



Lecture 19

Boosting as G.D.

. Froject test labels due Toony miduight • L20 -overview of Clustering

Lecture VII: Combining predictors - Part II

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

November, 2023

lovember, 2023

Boosting as descent in function space

Boosted predictors are additive models
AdaBoost is steepest descent on training set A statistical view of boosting

More surrogate losses, more boosting algorithms 🗲

- Why the e^{-yf} loss? other surrogate losses
- GRADIENTBOOST general φ

Practical remarks and theoretical results Practicalities [Theoretical results]

Extensions of boosting Boosting for multiclass and ranking •[Multiplicative updates algorithms]

Reading HTF Ch.: 15. Random Forests, 16. Ensembles of predictors, 8. Model inference and averaging, 10. Boosting, Murphy Ch.: 16.3, 16.3.1 Generalized Additive models (ignore the regularization, for "smoother" read "minimize loss function"), 16.4.1-5 [and optionally 8] Bosting, 3.2.4 Bayesian model averaging, 16.6.3, 16.6.1 Ensembles of predictors, Bach Ch.:

Idea of boosting Ada boost alg $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k}$ $S \in [0,1)$ $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k}$ $S \in [0,1)$ $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k}$ $S \in [0,1]$ $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k}$ $S \in [0,1]$ $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k}$ $\hat{L}_{\varphi}(f_{k}) \subseteq S^{k$

AdaBoost Algorithm

	AdaBoost	ALGORITHM	
	Assume	b contains functions b taking values in $[-1, 1]$ or $\{\pm 1\}$	
	Input	M, labeled training set D	
	Initialize	f = 0	
		$w_i^1 = \frac{1}{n}$ weight of datapoint x_i	
	for	$k = 1, 2, \dots M$ direction	k .
		1. "learn classifier for \mathcal{D} with weights $w^{kn} \Rightarrow$ predictor b^k	8 ≥ 1/2 - 5
Part I	- time	2. compute error $\overline{\varepsilon_{i}^{k}} = \sum_{i=1}^{n} w_{i}^{k} \frac{1-y_{i}b^{k}(x_{i})}{2}$ Weighted error	þ.
ors -		3. set $\beta^k = \frac{1}{2} \ln \frac{1 - \varepsilon^k}{\iota} > 0$	in the if which
edict	step	$\sum_{k=1}^{2} \frac{\varepsilon^{k}}{\varepsilon^{k}} = \frac{1}{\varepsilon^{k}} \frac{1}{\varepsilon^{k}$	le · ···
ng pi	510-	4. compute new weights $w_i = \frac{Z^k}{Z^k} w_i = \frac{k+1}{Z^k}$ where Z is	Bu-Con
nbini		the normalization constant that makes $\sum_{i} w_{i}^{n+1} = 1$	le iten
Cor	Output	$f(x) = \sum_{k=1}^{M} \beta^k b^k(x)$	malized
II N			Norman
Marina Meila: Lecture		$f(x) = f(x) + \beta b_{E}(x)$ constrained gradient ducent	

Gradient of
$$f \varphi$$

$$f = \left(f(x^{i}) \right) \text{ at stepk} \quad y = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$f = \left(y^{i} \right) \text{ if } y^{i}$$

$$\varphi(z) = e^{-z}$$

$$\varphi'(z) = -\varphi(z)$$

$$\frac{\partial}{\partial f}\varphi(yf) = -y\varphi(yf)$$

want
$$f^* \in \mathbb{R}^n$$
 to minimize \hat{L}_{φ} | decrease \hat{L}_{φ} iteratively
 $f_{(x)} \stackrel{*}{\underset{k}{\rightarrow}} \stackrel{*}{\underset{k}{\rightarrow} \stackrel{*}{\underset{k}{\rightarrow}} \stackrel{*}{\underset{k}{\rightarrow}} \stackrel{*}{\underset{k}{\rightarrow}} \stackrel{*}{\underset{k}{\rightarrow}} \stackrel{*}{\underset{k}{\rightarrow}}$

