

1/30/23

Lecture 18

Lecture VII 1,2 Boosting + alg-adabooit gradient boosl Hw6 due

Wide NN Boosting

NTK during training – empirical evidence



Figure 1: Convergence of the NTK to a fixed limit Figure 2: Networks function f_{θ} near convergence for two widths n and two times t. for two widths n and 10th, 50th and 90th per-

wide
1. NN ~ GP = init Kernel = NTK

$$\approx$$
 after traing T
 $k(x_1x') = \nabla_{\theta}f_{\theta}(x) \nabla_{\phi}f_{\theta}(x')$
 $2. K^{L, q} = L / K^{L, q}$ smoother $f \approx 0$
 \Rightarrow BAD for regression

Wide and deep neural networks for classification – Basic quantities and assumptions

[Radhakrishnan, Belkin, Ulher, 2022]

- ▶ This paper studies the limits of wide neural networks $m_l \to \infty$ for all l = 1 : L when the depth $L \rightarrow \infty$
- It is already known that for regression $L \to \infty$ is NOT OPTIMAL
- Since the NTK depends only of the activation function ϕ , the limit shall only depend on ϕ as well. of K
- ln particular, the limit depends on ϕ only through the following

$A = E[\phi(Z)]$	when $z \sim N(0,1)$
$A' = E[\phi'(Z)]$	when $z \sim N(0,1)$
$B = E[(\phi'(Z))^2]$	when $z \sim N(0,1)$

- Classifier $f(x) = \lim_{L \to \infty} \sup_{x \to \infty} YG^{-1}\kappa^{L}(X, x)$ with $G = [\kappa^{L}(x^{i}, x^{j})]_{j,j=1:n}$.
- Additional assumptions
 - Data $X \subseteq S^d_+$, vectors of norm 1 with all entries ≥ 0 .
 - Simplifying assumptions on NTK parameters (e.g. $\sigma_w = \sigma_b = 1$)

Case $A \neq 0$: Networks implement majority vote

Theorem (Proposition 1 in [Radhakrishnan, Belkin, Ulher, 2022]) If there is a function $0 < c(L) < \infty$ so that

$$\lim_{L \to \infty} \frac{\kappa^{L}(x, x')}{c(L)} = c_1 > 0 \text{ for any } x \neq x', \text{ and } \lim_{L \to \infty} \frac{\kappa^{L}(x, x)}{c(L)} \neq c_1,$$
(27)

then

$$\lim_{L \to \infty} f(x) = \operatorname{sgn} \sum_{i=1}^{n} y^{i} \qquad MAJORITY \ CLASSIFIER$$
(28)

What
$$\phi$$
's satisfy theorem? ReLU, all ϕ with $B \neq 1$.

Case A = A' = 0: Networks implement 1-nearest neighbor

Theorem (Theorem 3 in [Radhakrishnan, Belkin, Ulher, 2022]) Given x, assume w.l.o.g. that $x^T x^1 = \max_{i=1:n} x^T x^i$.

$$\lim_{L\to\infty}\frac{\kappa^L(x,x^i)}{\kappa^L(x,x^1)} = 0.$$
 (29)

and

$$\lim_{x \to \infty} f(x) = \operatorname{sgn} y^1 \qquad 1-nn \tag{30}$$

Case A = 0, $A' \neq 0$: Networks implement singular kernel classifier

Theorem (Theorem 1 in [Radhakrishnan, Belkin, Ulher, 2022])

$$\lim_{N \to \infty} \frac{\kappa^{L}(x, x')}{(A')^{2L}(L+1)} = \frac{R(\|x - x'\|)}{\|x - x'\|_{\alpha}^{\alpha}},$$
(31)

with $\alpha = -4 \frac{\log A'}{\log B'}$ and $R() \ge 0$, bounded, and $R(u) > \delta$ around 0.

- if $\alpha > 0$, $\frac{R(||x-x'||)}{||x-x'||^{\alpha}}$ is singular kernel
- Computationally not a problem: if data $x^{1:n}$ distinct, G_0 is well defined
- If $x = x^i$, set $f(x) = y^i$.

Optimality of singular kernel classifier

Theorem (Theorem 2 in [Radhakrishnan, Belkin, Ulher, 2022]) If A = 0, $A' \neq 0$ and $\alpha = -d$ then $\lim_{L\to\infty} f(x)$ is Bayes-optimal.

