
A Stochastic Approximation View of Boosting

C. Andy Tsao Yuan-chin Ivan Chang*

Department of Applied Math, National Dong Hwa University

Institute of Statistical Science, Academia Sinica*

Outline

Introduction



FHT's Interpretation



Stochastic Approximation Viewpoint

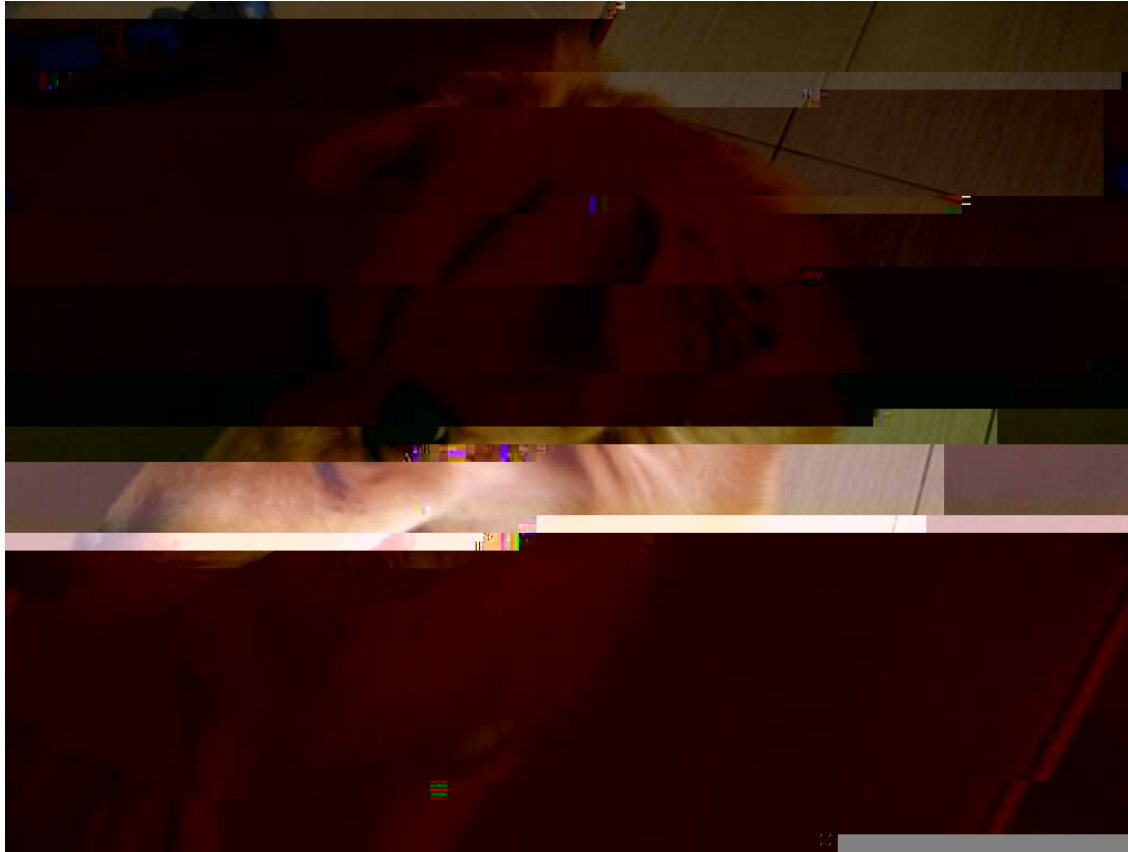


Results



Conclusion and Discussion

Murray at work



A not-so-Lean not-so-Mean Stat Machine

Intro: Supervised Learning

■ Training data $(\mathbf{x}_i; y_i)_{i=1}^N$; $\mathbf{x} \in \mathbf{X}; y \in \mathbf{Y} = \{f, \dots, g\}$:

Find *Machine (Classifier)* $H: \mathbf{X} \rightarrow \mathbf{Y}$

🎯 Testing Data $(\mathbf{x}'_j; y'_j)_{j=1}^M$.

■ Training Error


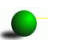
$$TE = \frac{1}{N} \sum_i \mathbf{1}_{[y_i \neq H(\mathbf{x}_i)]}$$

🎯 Generalization Error

$$GE = \frac{1}{M} \sum_j \mathbf{1}_{[y'_j \neq H(\mathbf{x}'_j)]}$$

Intro: Boosting

 *Weak (base) learner*

- Sequentially applying it to reweighted version of the training data
 -  Higher weights on the previous misclassified data
 -  Boosting iteration: T
- Weighted majority vote

Schapire (1990), Freund and Schapire (1997), Friedman, Hastie and Tibshirani (2000).

Brieman (2004), Jiang (2004), etc.

Intro: Discrete AdaBoost

1. Start with weights $D_t(i) = 1/N; i = 1$ to N :
2. Repeat for $t = 1$ to T
 - Obtain $h_t(x)$ from weak learner h using weighted training data wrt D_t
 - Compute $\epsilon_t = E_{D_t} 1_{[y \neq h_t(x)]}$; $\alpha_t = \log \frac{1-\epsilon_t}{\epsilon_t}$:
 - Update $i = 1$ to N ,

$$D_{t+1}(i) = \frac{1}{Z_t} D_t(i) \exp[-\alpha_t 1_{[y_i \neq h_t(x_i)]}];$$

where Z_t is the normalizer.