$$\begin{split} \widehat{L}_{p}(\underline{f}) &= \lim_{n \to \infty} \sum_{i=1}^{n} \varphi(y_{i} \widehat{f}(x_{i})) \\ &= \lim_{\substack{p \in Y \\ p \in T_{i}}} \sum_{j=1}^{n} \varphi(\varphi'(y_{j} \widehat{f}_{i})) = -\lim_{\substack{p \in Y \\ p \in T_{i}}} y_{i} \widehat{e}^{-y_{i}} \widehat{f}_{i}, \\ &= \lim_{\substack{p \in Y \\ p \in T_{i}}} \sum_{\substack{p \in Y \\ p \in T_{i}}} \sum_{p$$

AdaBoost is steepest descent on training set

We will show that boosting is a form of (stochastic) gradient descent on the surrogate loss \hat{L}_{ϕ} (we already know from Part I that ADABOOST pushes \hat{L}_{ϕ} asymptotically towards 0). Assume we want to minimize the surrogate loss \hat{L}_{ϕ} on the training set. For any finite \mathcal{D} , f and $b \in \mathcal{B}$ affect \hat{L}_{ϕ} only via the *n*-dimensional vectors of their values on \mathcal{D} (which we will abusively denote by f, b)

$$f = \begin{bmatrix} f(x^1) \\ f(x^2) \\ \dots \\ f(x^n) \end{bmatrix} \quad b = \begin{bmatrix} b(x^1) \\ b(x^2) \\ \dots \\ b(x^n) \end{bmatrix}$$
(4)

Thus, $\hat{L}_{\phi}(f)$ is a function of *n* variables, with partial derivatives

$$\frac{\partial \hat{L}_{\phi}}{\partial f(x^{i})} = \frac{\partial}{\partial f(x^{i})} \left[\frac{1}{n} \sum_{i=1}^{n} \phi(y^{i} f(x^{i})) \right] = \frac{1}{n} y^{i} \phi'(y^{i} f(x^{i})) = -\frac{1}{n} y^{i} e^{-y^{i} f(x^{i})}, \quad (5)$$

since $\phi'(z) = -e^{-z}$. Imagine a boosting step as trying to find a change βb in f which minimizes the loss $\hat{L}_{\phi}(f + \beta b)$. This minimization is equivalent to maximizing the decrease in loss $\hat{L}_{\phi}(f) - \hat{L}_{\phi}(f + \beta b)$.

The direction of descent

The change in L_{ϕ} along "direction" *b* with step size β is approximately

$$\hat{L}_{\phi}(f) - \hat{L}_{\phi}(f + \beta b) \approx -\left(\nabla_{f}\hat{L}_{\phi}(f)\right)^{T}(\beta b) = \sum_{i} \left[\left(\frac{1}{n}y^{i}e^{-y^{i}f(x^{i})}\right)(\beta b(x^{i}))\right] \propto \sum_{i}y^{i}b(x^{i})w_{i} \quad (6)$$

(denoting/recalling $w_i \propto e^{-y_i f(x^i)}$). The direction of steepest descent *b* is therefore the maximizer of

$$\underset{b \in \mathcal{B}}{\operatorname{argmax}} \sum_{i} w_{i} y_{i} b(x^{i})$$
(7)

where in the sum on the r.h.s we recognize the r of ADABOOST.

- If b(xⁱ) = ±1 values, then 1 − y_ib(xⁱ) = 1_[i error], and maximizing (7) is the same as minimizing the weighted training error L^w₀₁.
- If b takes real values, then y_ib(xⁱ) is the margin of example i, and maximizing (7) is a natural objectiv for many training algorithm. Exercise Can you find examples of algorithms/predictors which do/don't maximize the loss in (7)?

