

Ada Boost  $y_t \in \{\pm 1\}$   $x_t \in \mathbb{R}^d$ Input:  $D_n = \{(x_t, y_t)\}_{t=1}^n$   $m$ : total # of boosting iteration- Initialize weight  $W_0(t) = \frac{1}{n}$   $t=1, 2, \dots, n$ 

- At boosting stage  $m \geq 1$ . find a base learner  $h(\cdot; \hat{\theta}_m)$   $S, K, \theta_0$   
 $\Uparrow$   
 that minimizes we can use other Base Learner as well!

$$\hat{\theta}_m = \underset{\theta_m}{\operatorname{argmin}} - \sum_{t=1}^n \underbrace{\tilde{W}_{m-1}(t)}_{\text{weighted training error}} y_t h(x_t, \theta_m)$$

probability distribution

- Choose  $\hat{\alpha}_m \rightarrow \hat{\alpha}_m = \frac{1}{2} \ln \left( \frac{1 - \hat{\epsilon}_m}{\hat{\epsilon}_m} \right)$ 

- Update  $\tilde{W}_m(t) = \begin{cases} \frac{\tilde{W}_{m-1}(t)}{Z_m} e^{-\hat{\alpha}_m} & , \text{ classify correctly} \\ \frac{\tilde{W}_{m-1}(t)}{Z_m} e^{\hat{\alpha}_m} & , \text{ classify wrongly} \end{cases}$

- Go to minimize Weighted Training Error To attain  $\tilde{\theta}_{m+1}$ 

$\Downarrow$   
 $h(\cdot; \hat{\theta}_{m+1})$

Thm: If each base learner is slightly better than RANDOM GUESSING, i.e.,  $\hat{\epsilon}_j < \frac{1}{2} \Rightarrow$  weighted training Error

Then We have Training Error of Ensemble Learner  
 decrease to 0! (exponentially fast) //

$$h_m(x) = \sum_{j=1}^m \alpha_j h(x; \hat{\theta}_j)$$

$$\frac{1}{n} \sum_{t=1}^n \mathbb{1} \{ y_t h_m(x_t) \leq 0 \} \leq \exp \left\{ -2 \sum_{j=1}^m \left( \frac{1}{2} - \hat{\epsilon}_j \right)^2 \right\} \quad \forall m$$

Weighted Training Error

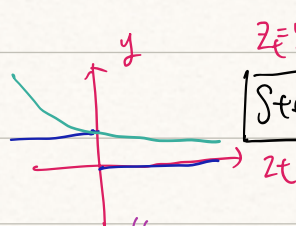
$$\hat{\epsilon}_j = \sum_{t=1}^n \tilde{w}_{j-1}(t) \mathbb{1} \{ y_t \neq h(x_t, \hat{\theta}_j) \} \leq \exp \{ -2m\gamma^2 \}$$

Objective of choosing  $\hat{\theta}_j = \argmin \underline{2\epsilon_j - 1}$

$$\Leftrightarrow \argmin \epsilon_j$$

if  $\frac{1}{2} - \hat{\epsilon}_j \geq \gamma$   
uniform Bound!

Pf:  $\frac{1}{n} \sum \mathbb{1} \{ y_t h_m(x_t) \leq 0 \} \leq \exp \{ -2 \sum (\frac{1}{2} - \hat{\epsilon}_j)^2 \}$



$$z_t = y_t h_m(x_t)$$

Step 1:

$$\frac{1}{n} \sum \mathbb{1} \{ y_t h_m(x_t) \leq 0 \} \leq \frac{1}{n} \sum_{t=1}^n \exp \{ -y_t h_m(x_t) \}$$

exp loss is an upper Bound of 0-1 loss! (convex & Smooth)

Can be very small!

Step 2: Claim:  $\frac{1}{n} \sum \exp(-y_t h_m(x_t)) = \frac{\sum_{j=1}^m z_j}{m}$

