methods for classification.

Proof of Theorem 1 The proof is broken down into some lemmas.

Lemma 1

$$\frac{1}{n} \sum_{i=1}^{n} 1_{\{h_T(x_i) \neq y_i\}} \leq \prod_{t=1}^{T} Z_t$$

Proof By unraveling the update rule we find

$$w_i^{T+1} = \frac{w_i^T \exp\left(-\alpha_T y_i f_T(x_i)\right)}{Z_T}$$

$$= \frac{w_i^{T-1} \exp\left(-y_i\left[\alpha_{T-1} f_{T-1}(x_i) + \alpha_T f_T(x_i)\right]\right)}{Z_{T-1} \cdot Z_T}$$

$$\vdots$$

$$= \frac{\frac{1}{n} \cdot \exp\left(-y_i \sum_{t=1}^{T} \alpha_t f_t(x_i)\right)}{Z_1 \cdot Z_2 \cdots Z_T}$$
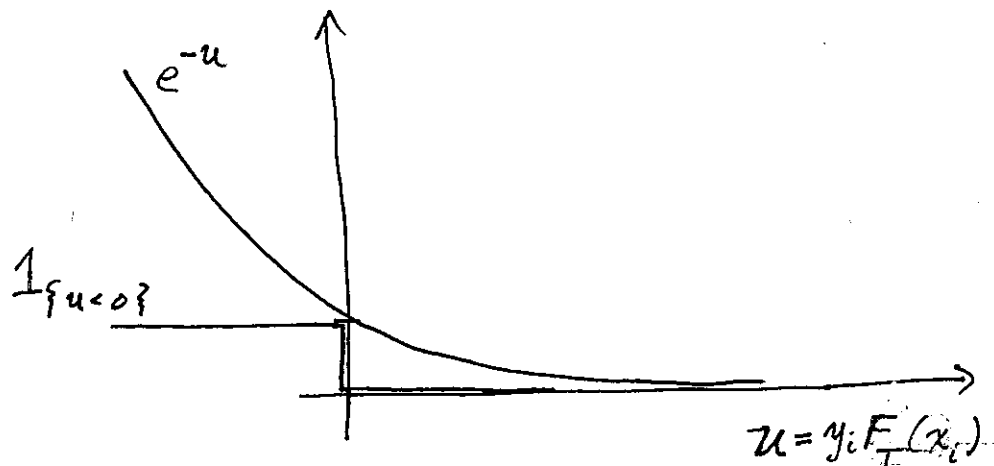
$\boxed{w_i^1}$

$$= \frac{\exp\left(-y_i F_T(x_i)\right)}{n \prod_{t=1}^{T} Z_t}$$

where $F_t = \sum_{s=1}^{t} \alpha_s f_s$

and $h_T(x) = \text{sign}\{F_T(x)\}$

Now use the bound

$$1_{\{h_T(x_i) \neq y_i\}} = 1_{\{y_i F_T(x_i) < 0\}} \leq \exp\left(-y_i F_T(x_i)\right)$$



Then

$$1 = \sum_{i=1}^{n} w_i^{T+1}$$

$$= \sum_{i=1}^{n} \frac{\exp\left(-y_i F_T(x_i)\right)}{n \cdot (\pi Z_t)}$$

$$\geq \frac{1}{(\pi Z_t)} \cdot \frac{1}{n} \sum_{i=1}^{n} 1_{\{h_T(x_i) \neq y_i\}}$$

and the lemma follows.