More generally (we will use this later), the direction b maximizes

$$\sum_{i} y_i b(x^i) [-\phi'(y_i f(x^i))] \tag{8}$$

Finding the direction *b* is equivalent with step 1 of the ADABOOST algorithm, training a weak classifier on the weighted data. The resulting *b* can be seen as the best approximate of the gradient of L_{ϕ} in \mathcal{B} .

Now let us do line minimization: find the optimal step size β in direction *b*. For this we take the derivative of $\hat{L}_{\phi}(f + \beta b)$ w.r.t β and set it to 0.

$$\frac{\partial \hat{L}_{\phi}(f+\beta b)}{\partial \beta} = \sum_{i} \pm b(x^{i})\phi'(y_{i}f(x^{i})) = -\sum_{i} \underbrace{y_{i}b(x^{i})e^{-y_{i}f(x^{i})-\beta y_{i}b(x^{i})}}_{W_{i}}$$
(9)

 β is the (unique) root of

$$\sum_{i} \underline{w_{i}} y_{i} \underline{b(x^{i})} e^{-\beta y_{i} b(x^{i})} = 0$$
(10)

If
$$b(x) \in \{-1,1\}$$

 $b(x) \in [-1,1]$
 $b(x) \in (-\infty,\infty)$
then line optimization gives β^k from ADABOOST
then line optimization gives β^k from ADABOOST approximately
then β amounts to a rescaling of b and is redundant.
 $y_i^{i}b(x_i^{i}) = \begin{cases} 1 & \text{iff correct} \\ -1 & \text{iff error} \\ 1 - \xi_k \\ correct \end{cases} = \beta = \begin{cases} \xi_k \\ \xi_k \\ \xi_k \end{cases}$

Marina Meila: Lecture VII: Combining predictors – Part II

Calculating β^k for binary **b**'s

Assume $b(x) \in \{\pm 1\}$. In this case $y^i b^{(x^i)} = \pm 1$ and we obtain

$$\frac{d\hat{L}_{\phi}(f+\beta b)}{d\beta} = \sum_{i \text{ corr}} w_i e^{-\beta} - \sum_{i \text{ err}} w_i e^{\beta} = 0 \quad (11)$$

$$0 = (1 - \sum_{i \text{ err}} w_i) - (\sum_{i \text{ err}} w_i) e^{2\beta} \quad (12)$$

$$\beta = \frac{1}{2} \ln \frac{1 - \varepsilon^k}{\varepsilon^k} \quad (13)$$

This is the β^k coefficient of step 4 of ADABOOST

Hence, the ADABOOST algorithm can be seen as minimizing the loss $L_{\phi}(f)$ by steepest descent in the function space span \mathcal{B} .

RealAdaBoost

The third case corresponds to the $\operatorname{REALADABOOST}$ in the FHT variant, described here for completeness

REAL ADABOOST ALGORITHM (in the FHT variant)

```
 \begin{array}{lll} \textbf{Assume} & \mathcal{B} \text{ contains real-valued functions} \\ \textbf{Input} & M, \text{ labeled training set } \mathcal{D} \\ \textbf{Initialize} & f = 0 \\ & w_i^1 = \frac{1}{n} \text{ weight of datapoint } x^i \\ \textbf{for} & k = 1, 2, \dots M \\ & \text{``learn classifier for } \mathcal{D} \text{ with weights } w^k \Rightarrow b^{k \text{''}} \\ & \text{ compute new weights } w_i^{k+1} = w_i^k e^{-y^i b^k (x^i)} \text{ and normalize them to sum to 1} \\ \textbf{Output} & f(x) = \sum_{k=1}^M b^k(x) \end{array}
```

A statistical view of boosting

It has been shown [?] (FHT) that boosting can also be seen as noisy gradient descent in function space when we replace the finite training set with the true data distribution. The loss function and gradient can be given a probabilistic interpretation. This point of view is useful in two ways:

- It shows that boosting is asymptotically minimizing a reasonable loss function, so that we can expect the performace/and algorithm behavior on finite samples to be a good predictor on its behaviour with much larger samples.
- 2. It is an interpretation that allows on to create a very large variety of boosting algorithms, like the LOGITBOST, GENTLE ADABOOST and GRADIENTBOOST, presented hereafter.