What activations \u03c6 satisfy this theorem?

$$\phi^{\mathrm{opt}}(z) = rac{1}{2^{d/4}} rac{z^3 - 3z}{\sqrt{6}} + \sqrt{1 - 2^{1 - d/2}} rac{z^2 - 1}{\sqrt{2}} + rac{1}{2^{d/4}} z \quad ext{ for } d \ge 2.$$
 (32)



Optimal singular kernels for d = 4, 8, 16, 32



phi

phi'

4 5

phi

0

----- phi

3

Summary



b when A ≠ 0, lim_{L→∞} κ^L(x, x') = 0 for x ≠ x', and f(x) = 0 is vanishingly small (useless for regression), but sgnf(x) can be optimal for classification
c Singular kernel α > d,α < d, majority vote kernel, and 1-nn kernel

Limits of some activation functions

 $\begin{array}{l} \phi^{\rm opt} \quad {\rm Bayes\ classifier} \\ \hline {\rm ReLU} \quad {\rm majority\ vote} \\ {\rm sigmoid} \quad \frac{1}{1+e^{-z}}-\frac{1}{2} \ 1{\rm -nearest\ neighbor} \end{array}$

Lecture VII: Combining predictors - Part I

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

December, 2023



There is more than one way to average predictors Bayesian averaging Bagging → Random Forest (Reduce Var) Stacking Mixtures of Experts Forward Fitting and Backfitting

Reducing bias: Boosting - Not averaging

AdaBoost Properties on the training set

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] Boosting, 3.2.4 Bayesian model averaging, 16.6.3, 16.6.1 Ensembles of predictors, Bach Ch.:

There is more than one way to average predictors

Classification will be the running exaple here, but most results hold for other prediction problems as well.

Denote $\mathcal{B} = \{b\}$ a base classifier family Averaging:

$$f(x) = \sum_{k=1}^{M} \beta^k b^k(x) \tag{1}$$

f is real-valued even if the b^k 's are ± 1 valued Why average predictors?

- to reduce bias
- to compensate for local optima (a form of bias)
- to reduce variance
- if b₁, b₂,...b_M make independent errors, averaging reduces expected error (loss). We say that b₁ and b₂ make independent errors ⇔
 P(b₁ wrong | x) = P(b₁ wrong | x, b₂ wrong)
- because the b^k functions are given: real world domain experts (weather prediction), a set of black-box classifiers (from a software package), a set of [expert designed] features (speech recognition, car/human recognition) each of them weakly informative of y
- ▶ because *B* is a set of (simple) "basis functions" and we need a more complex predictor in our task

Averaging is not always the same thing

Depending how we choose \mathcal{B} , $b^1, b^2, \dots b^M$ and $\beta^1, \beta^2, \dots \beta^M$, we can obtain very different effects.

We will examine

- Bayesian averaging (briefly)
- Mixtures of experts (briefly)
- Bagging (briefly)
- Backfitting (briefly)
- Stacking (see Murphy)
- Boosting

Reducing bias: Boosting

Task = classification

Base classifier family \mathcal{B} has large bias (e.g. linear classifier, decision stumps) but learning always produces b that is better (on the training set) than random guessing. Preconditions for boosting

1. Learning algorithm accepts weighted data sets. Training minimizes

$$\hat{L}_{01}^{\sf w}(b) \,=\, \sum_{i=1}^n w_i L_{01}(y^i, b({\sf x}^i)) \quad {
m with} \;\; \sum_{i=1}^n w_i \,=\, 1.$$

2. \mathcal{B} is a weak classifier family. For any \mathcal{D} and any weights $w_{1:n}$ there can be found $b \in \mathcal{B}$ such that the training error of b on \mathcal{D} is bounded below one half.

$$0 < \hat{L}_{01}^{w}(b) \leq \delta < \frac{1}{2}$$

 $3 = \frac{1}{2}$
 $3 = \frac{1}{2}$ data consider

B=16()3 Weak clamifier

First Idea of boosting: train a classifier b^1 on \mathcal{D} , then train a b^2 to correct the errors of b^1 , then b^3 to correct the errors of b^2 , etc. TRAIN: Text $f(x) = \sum_{k=1}^{\infty} b^k b(x)$ k = 1

M~n& K<1

Example



Marina Meila: Lecture VII: Combini

- Part I

AdaBoost Algorithm



Remarks

- If b(x) ∈ {±1} then yⁱb(xⁱ) ∈ {±1}, and ^{1-yⁱb(xⁱ)}/₂ = 1 if an error occurs and 0 otherwise. Thus, ε^k in step 2 adds up the weights of the errors. If b(x) ∈ [-1, 1] then the errors contribute different amounts to the loss depending on their margin.
- 2. In both cases, $\varepsilon^k \in [0, \delta], \, \delta < 0.5$ by the weak learner property of $\mathcal B$
- 3. $\beta^k > 0$ whenever $\varepsilon^k < 1/2$.
- 4. If $b \in \{\pm 1\}$, then step 4 can be written equivalently (up to a multiplicative constant)

$$w_i^{k+1} = \begin{cases} \frac{1}{Z^k} w_i^k & \text{if } b^k(x_i) = y_i \\ \frac{1}{Z^k} w_i^k e^{2\beta^k} & \text{if } b^k(x_i) \neq y_i \end{cases}$$
(7)

This form corresponds to the DISCRETEADABOOST algorithm, the first AdaBoost algorithm published, which assumed $b(x) \in \{\pm 1\}$. As we shall see later, modern boosting algorithms dispense with the assumption $b \in [-1, 1]$ too.