Where  $z_j = \sum \tilde{w}_{j-1} \exp(-y_t \alpha_t h(x_t, \hat{\theta}_j))$

$\hookrightarrow$  the normalization Factor

Pf: Note that  $\tilde{w}_0(t) = w_0(t) = \frac{1}{n} \quad \forall t$

$$\tilde{w}_1(t) = \begin{cases} \frac{\tilde{w}_0(t) \exp\{\alpha_1\}}{Z_1} \\ \frac{\tilde{w}_0(t) \exp\{-\alpha_1\}}{Z_1} \end{cases}$$

$$\Leftrightarrow \tilde{w}_1(t) = \frac{1}{Z_1} \tilde{w}_0(t) \exp \{ -\alpha_1 y_t h(x_t, \hat{\theta}_1) \}$$



$$\tilde{w}_2(t) = \frac{1}{Z_2} \tilde{w}_1(t) \exp \{ -\alpha_2 y_t h(x_t; \hat{\theta}_2) \}$$

$$= \frac{1}{n} \frac{1}{Z_1 Z_2} \exp \{ -y_t [ \alpha_1 h(x_t; \theta_1) + \alpha_2 h(x_t; \theta_2) ] \}$$

$$\Rightarrow \tilde{w}_m(t) = \frac{1}{n} \frac{1}{\prod Z_i} \exp \{ -y_t h_m(x_t) \}$$

$$\left( \sum_{t=1}^n \tilde{w}_m(t) = 1 \right)$$

$$\Rightarrow \prod_{i=1}^m Z_i = \frac{1}{n} \sum_{t=1}^n \exp \{ -y_t h_m(x_t) \} \quad \square$$

Step 3    If  $Z_j < 1$ , training Error Decays Exp. FAST

Pf:     $Z_m = \sum_{t=1}^n \tilde{w}_{m-1}(t) \exp(-\hat{\alpha}_m y_t h(x_t; \hat{\theta}_m))$

$$= e^{-\hat{\alpha}_m} \sum \tilde{w}_{m-1}(t) \cdot \mathbb{1} \{ y_t = h(x_t; \hat{\theta}_m) \}$$

In order to attain Training Error as small as possible, we minimize  $Z_m$ !

$$+ e^{\hat{\alpha}_m} \sum \tilde{w}_{m-1}(t) \mathbb{1} \{ y_t \neq h(x_t; \hat{\theta}_m) \}$$

$$= e^{-\hat{\alpha}_m} \cdot (1 - \hat{\epsilon}_m) + e^{\hat{\alpha}_m} \hat{\epsilon}_m$$

Notice that  $Z_m$  is convex in  $\hat{\alpha}_m$

$$\Rightarrow \frac{dZ_m}{d\hat{\alpha}_m} = 0 \Rightarrow e^{-\hat{\alpha}_m} (\hat{\epsilon}_m - 1) + e^{\hat{\alpha}_m} \hat{\epsilon}_m = 0$$

$$\Rightarrow \hat{\alpha}_m = \frac{1}{2} \ln \left( \frac{1 - \hat{\epsilon}_m}{\hat{\epsilon}_m} \right)$$

↘ show the choice of  $\hat{\alpha}_m$ .

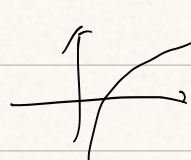
Step 4 : Substitute  $\hat{\alpha}_m$  to  $Z_m$  !

$$Z_m = e^{-\hat{\alpha}_m (1 - \hat{\epsilon}_m)} + e^{\hat{\alpha}_m \hat{\epsilon}_m} \quad \left| \quad \hat{\alpha}_m = \frac{1}{2} \ln \left( \frac{1 - \hat{\epsilon}_m}{\hat{\epsilon}_m} \right) \right.$$

$$= \left( \frac{1 - \hat{\epsilon}_m}{\hat{\epsilon}_m} \right)^{-\frac{1}{2}} (1 - \hat{\epsilon}_m) + \left( \frac{1 - \hat{\epsilon}_m}{\hat{\epsilon}_m} \right)^{\frac{1}{2}} \hat{\epsilon}_m$$

$$= 2 \sqrt{\hat{\epsilon}_m (1 - \hat{\epsilon}_m)}$$

Putting Everything together: Training Error  $\leq \prod_{i=1}^m Z_i$   
 Best Possible  $\hat{\alpha}_m$

  $\ln u < u-1$  (for the choice of  $\hat{\alpha}_m$ )  $= \prod_{i=1}^m 2 \sqrt{\hat{\epsilon}_i (1 - \hat{\epsilon}_i)}$

$$(1-u)^{\frac{1}{2}} = \exp \left( \ln (1-u)^{\frac{1}{2}} \right) = \prod_{i=1}^m \sqrt{1 - (1 - 2\hat{\epsilon}_i)^2}$$

$$= \exp \left( \frac{1}{2} \ln (1-u) \right) \Rightarrow$$

$\ln u < u-1$   
 $\Leftrightarrow \ln (1-u) < -u$

$$\Rightarrow < \exp \left( -\frac{1}{2} u \right)$$

$$\leq \prod_{i=1}^m \exp \left( -\frac{1}{2} (1 - 2\hat{\epsilon}_i)^2 \right)$$

$$= \exp \left\{ -\frac{1}{2} \sum_{i=1}^m (1 - 2\hat{\epsilon}_i)^2 \right\}$$

$$= \exp \left\{ -2 \sum_{i=1}^m \left( \frac{1}{2} - \hat{\epsilon}_i \right)^2 \right\}$$

Counter Example : ( Decision Stump CANNOT DO very well )







$$\Rightarrow \hat{\xi}_j = \frac{1}{2} \text{ for } \forall j$$

since the weight will never update.