Assume

- we do boosting "at the distribution level", i.e using P_{XY} instead of the empirical distribution given by \mathcal{D} .
- The loss function is $L_{\phi}(f) = E[e^{-yf(x)}]$.
 - The notation E[] denotes expectation w.r.t the joint P_{XY} distribution.
- learning a classifier means "find the best possible minimizer to $L_{\phi}(f)$ "

Is L_{ϕ} a good loss? 1)

Proposition

Denote $p_x = P_{XY}(y = 1|x)$. The loss $L_{\phi}(f)$ is minimized by

$$f^*(x) = \frac{1}{2} \ln \frac{P_{XY}(y=1|x)}{P_{XY}(y=-1|x)} = \frac{1}{2} \ln \frac{p_x}{1-p_x}$$

And $p_x = \frac{e^{f(x)}}{e^{f(x)} + e^{-f(x)}}$ the logistic function.

Exercise Does the expression of p_x look familiar? What is the connection?

Proof Since we are minimizing over all possible f's with no restrictions, we can minimize separately for every f(x). Hence, let x be fixed

$$E_{P_{Y|X=x}}[e^{-yf(x)}] = P(y=1|x)e^{-f(x)} + P(y=-1|x)e^{f(x)}$$

and the gradient is

$$\frac{\partial E[e^{-yf(x)}|x]}{\partial f(x)} = -P(y=1|x)e^{-f(x)} + P(y=-1|x)e^{f(x)}$$

By setting this to 0 the result follows.

In summary f^* is the Bayes optimal predictor for L_{ϕ} . But by the Proposition, f^* is also Bayes optimal for L_{01} . (Good!)

Steepest descent on $L_{\phi}(f)$ is (like) REALADABOOST

Proposition

The REAL ADABOOST (with "learn a classifier" defined at the distribution level) fits an additive logistic regression model f by iterative descent on $L_{\phi}(f)$.

Proof The proof is similar to that for the training set case.

Suppose we have a current estimate f(x) and seek to improve it by minimizing $L_{\phi}(f + b)$ over b. In the proof we assume that b is an arbitrary function, while in practice b will be chosen to best approximate the ideal f within the class \mathcal{B} .

Denote by $p_x = P[y = 1|x]$ (the true value) and by \hat{p}_x the "estimate"

$$\hat{\rho}_{x} = \frac{e^{f(x)}}{e^{f(x)} + e^{-f(x)}}$$
(14)

Assume again x is fixed. Then,

$$L_{\phi}(f+b) = E[e^{-yf(x)-yb(x)}]$$

= $e^{-f(x)}e^{-b(x)}p_x + (1-p_x)e^{f(x)}e^{b(x)}$

Taking the derivative and setting it to 0 we obtain the new step:

$$b(x) = \frac{1}{2} \ln \frac{p_x e^{-f(x)}}{(1-p_x)e^{f(x)}} = \frac{1}{2} \left[\ln \frac{p_x}{1-p_x} - \ln \frac{\hat{p}_x}{1-\hat{p}_x} \right]$$
(15)

Note that if one could exactly obtain the b prescribed by (15) the iteration would not be necessary.

- Part II

(Proof, continued) More interesting than the exact form of *b* above is the optimization problem that leads to it. Denote $w(x, y) = e^{-yf(x)}$. Then, *b* is the solution of

$$b = \operatorname*{argmin}_{b \in \mathcal{B}} E_{P_{XY}w(X,Y)}[e^{-Yb}]$$
(16)

where $P_{XY}w(X, Y)$ denotes the (unnormalized) **twisted distribution** obtained by multiplying the original data distribution with w(x, y). (Of course, one may have to put some restrictions on P_{XY} and \mathcal{B} in order to obtain a proper distribution.) Finally, note that the new f is f + band the new weights are $w(x, y)e^{-yb(x)}$ which finishes the proof.