**Lemma** $Z_t = \sqrt{1 - 4\gamma_t^2}$

**Proof** $Z_t = \sum_{i=1}^{n} \underbrace{w_i^t \exp(-\alpha_t y_i f_t(x_i))}_{w_i^{t+1}}$

$$= \sum_{i: f_t(x_i) = y_i} w_i^t \exp(-\alpha_t) + \sum_{i: f_t(x_i) \neq y_i} w_i^t \exp(\alpha_t)$$

$$= (1 - r_t) e^{-\alpha_t} + r_t e^{\alpha_t}$$

Now recall

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 - r_t}{r_t}\right) = \ln\sqrt{\frac{1 - r_t}{r_t}}$$

Then

$$Z_t = (1 - r_t)\sqrt{\frac{r_t}{1 - r_t}} + r_t\sqrt{\frac{1 - r_t}{r_t}}$$

$$= 2\sqrt{r_t(1 - r_t)}$$
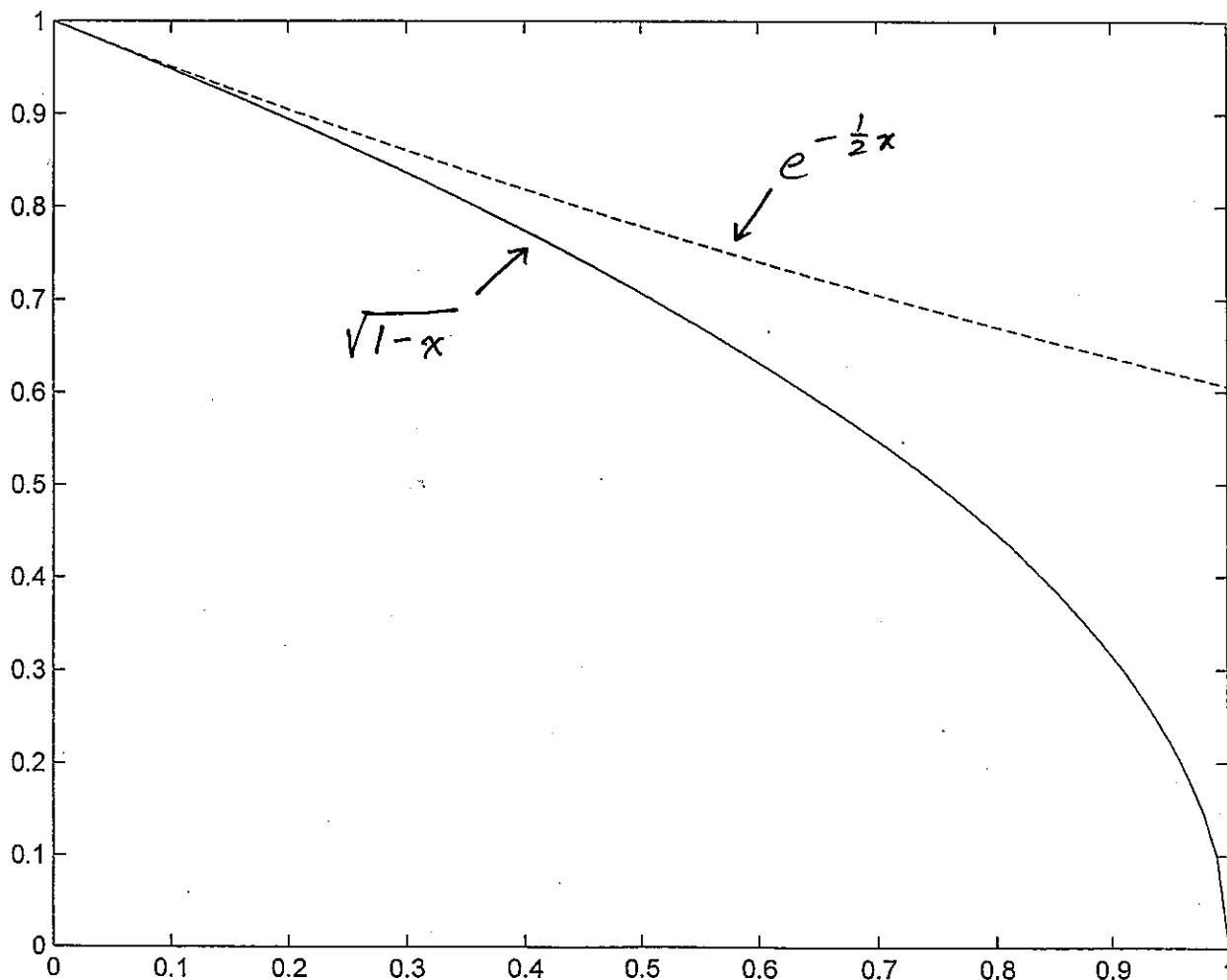
Now substitute

