

Boosting

Weikai Xu

School of Information Science and Technology
ShanghaiTech University

Mar 15th, 2020

Outline

- 1 Basic Algorithm and Core Theory
- 2 Other Ways of Understanding AdaBoost

Outline

- 1 Basic Algorithm and Core Theory
- 2 Other Ways of Understanding AdaBoost

A Formal Description of Boosting

- ① given training set $(x_1, y_1), \dots, (x_m, y_m)$
- ② $y_i \in \{-1, +1\}$ correct label of instance $x_i \in X$
- ③ for $t = 1, \dots, T$:
 - ① construct distribution D_t on $\{1, \dots, m\}$
 - ② find weak classifier

$$h_t : X \rightarrow \{-1, +1\}$$

with small error ϵ_t on D_t :

$$\epsilon_t = \Pr_{i \sim D_t}[h_t(x_i) \neq y_i]$$

- ④ output final classifier H_{final}

AdaBoost

- ① constructing D_t
 - ① $D_1(i) = 1/m$
 - ② given D_t and h_t :

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{cases} \\ &= \frac{D_t(i)}{Z_t} \exp(-\alpha_t y_i h_t(x_i)) \end{aligned}$$

Where $Z_t =$ normalization constant, $\alpha_t = \frac{1}{2} \ln(\frac{1-\epsilon_t}{\epsilon_t}) > 0$

- ② Final Classifier:

$$H_{final}(x) = \text{sign}\left(\sum_t \alpha_t h_t(x)\right)$$

Analyzing the training error

1 Theorem

- 1 write ϵ_t as $1/2 - \gamma_t$
- 2 then

$$\begin{aligned}\text{training error}(H_{\text{final}}) &\leq \prod_t [2\sqrt{\epsilon_t(1 - \epsilon_t)}] \\ &= \prod_t \sqrt{1 - 4\gamma_t^2} \\ &\leq \exp(-2 \sum_t \gamma_t^2)\end{aligned}$$

2 so: if $\forall t : \gamma_t \geq \gamma > 0$, then $\text{training error}(H_{\text{final}}) \leq e^{-2\gamma^2 T}$

3 AdaBoost is adaptive:

- 1 does not need to know γ or T a prior
- 2 can exploit $\gamma \gg \gamma$

Proof

- ① let $f(x) = \sum_t \alpha_t h_t(x) \Rightarrow H_{final}(x) = \text{sign}(f(x))$.
- ② Step 1: unwrapping recurrence:

$$\begin{aligned} D_{final}(i) &= \frac{1}{m} \frac{\exp(-y_i \sum_t \alpha_t h_t(x_i))}{\prod_t Z_t} \\ &= \frac{1}{m} \frac{\exp(-y_i f(x_i))}{\prod_t Z_t} \end{aligned}$$

Proof

① Step 2: training error(H_{final}) $\leq \prod_t Z_t$

② Proof:

$$\begin{aligned}\text{training error}(H_{final}) &= \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i \neq H_{final}(x_i) \\ 0 & \text{else} \end{cases} \\ &= \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i f(x_i) \leq 0 \\ 0 & \text{else} \end{cases} \\ &\leq \frac{1}{m} \sum_i \exp(-y_i f(x_i)) \\ &= \sum_i D_{final}(i) \prod_t Z_t \\ &= \prod_t Z_t\end{aligned}$$

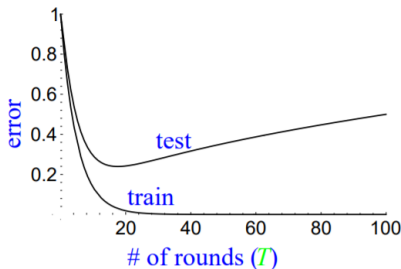
Proof

① Step 3: $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$

② Proof:

$$\begin{aligned} Z_t &= \sum_i D_t(i) \exp(-\alpha_t y_i h_t(x_i)) \\ &= \sum_{i: y_i \neq h_t(x_t)} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t(x_t)} D_t(i) e^{-\alpha_t} \\ &= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t} \\ &= 2\sqrt{\epsilon_t(1 - \epsilon_t)} \end{aligned}$$

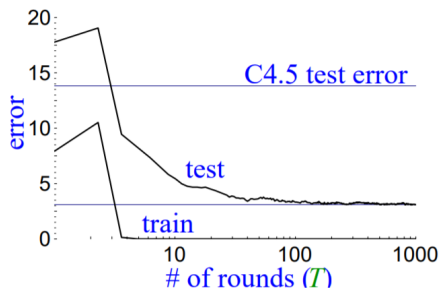
How Will Test Error Behave?(Guess)



expect:

- 1 training error to continue to drop (or reach zero)
- 2 test error to increase when H_{final} becomes “too complex”

Actual Typical Run



- 1 test error does not increase, even after 1000 rounds
- 2 test error continues to drop even after training error is zero!