Hence, the REAL ADABOOST algorithm can be seen as a form of "noisy gradient" algorithm at the distribution level. (Note that the minimization in equation (16) is over both direction and scale of f.)

What next?
J. What about Lor?
J. What about Lor?
J. What about Lor?
J. Convergence
$$Lq \leq \delta k^{k}$$

V. When stop?
Consistency $f^{k} \rightarrow f^{*}$ Bayes optimal
 $n \rightarrow \infty$
 $k_{(n)} \infty$
If f not restricted $\Leftrightarrow f^{*} = \operatorname{argmin} Lq(f) \leftarrow Bayes qut for q$
 $\Rightarrow f = \operatorname{argmin} Lor(f) \leftarrow -n \longrightarrow Od$

1)

- We saw that L_φ is statistically motivated. Now we will see that it is computationally motivated as well.
- ▶ Recall: The "true" classification loss L_{01} is nonsmooth (has 0/no derivatives), non-convex. For training, one uses surrogate losses of the form $L_{\phi}(y, f) = \phi(yf)$.
- Want the following properties for

 - φ(z) is smooth (has continuous derivatives of any order if f has them); (this lets us use continuous optimization techniques to fit the classifier)
 - $\phi(z)$ is convex (this leads to global optimization, which has been recognized as beneficial in practice; it also allows to prove bounds, rates of convergence and so on)
 - $\phi(z)$ is monothone (decreasing) (thus, when z > 0, driving the margins to increase even if the classification is correct).

These properties are satisfied by $\phi(z) = e^{-z}$

Surrogate losses and boosting algorithms bad for outliers!



(sometimes good to have L(z) decrease for all z < 0, and sometimes bad – causes overfitting)

... and of boosting algorithms

- GENTLEADABOOST: approx Newton, $\phi = e^{-z}$
- ► LEAST-SQUARESBOOST: φ = (1 − z)² many operations in closed form
- ► LOGITBOOST $\phi = \ln(1 + e^{-z})$ slower (almost linear) decrease for $z \ll 0$
- ANYBOOST, GRADIENTBOOST work with any ϕ

Any Boost · Choose B

- · choose que descent method (eg Newton) GD
 - ⇒ My Boasting Alg



GradientBoost

GradientBoost Algorithm Given \mathcal{B} contains real-valued functions, loss L_{ϕ} , ϕ differentiable Input M, labeled training set \mathcal{D} **Initialize** $f^0(x) = \beta_0 = \operatorname{argmin}_{\beta \in \mathbb{R}} \hat{L}_{\phi}(\beta)$ steepest 1. compute $r^{i} = -y^{i}\phi'(y^{i}f(x^{i}))$ for i = 1 : n direction 2. fit $b^{k}(x)$ to data $\{(x^{i}, r^{i})\}$ for $k = 0, 1, 2, \dots, M - 1$ 3. find $\beta_k = \operatorname{argmin}_{\beta \in \mathbb{R}} L_{\phi}(f^k + \underline{\beta}b^k)$ (univariate optimization) \Leftrightarrow step size update $f^{k+1}(x) = f^k(x) + \beta^k b^k(x)$ Can be used for either classification or regression with 3 = CART with Works with any \$\phi\$ If \$\phi\$ convex stands If \u03c6 convex, step 3 is convex optimization (efficient)

- ▶ Proposed first as ANYBOOST, later specialized for $\vec{\mathcal{B}}$ =decision/regression trees, with other tweaks and new name GRADIENTBOOST
- When B = CART
 - step 3 optimizes over every leaf separately
 - depth of trees J represents maximum number of interactions in f; should not be too large (B must be weak)

Practical aspects

Overfitting in noise When the classes overlap much (many examples in \mathcal{D} hard/impossible to classify correctly) boosting algorithms tend to focus too much on the hard examples, at the expense of overall classification accuracy. The same happens for outliers. Observe also that the loss function(s) in the previous figure, which penalize more as the margin becomes more negative.