Other types of Error

① Weighted Error Relative to Update weights.  $\tilde{W}_m(t)$

Claim: 
$$\sum_{t=1}^n \tilde{W}_m(t) \mathbb{1} \{ y_t \neq h(x_t; \hat{\theta}_m) \} = \frac{1}{2}$$

Rmk: The Base Learner  $\hat{\theta}_m (h(\cdot; \hat{\theta}_m))$  is useless in the next boosting iteration.

only NEXT ONE

since  $h(\cdot; \hat{\theta}_m)$  is attain by:

$$\hat{\theta}_m = \arg \min_{\theta} \sum \tilde{W}_{m-1}(t) \mathbb{1} \{ y \neq h(x_t; \theta) \}$$

Pf: 
$$2 \mathbb{1} \{ y_t \neq h(x_t; \hat{\theta}_m) \} - 1 = -y_t h(x_t; \hat{\theta}_m)$$

$\Rightarrow$  Our conclusion  $\Leftrightarrow \sum \tilde{W}_m(t) y_t h(x_t; \hat{\theta}_m) = 0$   $\hat{\theta}_{m+1} \neq \hat{\theta}_m$

weighted agreement  $\Downarrow$  choose a different decision stump.

$= \sum_{\text{correct}} \tilde{W}_m(t) y_t h(x_t; \hat{\theta}_m)$

$+ \sum_{\text{wrongly}} \tilde{W}_m(t) y_t h(x_t; \hat{\theta}_m)$

$= \sum_{\text{correct}} \frac{\tilde{W}_{m-1} e^{-\alpha_m} y_t h(x_t; \hat{\theta}_m)}{Z_m}$

However, it only guarantees that we don't choose the same decision stump IN TWO CONSECUTIVE ITERATION!

10th decision stump can be the same with 20th decision stump.

$$+ \sum_{\text{wrongly}} \frac{\tilde{W}_{m+1}(t) e^{\alpha_m} y_t h(x_t; \hat{\theta}_m)}{Z_m}$$

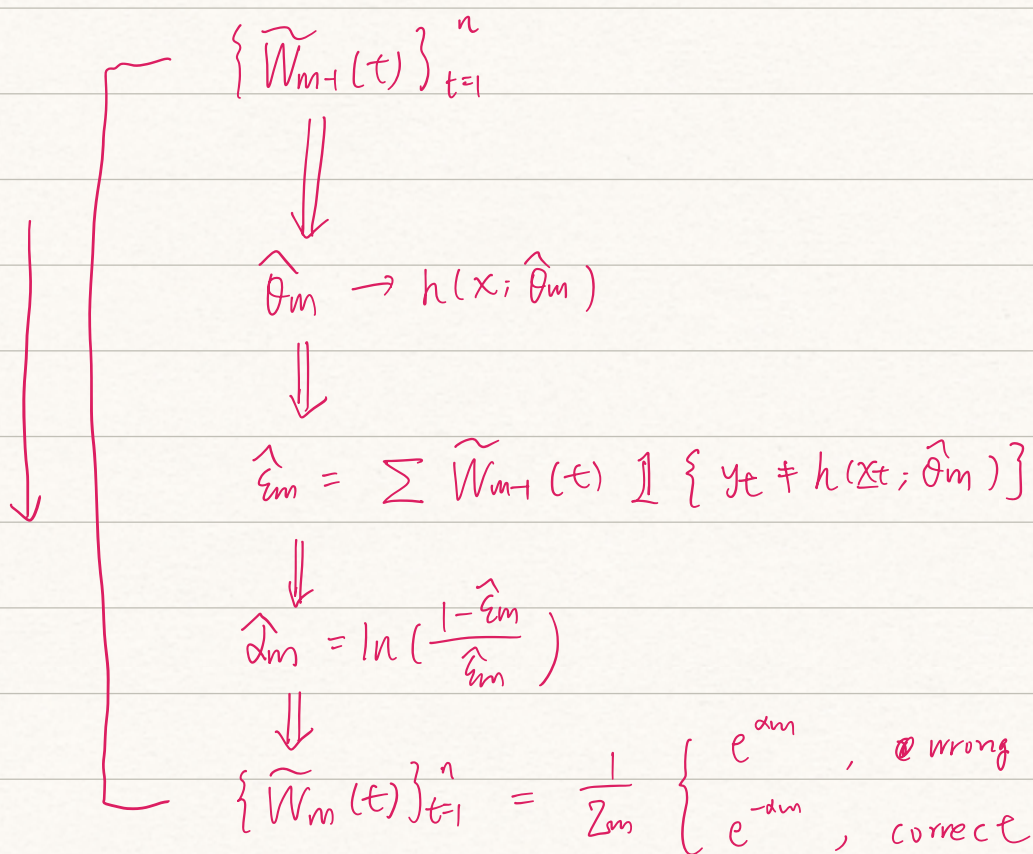
$$= \frac{1}{Z_m} \left( \sum_{\text{correct}} e^{-\alpha_m} \tilde{W}_{m+1}(t) - \sum_{\text{wrongly}} e^{\alpha_m} \tilde{W}_{m+1}(t) \right)$$