$$r_t = \frac{1}{2} - \gamma_t$$

$$\Rightarrow \quad z_t = 2\sqrt{\left(\tfrac{1}{2}-\gamma_t\right)\left(\tfrac{1}{2}+\gamma_t\right)}$$

$$= 2\sqrt{\tfrac{1}{4}-\gamma_t^2}$$

$$= \sqrt{1-4\gamma_t^2}$$

## Lemma

$$\sqrt{1-x} \leq e^{-\tfrac{1}{2}x}$$

## Proof

Formally, $\sqrt{1-x}$ is concave, $e^{-\frac{1}{2}x}$ is convex, so it suffices to show their slopes (derivatives) are both $= -\frac{1}{2}$ at $0$.

Putting it all together, we obtain

$$\frac{1}{n} \sum_{i=1}^{n} \mathbb{I}_{\{h_T(x_i) \neq y_i\}} \leq \frac{1}{n} \sum_{i=1}^{n} \exp\left(-y_i F_T(x_i)\right)$$

$$= \prod_{t=1}^{T} Z_t$$

$$= \prod_{t=1}^{T} \sqrt{1 - 4\gamma_t^2}$$

$$\leq e^{-2\sum_{t=1}^{T} \gamma_t^2}$$

Exercise / View $Z_t$ as a function of $\alpha_t$, and find the value of $\alpha_t$ that minimizes $Z_t$.

## Solution

Earlier we showed

$$Z_t = (1 - r_t)e^{-\alpha_t} + r_t e^{\alpha_t}.$$

This is a convex, differentiable function of $\alpha_t$.

It is minimized by setting

$$0 = \frac{\partial Z_t}{\partial \alpha_t} = -(1 - r_t)e^{-\alpha_t} + r_t e^{\alpha_t}$$

$$\Rightarrow \quad e^{2\alpha_t} = \frac{1 - r_t}{r_t}$$

$$\Rightarrow \quad \alpha_t = \frac{1}{2} \ln\left(\frac{1 - r_t}{r_t}\right)$$

In conclusion, each $\alpha_t$ is choosen to minimize the corresponding term $Z_t$ in the bound $\prod_{t=1}^{T} Z_t$.

That is, the bound is minimized <u>incrementally</u> (not <u>globally</u>)