Choice of features Often times, the base class \mathcal{B} consists of function of the form $b(x) = x_j - a$, which perform a split on coordinate x_j at point $x_j = a$. They have the advantage that they can be learned and evaluated extremely fast. One can also augment the coordinate vector x with functions of the coordinates (e.g. $x \rightarrow [x_1 \dots x_d x_1 x_2 x_1 x_3 \dots])$ essentially creating a large set of features, which corresponds to finite but very large \mathcal{B} . In such a situation, the number of features d can easily be larger than M the number of b's in the final f. Thus, boosting will be implicitly performing a feature selection task.

When to stop boosting?

Loi (f, D') validation loss

The idea of Cross-Validation (CV) is to use an idependent sample from P_{XY} , denoted \mathcal{D}' and called the validation set to estimate the expected loss $L_{01}(f)$. When overfitting starts, $L_{01}(f^k)$ will start increasing with k. Boosting is the stopped at the value M that minimizes $\hat{L}_{01}(b^k; \mathcal{D}')$ (denoted L_{CV} below to simplify notation) ADABOOST WITH CROSS-VALIDATION

Given Training set D of size n, validation set D of size n', base classifier BInitialize

1. while L_{CV} decreases (but for at least 1 step)

- do a round of boosting on \mathcal{D}
- For i' = 1 : n' compute $f(x^{i'}) \leftarrow f(x^{i'}) + \beta^k b^k(x^{i'})$ ⇒ compute $f^{CV} = \frac{1}{2} \sum_{i=1}^{N} f_{i}(x^{i'}) + \beta^k b^k(x^{i'})$

• compute
$$L_{01}^{cv} = \frac{1}{n'} \sum_{i'} 1_{[y''f(x'') < 0]}$$

$$L_{oi}(\mathcal{D}', f^{k+1}) = \frac{1}{n'} \sum_{i \in \mathcal{D}} L_{oi}(f^{k+1}_{(k)}, y^{i})$$

How to choose a loss ϕ ?

Can we analyze which loss functions ϕ are "better"? Can we offer some guarantees in terms of generalization bounds? The answers are in [?] Bartlett, Jodan & McAuliffe,"Convextity, classification and risk bounds", 2005 (BJM).

We will restrict ourselves to convex, almost everywhere differentiable losss ϕ that are upper bounds of the 0-1 loss.

Let p = P[Y = 1|X], z = f(X). Then the expected loss of classification at X is

$$C_{p}(z) = p\phi(z) + (1-p)\phi(-z)$$
 (17)

and the optimal loss is

$$H(p) = \min b^{z} C_{p}(z) \text{ attained for } f^{*} = \frac{1}{2} \ln \frac{p}{1-p}$$
(18)

Let H^- denote the smallest loss for a misclassification

$$H^{-}(p) = \inf_{\operatorname{sgn} z = -\operatorname{sgn}(2p-1)} C_{p}(z)$$
(19)

Intuitively, we are minimizing ϕ instead of the "true" misclassification loss, and we want to measure how much we can be off when doing this. The following results say that we can bound the "true" loss $L_{01}(f)$ in terms of the ϕ -loss L_{ϕ} . We say ϕ is classification calibrated if $H^-(p) > H(p)$ for all $p \neq 1/2$. For ϕ convex, we have

that ϕ is classification calibrated iff ϕ differentiable at 0 and $\phi'(0) < 0$.