An interpretation of the weights



- weight of example i at step k is proportional to e^{-y_ir^{k-1}(x_i)} the exponential of its negative margin
- Examples that have been hard to classify get exponentially high weights.
- Examples that are classified with high margins get vanishingly small weights.

$$y^{i}f(x^{i}) = margin$$



Hence, the r.h.s of (10) is the average over the data set of the **exponential loss** L_{ϕ} . The function ϕ decreases with the margin, thus decreasing \hat{L}_{ϕ} will produce a better classifier (on the training set). In this sense, L_{ϕ} is an alternative loss function for classification.

Lo < Xth

4<1



\hat{L}_{ϕ} decreases exponentially with M

For simplicity, we show this in the special case $b(x) \in \{\pm 1\}$ for all $b \in \mathcal{B}$.

$$Z^{k} = \sum_{i=1}^{n} w_{i}^{k} e^{-\beta^{k}(y_{i}b^{k}(x_{i}))}$$
(13)

$$= e^{\beta^{k}} \underbrace{\sum_{i=err} w_{i}^{k} + e^{-\beta^{k}}}_{e^{k}} \underbrace{\sum_{i=corr} w_{i}^{k}}_{1-e^{k}}$$
(14)

$$= e^{\beta^{k}}\varepsilon^{k} + (1 - \varepsilon^{k})e^{-\beta^{k}}$$
(15)

$$= \sqrt{\frac{1-\varepsilon^{k}}{\varepsilon^{k}}}\varepsilon^{k} + \sqrt{\frac{\varepsilon^{k}}{1-\varepsilon^{k}}}(1-\varepsilon^{k}) = 2\sqrt{(1-\varepsilon^{k})}\varepsilon^{k} \leq \gamma$$
(16)

where $\gamma < 1$ depends on δ the maximum error. It follows that

$$\hat{L}_{\phi}(f^k) = \prod_{k=1}^k Z^{k'} \leq \gamma^k \tag{17}$$

The training set error \hat{L}_{01} decreases exponentially with M

Note that $\phi(z) \ge 1_{z<0}$ for all z (see also figure ??). Therefore

$$\hat{\mathcal{L}}(f^k) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[y_i f^k(x_i) < 0]}$$
(18)

$$\leq \frac{1}{n} \sum_{i=1}^{n} e^{-y_i f^k(x_i)} = \hat{L}_{\phi}(f^k) \leq \gamma^k$$
 (19)

In other words, the training error $\hat{L}(f^k)$ is bounded by a decaying exponential. Moreover, since $\hat{L}(f^k) \in \{0, 1/n, 2/n, \dots 1\}$, it follows that after a finite number of steps, when $\gamma^{k^0} < 1/n$, the training error will become 0 and the training data will be perfectly classified!

The test set error and overfitting

- ▶ Do NOT take $M = k^0$. The number of steps M for good generalization error is often much larger than k^0 (and sometimes smaller).
- Below is a typical plot of \hat{L} and L (which can be estimated from an independent sample) vs the number of boosting iterations.



Lecture VII: Combining predictors - Part II

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

November, 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

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.:

q = surrogate Loss

Meila: Lecture VII: Combining

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)

$$f = \begin{bmatrix} f(x^1) \\ f(x^2) \\ \cdots \\ f(x^n) \end{bmatrix} \quad b = \begin{bmatrix} b(x^1) \\ b(x^2) \\ \cdots \\ b(x^n) \end{bmatrix}$$
(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 $\hat{L}_{\phi}(f) - \hat{L}_{\phi}(f + \beta b)$.



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

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

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

want $f^* \in \mathbb{R}^n$ to minimize \hat{L}_{φ} decrease \hat{L}_{φ} iteratively $f_{(x)} \stackrel{?}{\underset{k}{\to}} \stackrel{restricted !!}{f^{k+1}} = f^k + pe_k \stackrel{le}{\underset{q}{\to}} \stackrel{restricted !!}{f^{k+1}} = f^k + pe_k \stackrel{le}{\underset{q}{\to}} \stackrel{restricted !!}{f^{k+1}}$

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} .

The line minimization

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.

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

 β is the (unique) root of

$$\sum_{i} w_i y_i b(x^i) e^{-\beta y_i b(x^i)} = 0$$
⁽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.

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}
```

Marina Meila: Lecture VII: Combining predictors – Part II