$$= \frac{1}{Z_m} \left( e^{-\alpha_m} (1 - \hat{\epsilon}_m) - e^{\alpha_m} \hat{\epsilon}_m \right) = 0$$

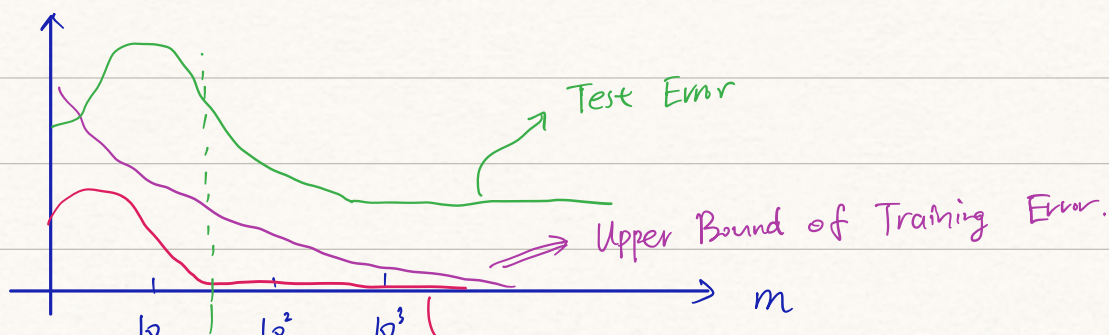
$$\frac{dZ_m}{d\alpha_m} = 0 \Rightarrow \text{choose best possible } \hat{\alpha}_m \text{ for } Z_m$$

#

EMPHASIZE THE SEQUENCE



Test Error





# of iteration / Boosting Rounds.  
↓  
Training Error

Rmk: 1. Test error continue ↓ even after training error = 0

2. Test Error Does NOT ↑ even after many  
Boosting iteration! ⇒ Very Robust!

↓  
Simulation! → 1 vs 7/8 Pattern Recognition

---

Heuristic Explanation: (for ①) Normalized Ensemble Classifier

$$\tilde{h}_m(x) = \frac{\sum \hat{\alpha}_j h(x; \hat{\theta}_j)}{\sum \alpha_i} \in [-1, 1]$$

Define Voting Margin for each  $t$ :

↓

$$\text{margin}(t) := y_t \tilde{h}_m(x_t) \in [-1, 1]$$

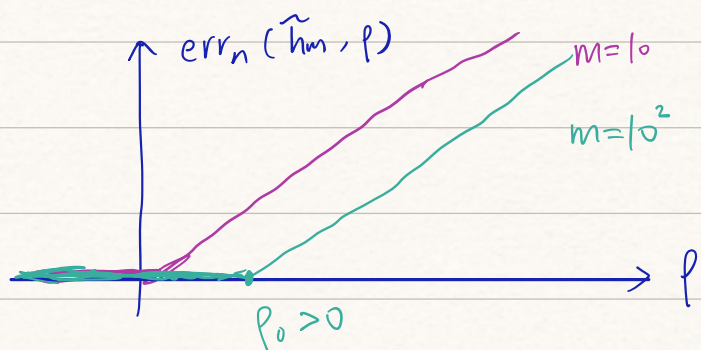
Rmk:  $\text{margin}(t) > 0 \Leftrightarrow$  example  $\tilde{x}_t$  is classified correctly  
by  $h_m(\cdot)$  ( $\tilde{h}_m(\cdot)$ )

Def:  $p$ -training error:  $\text{err}_n(\tilde{h}_m; p) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{\text{margin}(t) \leq p\}$

Rmk: Training Error =  $\text{err}_n(\tilde{h}_m, 0)$

From training error theorem:  $\text{err}_n(\tilde{h}_m; 0) = 0$  after finitely  $m$

But even after  $\text{err}_n(\tilde{h}_m, 0) = 0$ , boosting still decreases  $\text{err}_n(\tilde{h}_m, p)$  for some  $p > 0$



Why? Because we are implicitly minimize the 0-1 loss!