# Boosting as Functional Gradient Descent

We will now generalize Adaboost to an entire class of boosting algorithms.
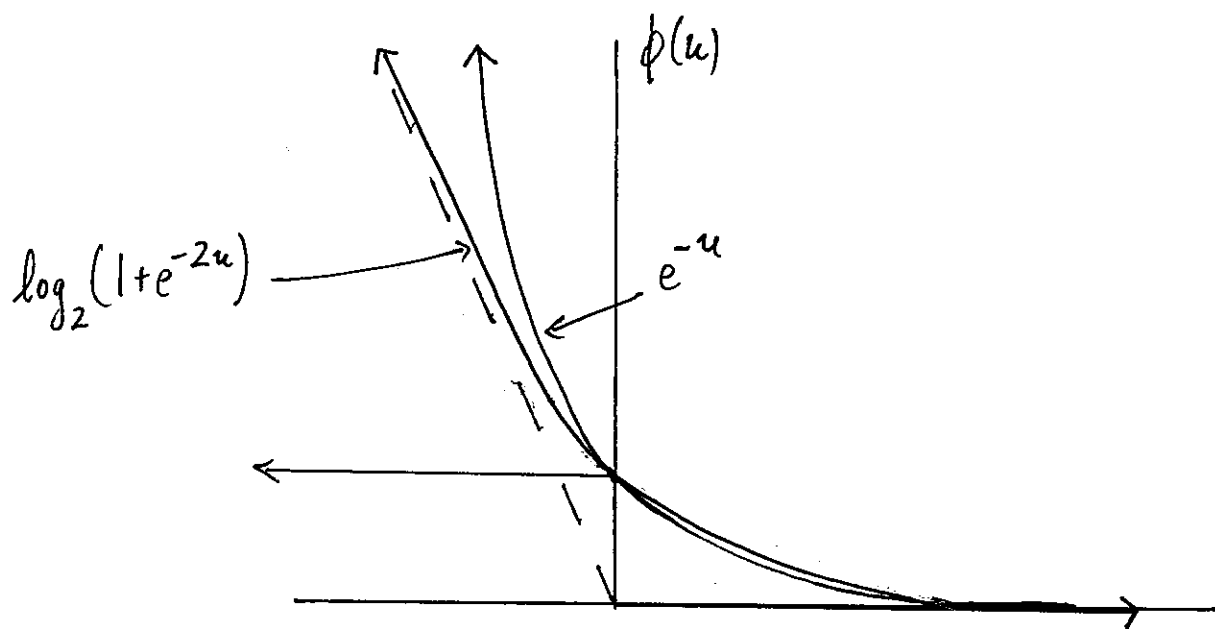
Recall that in bounding the Adaboost training error we used the bound

(G) $\qquad 1_{\{u < 0\}} \leq \qquad\qquad , \quad u = y_i F_t(x_i)$

We will generalize Adaboost by using the bound

$$1_{\{u < 0\}} \leq$$

where $\phi(u)$ is called a _____ function.

For computational reasons, the loss is often choosen
to be

(H)
· 
· 
· 

Examples

- exponential
- logistic
- hinge (not differentiable, decreasing)
- squared error (not decreasing)

Why use different losses?

The logistic loss, for example, doesn't penalize
misclassified points as severely, and therefore
may be less susceptible to ——————.

Let's assume that $\phi$ is convex, differentiable,
and decreasing.

On the $t^{th}$ iteration of boosting we have the ensemble

$$F_t(x) = \sum_{s=1}^{t} \alpha_s f_s(x)$$

and the bound

$$\frac{1}{n} \sum_{i=1}^{n} 1_{\{y_i F_t(x_i) < 0\}} \leq \frac{1}{n} \sum_{i=1}^{n} \phi(y_i F_t(x_i))$$

View this bound as an objective function to be minimized with respect to $F_t$.

Boosting may be viewed as _functional_, _gradient_ _descent_ applied to the bound.
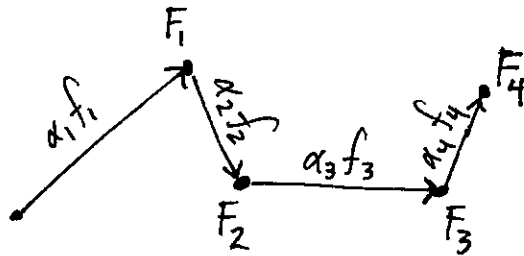
In particular, suppose $\alpha_1, f_1, \ldots, \alpha_{t-1}, f_{t-1}$ are given, and set

$$B_t(\alpha, f) = \frac{1}{n} \sum_{i=1}^{n} \phi(y_i F_{t-1}(x_i) + y_i \alpha f(x_i))$$

Then set

1) $f_t$ = function $f \in \mathcal{F}$ (base class) for
which the __directional__
__derivative__ of $B_t$ in the
direction $f$ is __minimized__ .

2) $\alpha_t$ = stepsize $\alpha > 0$ in the
direction $f_t$ for which
$B(\alpha, f_t)$ is __minimized__ .