Proposition

(Theorem 4 in BJM) If ϕ is classification calibrated and convex, then for any classifier F

 $\psi(L_{01}(f) - L_{01}^*) \le L_{\phi}(f) - L_{\phi}^*$ (20)

where L^*_{ϕ} , L^*_{01} represent respectively the optimal ϕ -loss and optimal classification loss on the given data distribution and ψ is

$$\psi(\theta) = \phi(0) - H(\frac{1+\theta}{2})$$
 (21)

 $\label{eq:constraint} \begin{array}{|c|c|c|} \hline \text{Loss function } \phi(z) & \text{Transform function } \psi(\theta) \\ \hline \text{exponential: } e^{-z} & 1 - \sqrt{1 - \theta^2} \\ \text{truncated quadratic: } (\max(1-z,0))^2 & \theta^2 \\ \text{hinge: } \max(1-z,0) & |\theta| \\ \hline \text{In BJM there are also more general theorems that do not assume } \phi \text{ is convex.} \\ \hline \text{Furthermore, a convergence rate bound is given, which depends on: the noise in the labels, a complexity parameter of the function class \mathcal{B}, the curvature of ϕ. By optimizing this expression w.r.t to \mathcal{B} and ϕ one can theoretically choose the loss function and/or the base classifier.} \end{array}$

[?] proved that a large family of boosting algorithms is consistent. Below is an informal version of their main theorem.

Consistency is defined as $L_{01}(f^{t_n}) \to L_{01}^*$ for $n \to \infty$ and a t_n a certain sequence of stopping times that tends to infinity with n.

The conditions for consistency are as follows:

- 1. The function ϕ is convex, lower bounded (by 0) and calibrated as in section 2
- 2. The boosting step satisfies a weak leaning condition

$$\hat{\mathcal{L}}_{\phi}(f^{k}) \leq \gamma \inf_{b} \hat{\mathcal{L}}_{\phi}(f^{k-1} + \alpha b) + (1 - \gamma)\hat{\mathcal{L}}_{\phi}(f^{k-1})$$
(22)

- The set B is rich enough that the Bayes loss L^{*}_φ can be attained by convex combinations in B. There is a sequence f^k of k terms from B whose L_φ tends to L^{*}_φ
- 4. The empirical loss \hat{L}_{ϕ} converges to L_{ϕ} when $n \to \infty$ uniformly over all f which are t_n combinations
- 5. The empirical risks of \bar{f}^k converge, i.e max $\{0, |\hat{L}_{\phi}(\bar{f}^n) L_{\phi}(\bar{f}^n)|\} \rightarrow 0$, a.s. when $n \rightarrow \infty$
- 6. Algorithmic convergence $\max\{0, |\hat{L}_{\phi}(f^{t_n}) \hat{L}_{\phi}(\bar{f}^n)|\} \to 0$, a.s. when $n \to \infty$. In other words, the boosting algorithm produces an approximate minimizer of the L_{ϕ} risk if run for t_n iterations

$$t_n = n^{1-\varepsilon}$$
 for some $\varepsilon \in (0,1)$ (e.g. increases slowly with n).

M for $|\mathfrak{A}| = n$

7.

- Part II

...and a loss bound

The following result applies to any classifier, but it was developed in response to the idea that "boosting increases the margin" (which we now know is often, but not always true). It proves essentially that large margins counter overfitting. As with all worst-case bounds, the bounds are usually not realistic or practically applicable.

Proposition (to find citation)

Let \mathcal{F} be a model class of VC-dimension h, with $f(x) \in [-1, 1]$ for all x and for all $f \in \mathcal{F}$. Let $\delta > 0$ and $\theta \in (0, 1)$. Then, with probability w.p. $> 1 - \delta$ over training sets

$$L_{01}(f) \leq \frac{1}{n} |\{i \mid y^{i} f(x^{i}) \leq \theta\}| + \tilde{\mathcal{O}}\left(\sqrt{\frac{h}{n\theta^{2}}}\right)$$
(23)

for any $f \in \mathcal{F}$.

This theorem upper bounds the true loss $L_{01}(f)$ using the number of margin errors for an arbitrary margin θ . Note that for $\theta = 0$ a margin error is also a classification error, and for $\theta > 0$ the number of margin errors is greater or equal to that of classification errors. Hence, the first term of the bound increases with θ , while the second term decreases. So, if most examples can be classified with a large margin (not necessarily 1), then the bound of Theorem 4 can be tighter.