Training Error Thm:  $\Downarrow$  we actually minimize this!

$$0-1 \text{ loss} \leq \text{exp. loss} = \exp(-\text{margin}(t))$$

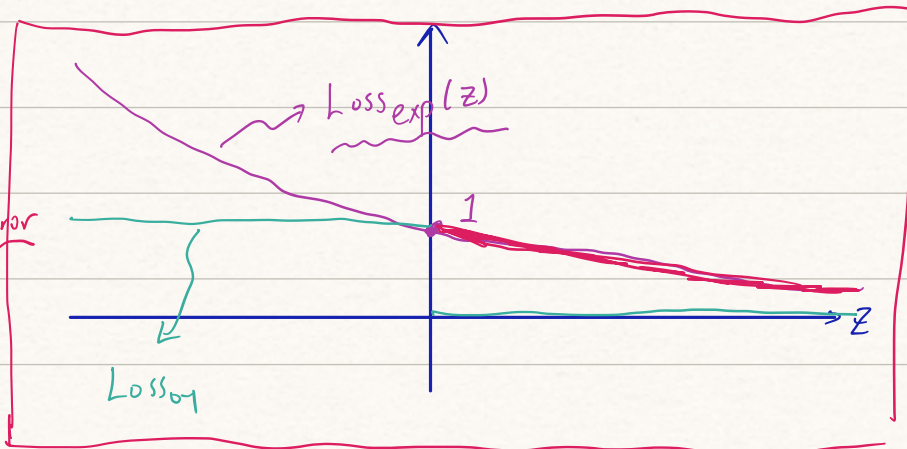
★ When we do update,

we actually minimize

$\text{Loss}_{\text{exp}}(\text{exp. Error})$

0-1 Training Error decrease

$p$ -Training Error decrease



Resistance to overfitting

- ① "Complexity" of ensemble classifier  $h_m \uparrow$  rapidly with  $m$
- ② Sequential optimizing instead of Joint (Greedy)



Define: Expected Risk.

$$R(f) = \mathbb{E} [\text{loss}_{0-1}(y, f(x))] = \Pr(y \neq f(x))$$

Empirical Risk

$$\hat{R}_n(f) = \frac{1}{n} \sum_{t=1}^n \text{Loss}_{0-1}(y_t, f(x_t))$$

Thm: Freund & Schapire (1995)

With prob  $\geq 1 - \delta$ ,  $R(h_m) \leq \hat{R}_n(h_m) + \underbrace{\tilde{O}\left(\sqrt{\frac{md}{n}}\right)}_{\substack{\text{according to } m \\ \downarrow \\ \text{very loose}}}$

$\left\{ \begin{array}{l} m: \# \text{ of boosting rounds} \\ n: \# \text{ of training sample} \\ d: \text{VC dimension of base learner} \end{array} \right.$

$\hookrightarrow$  a measure of complexity.

Rmk: This Bound is not predictive of Good (Robust) test error because of its dependence on  $m$ !



Not Good Thm!

Not related to  $m$ !



☆☆☆

Thm: With Prob  $\geq 1 - \delta$ ,  $R(h_m) \leq \text{err}_n(h_m; P) + \underbrace{\tilde{O}\left(\sqrt{\frac{d}{n p^2}}\right)}_{\substack{\text{fixed with } m \text{ grows?}}}$

$\parallel$   
 $\frac{1}{n} \sum_{t=1}^n \mathbb{1}\{y_t \hat{h}_m(x_t) \leq p\}$

$\Rightarrow$  if we choose  $\forall p < \gamma = \min(\frac{1}{2} - \epsilon_j)$

then  $\text{err}_n(\tilde{h}_m, p) \rightarrow 0$  exponentially fast !

★

Pmk: Essentially this bound says that there is no  
more dependence of the Expected Risk on  $m$



implies that Adaboost ( Boosting is  
Robust to overfitting )

*exp. Loss*