__Step 1__ | The directional derivative of $B_t$
in the direction $f$ is

(I)
$$\left. \frac{\partial B_t(\alpha, f)}{\partial \alpha} \right|_{\alpha = 0} =$$

Minimizing this is equivalent to minimizing

$$-\sum_{i=1}^{n} y_i f(x_i) \cdot \frac{\phi'(y_i F_{t-1}(x_i))}{\underbrace{\sum_{j=1}^{n} \phi'(y_j F_{t-1}(x_j))}_{w_i^t}}$$

since $\phi' < 0$

$$= \sum_{i=1}^{n} w_i^t \mathbb{1}_{\{f(x_i) \neq y_i\}} - \sum_{i=1}^{n} w_i^t \mathbb{1}_{\{f(x_i) = y_i\}}$$

$$= 2\left(\sum_{i} w_i^t \mathbb{1}_{\{f(x_i) \neq y_i\}}\right) - 1$$

To minimize this expression with respect to $f$,

Ⓙ we can use the _____ _____ .

Step 2 |

$$\alpha_t := \arg\min_{\alpha} B_t(\alpha, f_t)$$

$$= \arg\min_{\alpha} \frac{1}{n} \sum_{i=1}^{n} \phi\left(y_i F_{t-1}(x_i) + y_i \alpha f_t(x_i)\right)$$

## Generalized Boosting Algorithm

Given $(x_1, y_1), ..., (x_n, y_n)$, $y_i \in \{-1, 1\}$, convex loss $\phi$

Initialize $w_i^1 = \frac{1}{n}$

For $t = 1, ..., T$

- Apply base learner with weights $w^t$ to produce classifier $f_t$

- Set

$$\alpha_t = \arg\min_\alpha \frac{1}{n} \sum_{i=1}^n \phi\left(y_i F_{t-1}(x_i) + y_i \alpha f_t(x_i)\right)$$

- update

$$w_i^{t+1} = \frac{\phi'(y_i F_t(x_i))}{\sum_{j=1}^n \phi'(y_j F_t(x_j))}$$

End

Output

$$h_T(x) = \text{sign}\left\{ F_T(x) \right\} = \text{sign}\left\{ \sum_{t=1}^T \alpha_t f_t(x) \right\}$$

Since $\phi$ is convex, $\alpha$ is the solution of a convex, univariate optimization problem and can be found efficiently using Newton's method.

When $\phi(u) = e^{-u}$, the algorithm simplifies to Adaboost. In this case

(K)
- $W_i^t$ has a nice _____ formula since

$$\phi'(a+b) = \phi'(a) \cdot \phi'(b)$$

- $\alpha_t$ has a closed form solution.

When $\phi(u) = \log_2(1 + e^{-2u})$, the algorithm is called_____. For computational efficiency, Friedman, Hastie, and Tibshirani suggest using only one step of Newton's method at each round.

(L) Why did we assume $\phi$ to be decreasing?

$\boxed{\text{Key}}$  A. ensemble, weighted, base

$$h_T(x) = \text{sign}\left\{ \sum_{t=1}^{T} \alpha_t f_t(x) \right\}, \qquad \text{confidence}$$

B. base, base learner, $\sum_{i=1}^{n} w_i \mathbb{1}\{f(x_i) \neq y_i\}$

C. base learner, increase, decrease

$$h_T(x) = \text{sign}\left\{ \sum_{t=1}^{T} \alpha_t f_t(x) \right\}$$

D. + : increased stability, performance, − : lose interpretability

stumps, 1, kernel

E. $\exp(-2\gamma^2 T)$       F. weak learning hypothesis,

weak, random guessing, exponentially fast

F. $\infty$, model selection, overfitting

G. $e^{-u}$, $\varphi(u)$, loss

H. convex, differentiable, decreasing; outliers

I. $\frac{1}{n}\sum_{i=1}^{n} y_i f(x_i) \cdot \varphi'(y_i F_{t-1}(x_i))$        J. base learner

K. recursive, Logitboost

L. So that $\varphi' < 0$       (step 1)