3. Output the classifier $\text{sgn}\left[\sum_{t=1}^T \alpha_t h_t(x)\right]$

Intro: Theories and Observation-1

Under *weak base hypothesis assumption*, TE goes to zero exponentially as $T \rightarrow \infty$

PAC bound for GE: Let d : VC dim of the weak base hypothesis space.

$$GE = TE + O\left(\sqrt{\frac{d}{N}}\right)$$

Margin. Schapire, et al (1998). $m(x_i; y_i) = \frac{y_i - h_t(x_i)}{a_t}$.

$GE = \frac{1}{\hat{P}[m(x; y) \geq \theta]} + O\left(\sqrt{\frac{d}{N\theta^2}}\right)$ for any $\theta > 0$ with high probability.

Intro: Theories and

FHT's Interpretation

Friedman, Hastie and Tibishirani (2000).

Result 1 *The Discrete AdaBoost (population version) builds an additive logistic regression model via Newton-like updates for minimizing $\mathbf{E}(e^{-Y F(X)})$:*

SWOT

Strength

Easier for statisticians than MLs6 -0 1190.R143 Do Q 0 0

• For $c > 0$, $f(x) = \text{sgn}(E_w(y|x))$ minimizes (1),

• Given $f(x) = \pm 1$,

$$c = \arg \min_c J(F + cf) = \arg \min_c E_w e^{-cy - (f)x}$$

$$= \frac{1}{2} \log\left(\frac{1 - \text{err}}{\text{err}}\right); \quad \text{err} = E_w \mathbf{1}_{[y - (f)x < 0]} \quad 0.250010 \quad ($$

Alternative expression

$$\mathbf{F}_{t+1}(\mathbf{x}) = \mathbf{F}_t(\mathbf{x}) + \eta_t \operatorname{sgn}(\mathbf{E}_{w_t}(\mathbf{Y} \mathbf{j} \mathbf{x}))$$

where

$$\eta_t = \log\left(\frac{1}{t}\right)$$



Remarks on FHT

● $f(x)$ is fixed only when x is given.

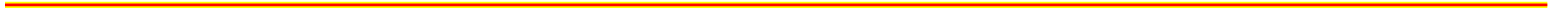
🇺🇸 $\mathbf{E}_{X,Y}(e^{-YF(X)}) = \mathbf{E}_X \left[\mathbf{E}_{Y|X}(e^{-YF(X)}) \right]$

🇺🇸 Convergence of the Newton-like updates via minimizing

$$(F) = \mathbf{E}_{Y|X} \mathbf{L}(Y; F(x))?$$

■ $\mathbf{L}(Y; F) = e^{-YF}$

Friedman (2001).



SA: Robbins-Monro Algorithm

Goal: Find

$$F_* = \arg \min_F (F) \quad [= E_{Y|x} L(Y; F(x))]$$

$$F_* ! \quad ' (F_*) = E_{Y|x} L' (Y; F_*) = 0$$

where

$$L' (Y; F) = \frac{\partial}{\partial F} L(Y; F)$$

Duflo (1997).

Robbins and Monro (1951), Kiefer and Wolfowitz (1952).

RM-D Algorithm

1 Choose F_0 arbitrarily.

2 Iterate $t = 0; 1; \dots$

$$\mathbf{g}_{t+1} = \mathbf{E} L'(\mathbf{Y}_{t+1}; \mathbf{F}_t)$$

$$\mathbf{F}_{t+1} = \mathbf{F}_t + \eta \mathbf{g}_{t+1}; \quad \eta = 0$$

RM-D with Exponential Criterion

- Choose F_0 arbitrarily.

- Iterate $t = 0; 1; \dots; Y_t \stackrel{iid}{\sim} Y | \mathbf{x}$

$$g_{t+1} = \mathbf{E} Y_{t+1} \exp(Y_{t+1} F_t)$$

$$F_{t+1} = F_t + g_{t+1}$$

$$= F_t + \mathbf{E}_{w_t} (Y_{t+1} \mathbf{j} \mathbf{x}); \quad 0$$

where $w_t(\mathbf{x}; y) = \exp(y F_t(\mathbf{x}))$:

RM-D convergence

Proposition 1 *If f is a conti. real ftn such that $(F_*) =$ and for all F*

$(f(F) - f(F_))(F - F_*) > 0$ and $\|f(F)\| \leq K(1 + \|F\|)$; for some K :*

Suppose $\{f_t\}$; $\{g_t\}$ are seq's of reals and $f_t \geq 0$. Define

$$F_{t+1} = F_t + g_t (f(F_t) - f_*):$$

If $f_t \geq 0$ s.t.

$$\sum_t f_t = 1; \quad \sum_t g_t < 1;$$

then $F_t \rightarrow F_$ for any initial F_0*

 $(F) = E_{Y|F} L'(Y; F).$

 t “step-size”.