A Better Story: The Margins Explanation

- ① key idea:
 - ① training error only measures whether classifications are right or wrong.
 - ② should also consider confidence of classifications
- ② recall: H_{final} is weighted majority vote of weak classifiers
- ③ measure confidence by margin = strength of the vote
= (fraction voting correctly) - (fraction voting incorrectly)

Theoretical Evidence: Analyzing Boosting Using Margins

- ① Theorem: large margins \Rightarrow better bound on generalization error (independent of number of rounds)
 - ① proof idea: if all margins are large, then can approximate final classifier by a much smaller classifier
- ② Theorem: boosting tends to increase margins of training examples (given weak learning assumption)
 - ① proof idea: similar to training error proof
- ③ so:
although final classifier is getting larger,
margins are likely to be increasing,
so final classifier actually getting close to a simpler classifier,
driving down the test error

Outline

- 1 Basic Algorithm and Core Theory
- 2 Other Ways of Understanding AdaBoost

Game Theory

- ① game defined by matrix $M = \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}$
- ② row player chooses row i
- ③ column player chooses column j (simultaneously)
- ④ row player's goal: minimize loss $M(i, j)$
- ⑤ usually allow randomized play: players choose distributions P and Q over rows and columns
- ⑥ learner's (expected) loss

$$\begin{aligned} &= \sum_{i,j} P(i)M(i,j)Q(j) \\ &= P^T M Q = M(P, Q) \end{aligned}$$

The Minmax Theorem

- ① von Neumann's minmax theorem:

$$\begin{aligned}\min_P \max_Q M(P, Q) &= \max_Q \min_P M(P, Q) \\ &= v \\ &= \text{"value" of game M}\end{aligned}$$

- ② in words:

- ① $v = \min \max$ means:

- ① row player has strategy P^*
such that \forall column strategy Q , loss $M(P^*, Q) \leq v$

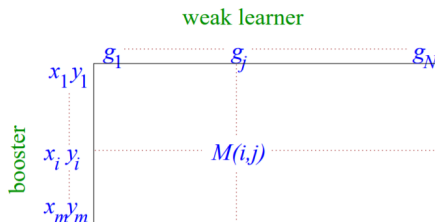
- ② $v = \max \min$ means:

- ① column player has strategy Q^*
such that \forall row strategy P , loss $M(P, Q^*) \geq v$

The Boosting Game

- ① let $\{g_1, \dots, g_N\}$ = space of all weak classifiers
- ② row player \rightarrow sampler
- ③ column player \rightarrow weak learner
- ④ matrix M :
 - ① row \rightarrow example (x_i, y_i)
 - ② column \rightarrow weak classifier g_i
 - ③

$$M(i, j) = \begin{cases} 1 & \text{if } y_i = g_j(x) \\ 0 & \text{else} \end{cases}$$



Boosting and the Minmax Theorem

① if:

① \forall distributions over examples
 $\exists h$ with accuracy $\geq \frac{1}{2} + \gamma$

② then:

① $\min_P \max_j M(P, j) \geq \frac{1}{2} + \gamma$

③ by minmax theorem:

① $\max_Q \min_i M(i, Q) \geq \frac{1}{2} + \gamma > \frac{1}{2}$

④ which means:

① \exists weighted majority of classifiers which correctly classifies all examples with positive margin(2γ)

⑤ optimal margin \leftrightarrow “value” of game

AdaBoost and Game Theory

- ① AdaBoost is special case of general algorithm for solving games through repeated play
- ② can show
 - ① distribution over examples converges to (approximate) minmax strategy for boosting game
 - ② weights on weak classifiers converge to (approximate) maxmin strategy
- ③ different instantiation of game-playing algorithm gives on-line learning algorithms (such as weighted majority algorithm)

Boost and Exponential Loss

- ① many (most?) learning algorithms minimize a “loss” function
 - ① e.g. least squares regression
- ② training error proof shows AdaBoost actually minimizes

$$\prod_t Z_t = \frac{1}{m} \sum_i \exp(-y_i f(x_i))$$

where $f(x) = \sum_t \alpha_t h_t(x)$

- ③ on each round, AdaBoost greedily chooses α_t and h_t to minimize loss.
- ④ we can prove that AdaBoost provably minimizes exponential loss.

Coordinate Descent

- ① $\{g_1, \dots, g_N\}$ = space of all weak classifiers
- ② want to find $\lambda_1, \dots, \lambda_N$ to minimize

$$L(\lambda_1, \dots, \lambda_N) = \sum_i \exp(-y_i \sum_j \lambda_j g_j(x_i))$$

- ③ AdaBoost is actually doing coordinate descent on this optimization problem:
 - ① initially, all $\lambda_j = 0$
 - ② each round: choose one coordinate λ_j (corresponding to h_t) and update (increment by α_t)
 - ③ choose update causing biggest decrease in loss
- ④ powerful technique for minimizing over huge space of functions

Functional Gradient Descent

- ① want to minimize

$$L(f) = L(f(x_1), \dots, f(x_m)) = \sum_i \exp(-y_i f(x_i))$$

- ② say have current estimate f and want to improve
- ③ to do gradient descent, would like update

$$f \leftarrow f - \alpha \nabla_f L(f)$$

- ④ but update restricted in class of weak classifiers

$$f \leftarrow f + \alpha h_t$$

- ⑤ so choose h_t “closest” to $-\nabla_f L(f)$
- ⑥ equivalent to AdaBoost