

---

# Bayes Consistency of Boosting: Population versus Sample

W. Drago Chen and C. Andy Tsao

Department of Applied Math, National Dong Hwa University

Conference on MDT, Statistical Inference and Applications

In Honor of Deng-Yuan Huang

# Outline

---

Introduction

Convergence and ConsCod

n

---

# Intro: Supervised Learning

---

Training data  $(x_i, y_i)$ ,  $\mathcal{X}$   $\mathcal{Y} = \{\pm 1\}$ ;

Testing Data  $(x_i, y_i) \rightsquigarrow (X_i, Y_i) \sim_{i.i.d.} X, Y$

Find **Machine(Classifier)**  $\in$

$: \mathcal{X} \rightarrow \mathcal{Y}$

Training Error

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

Testing Error

$$= \mathbb{E} \left[ \mathbb{1}_{f(X) \neq Y} \right]$$

Goal

$$\min_{f \in \mathcal{H}} \mathbb{E} \left[ \mathbb{1}_{f(X) \neq Y} \right]$$

Challenge

$$\min_{f \in \mathcal{H}} \mathbb{E} \left[ \mathbb{1}_{f(X) \neq Y} \right]$$

Goal

$$\min_{f \in \mathcal{H}} \mathbb{E} \left[ \mathbb{1}_{f(X) \neq Y} \right]$$

# Intro: Boosting

---

Ensemble classifiers.

Weak (base) learner

Sequentially applying it to reweighted version of the training data

- Higher weights on the previous misclassified cases
- Boosting iteration:

Weighted majority vote

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

Breiman (2004), Jiang (2004), Meir and Rätsch (2003)

---

# Intro: Discrete AdaBoost

---

1. Start with weights  $(w_i) = 1$   $i = 1$  to  $n$

2. Repeat for  $t = 1$  to  $T$

Obtain  $(h_t)$  from weak learner using weighted training data wrt

Compute  $\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$

---

# Convergence and Consistency

---

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f(x) dP_n(x) = \int_{\mathcal{X}} f(x) dP(x)$$

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f(x) dP_n(x) = \int_{\mathcal{X}} f(x) dP(x)$$

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f(x) dP_n(x) = \int_{\mathcal{X}} f(x) dP(x) \text{ where } F_{P_n} \rightarrow F_P \text{ and } \int_{\mathcal{X}} f(x) dP(x) = \int_{\mathcal{X}} f(x) dP(x)$$



# Intro: Theories

---

Bayes consistent

(Population Version, Breiman (2004)).

Process Consistent

(Sample Version, Jiang (2004))

Regularization needed, say, early stopping, restriction on base learners, particularly for noise data.

On the other hand

Relatively immune to overfitting in practical apps

Mease and Wyner (2007, JMLR). Evidence Contradictory to Statistical View.

- Relatively immune to overfitting (Convergence)
- No regularization needed for some noisy data sets

# “Statistical View”: FHT’s Insights

---

Friedman, Hastie and Tibishirani (2000).

The Discrete AdaBoost (population version) builds an additive logistic regression model via Newton-like updates for minimizing  $\sum_{i=1}^n \ell_{\eta}(-Y_i \eta(X_i))$

Exponential Criterion

$$\sum_{i=1}^n \ell_{\eta}(-Y_i \eta(X_i)) = \sum_{i=1}^n \ell_{\eta}(Y_i \eta(X_i)) \approx \sum_{i=1}^n \ell_{\eta}(Y_i \eta(X_i)) = \sum_{i=1}^n \mathbf{1}_{[Y_i \eta(X_i) < 0]}$$

Easier for statisticians than ML approach

Motivate boosting-like algorithm



# Closer Look

---

Goal: Predicting  $Y \in \{\pm 1\}$  by the sign of estimated  $\hat{f}_t$ .  
 $f_t: \mathcal{X} \rightarrow \mathcal{R}$ .

$$f_t(X) = \sum_{Y \in \{\pm 1\}} Y \mathbf{1}_{[Y = f_t(X)]} \approx \sum_{Y \in \{\pm 1\}} Y \mathbf{1}_{[Y = \hat{f}_t(X)]}$$

Minimize  $\sum_{X \in \mathcal{X}} \sum_{Y \in \{\pm 1\}} Y f_t(X) \mathbf{1}_{[Y = \hat{f}_t(X)]}$ . Update  $\hat{f}_t$  by  $\hat{f}_t + \eta$  with  $\eta = \pm 1$   $\mathcal{R}$

For fixed  $\eta$  and  $\mathcal{X}$ , expand at  $\hat{f}_t = 0$

$$\hat{f}_{t+1}(X) = \hat{f}_t(X) + \eta \cdot \mathbf{1}_{\{\hat{f}_t(X) < 0\}}$$

---

# Motivating Questions

---

Convergence: Whether this iterative update converge?