 Common choice

Table 1: Testing errors of SABoost, AdaBoost with Decision Stumps and RBF-Network as weak base learners.

METHOD	SABOOST	ADABOOST	ADABOOST
DATA SET	$\gamma = 1.0$	DS	RBF-N
IONOSPHERE	10.13 ± 4.77	13.98 ± 5.46	—
BUPA	28.92 ± 7.00	32.56 ± 8.10	—
PIMA	23.90 ± 5.51	25.60 ± 4.52	—
WDBC	4.18 ± 2.66	3.01 ± 2.25	—
DIABETES	25.06 ± 1.90	25.52 ± 1.99	26.47 ± 2.29
GERMAN	26.50 ± 2.33	26.81 ± 2.47	27.45 ± 2.50
HEART	18.16 ± 3.23	20.01 ± 3.76	20.29 ± 3.44
SPLICE	13.06 ± 0.83	14.22 ± 1.09	10.14 ± 0.51
TWONORM	9.59 ± 1.18	5.85 ± 0.57	3.03 ± 0.28
WAVEFORM	14.73 ± 2.43	13.18 ± 5.69	10.84 ± 0.58
BREAST	29.87 ± 4.98	30.03 ± 5.05	30.36 ± 4.73

There are 50 replications and 1000 iterations for Ionosphere, Bupa liver-disorder, Pima Indian-Diabetes, and WDBC breast-cancer using random sampling. For Breast, Diabetes, German, Heart, Splice Twonorm, and Waveform, we follow the original partitions in IDA Benchmark Repository. Thus, there are 100 replications and 200 iterations for each of them except Splice. There are only 20 combinations in Splice.

Table 2: SA Boosting with different step sizes.

DATA SET	$\gamma = 0.5$	$\gamma = 1.5$	HYBRID ($\gamma = 0.5$)
IONOSPHERE	13.03 ± 1.53	10.97 ± 4.86	13.80 ± 4.76
BUPA	29.47 ± 7.12	29.04 ± 6.95	30.58 ± 8.60
PIMA	24.61 ± 4.77	24.24 ± 5.47	24.81 ± 4.40
WDBC	4.77 ± 2.89	3.97 ± 2.59	2.80 ± 2.09
DIABETIS	25.14 ± 1.65	25.38 ± 1.95	25.19 ± 1.84
GERMAN	25.94 ± 2.31	26.73 ± 2.43	26.14 ± 2.40
HEART	17.58 ± 3.37	18.64 ± 3.32	18.76 ± 3.19
SPLICE	13.04 ± 0.83	13.57 ± 0.89	13.67 ± 0.85
TWONORM	6.06 ± 0.42	6.96 ± 0.41	5.64 ± 0.47
WAVEFORM	15.73 ± 0.89	13.86 ± 2.43	13.04 ± 0.56
BREAST	29.94 ± 4.93	30.07 ± 5.11	30.03 ± 5.05

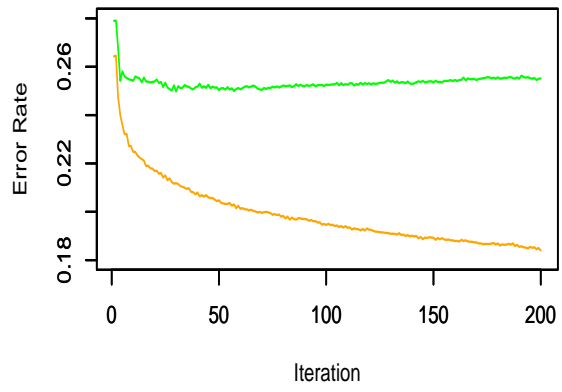


Table 4: Error rates (in %) of AdaBoost, SABoost ($\gamma = 0.5, 1$) with different base learner (SVM with Linear Kernel Polynomial Kernel with degree=3 and coef=1 and RBF kernel default parameter). Testing error(GE) and Training Error (TE).

BASE LEARNER: SVM-LINEAR		
METHODS	GE (TE)	$ TE - GE $
ADABOOST	34.7 (22.1)	0.122 ± 0.20
SABOOST $\gamma = 0.5$	31.8 (25.9)	0.059 ± 0.06
SABOOST $\gamma = 1$	35.5 (26.3)	0.091 ± 0.09
SVM-LINEAR GE	46.64 ± 4.84	
BASE LEARNER: SVM-POLYNOMIAL		
ADABOOST	23.06 (18.29)	3.31 ± 0.86
SABOOST $\gamma = 0.5$	23.81 (21.99)	2.06 ± 0.39
SABOOST $\gamma = 1$	22.96 (21.35)	1.74 ± 0.39
SVM-POLY GE	22.49 ± 1.86	
BASE LEARNER: SVM-RBF		
ADABOOST	23.63 (2.41)	3.34 ± 0.76

SABOOST $\gamma =$

SABOOST $\gamma =$.

Conclusion/Discussion

- SA provides new insight/tools to boosting-like algorithms.
 - Understandings
 - WBHA might not help.
 - “Theoretical” explanation for shrinkage
 - “Optimal” w
-