Consistency: Does it converge to the optimal Bayes with respect to  $\rho_{\eta}(Y(X)) = 1_{[Y(X)]}$ ?

Mease and Wyner (2007). Evidence Contradictory to Statistical View. of Boosting

# Questions Solved?

---

“Statistic View”: AdaBoost as a conditional risk minimizer wrt some approximate losses

- AdaBoost can overfit





# Normal-normal setting

---

Let  $X \sim (\theta)$  and  $(\theta) \sim (\quad)$ , w/ known  $\sigma^2$  and  $\mu$   
 Posterior  $(\theta | X) \sim (\quad)$ , where

$$= \frac{1}{\sigma^2 + \tau^2} \left( \frac{\tau^2}{\sigma^2} + \frac{\sigma^2}{\tau^2} \right) = \frac{\tau^2 + \sigma^2}{\sigma^2 \tau^2}$$

$$= \frac{1}{\sigma^2} + \frac{1}{\tau^2} = \frac{\tau^2 + \sigma^2}{\sigma^2 \tau^2}$$

And marginal density of  $X$

$$f(x) = \frac{1}{\sqrt{2\pi(\sigma^2 + \tau^2)}} \exp \left\{ -\frac{(x - \mu)^2}{2(\sigma^2 + \tau^2)} \right\}$$

# Iterative Bayes

# $B$ : Derivation

---

Follow the steps similar to FFFim(2.197(t)0.1470.197(t5266870

---

# Iterative Bayes

# $B$ : Iteration

$$\begin{aligned} \bar{\mu}_{B,t+1}(\cdot) &= \bar{\mu}_{B,t}(\cdot) + \mu_t(\cdot) \\ &= \bar{\mu}_{B,t}(\cdot) + \frac{(\sqrt{\gamma}) - \text{PIB},t(\cdot) [1 - (\sqrt{\gamma})]}{(\sqrt{\gamma}) + \text{PIB},t(\cdot) [1 - (\sqrt{\gamma})]} \end{aligned}$$

Does  $\bar{\mu}_{B,t}$  converge?

Does  $\bar{\mu}_{B,t}$  to the optimal Bayes procedure wrt  $\rho$  ?

**Theorem 1.** For any initial  $\pi_{i|B}^0(\cdot)$ , as  $i$  goes to infinity

$$\pi_{i|B}^0(\cdot) \rightarrow \pi(\cdot) = \frac{1}{2} \ln \left( \frac{(\sqrt{\cdot})}{1 - (\sqrt{\cdot})} \right)$$



**Lemma 1** (Fixed Point Theorem). If

# : Derivation

---

$$l(\cdot) \leftarrow l(\cdot) + \frac{1}{2} \ln \left( \frac{1 - \text{err}}{\text{err}} \right) (\cdot) \quad (1)$$

---

# FHT's AdaBoost: Convergence

---

By calculation, the iteration becomes

$$\begin{aligned} \left( \right) &\leftarrow \left( \right) + \frac{\left( \right)}{2} \left[ \ln \left( \frac{\left( \sqrt{\left( \right)} \right)}{1 - \left( \sqrt{\left( \right)} \right)} \right) - 2 \left( \right) \right] \\ &= \frac{1}{2} \ln \left( \frac{\left( \sqrt{\left( \right)} \right)}{1 - \left( \sqrt{\left( \right)} \right)} \right) \end{aligned}$$

**Remark 1.** One-step convergence

# Bayes Risk $\int_{\mathcal{X}} \ell(\theta; X) p(\theta) d\mu(X)$

Difficulty of the problem

Overfitting

$$\ell(\theta; X) = 2 \left( \sqrt{\frac{1}{2} + \frac{1}{2} \cos(X - \theta)} \right) - 1 \text{ and}$$

$$p(\theta) = \frac{1}{2\pi} \exp\left(-\frac{\theta^2}{2}\right)$$

Let  $\mu = 0$  and assume  $\sigma^2 = 1$

$$\begin{aligned} \int_{\mathcal{X}} \ell(\theta; X) p(\theta) d\mu(X) &= \int_{-\pi}^{\pi} \left( 2 \sqrt{\frac{1}{2} + \frac{1}{2} \cos(X - \theta)} - 1 \right) \frac{1}{2\pi} \exp\left(-\frac{\theta^2}{2}\right) d\theta \\ &= 2 \int_{-\pi}^{\pi} \sqrt{\frac{1}{2} + \frac{1}{2} \cos(X - \theta)} \frac{1}{2\pi} \exp\left(-\frac{\theta^2}{2}\right) d\theta - \int_{-\pi}^{\pi} \frac{1}{2\pi} \exp\left(-\frac{\theta^2}{2}\right) d\theta \end{aligned}$$

where  $\theta \sim \mathcal{N}\left(0, \frac{2 + \tau^2}{\tau}\right)$



# Summary

---

---

# Concluding Remark

---

---