



Boosting: Theory

Machine Learning: Jordan Boyd-Graber
University of Maryland

SLIDES ADAPTED FROM ROB SCHAPIRE

Training Error

First, we can prove that the training error goes down. If we write the the error at time t as $\frac{1}{2} - \gamma_t$,

$$\hat{R}(h) \leq \exp \left\{ -2 \sum_t \gamma_t^2 \right\} \quad (1)$$

- If $\forall t: \gamma_t \geq \gamma > 0$, then $\hat{R}(h) \leq \exp \{-2\gamma^2 T\}$

Adaboost: do not need γ or T a priori

Training Error Proof: Preliminaries

Repeatedly expand the definition of the distribution.

$$D_{t+1}(i) = \frac{D_t(i) \exp\{-\alpha_t y_i h_t(x_i)\}}{Z_t} \quad (2)$$

$$\frac{D_{t-1}(i) \exp\{-\alpha_{t-1} y_i h_{t-1}(x_i)\} \exp\{-\alpha_t y_i h_t(x_i)\}}{Z_{t-1} Z_t} \quad (3)$$

$$\frac{\exp\{-y_i \sum_{s=1}^t \alpha_s h_s(x_i)\}}{m \prod_{s=1}^t Z_s} \quad (4)$$

Training Error Intuition

- On round t weight of examples incorrectly classified by h_t is increased
- If x_i incorrectly classified by H_T , then x_i wrong on (weighted) majority of h_t 's
 - If x_i incorrectly classified by H_T , then x_i must have large weight under D_T
 - But there can't be many of them, since total weight ≤ 1

Training Error Proof: It's all about the Normalizers

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [y_i g(x_i) \leq 0] \quad (5)$$

(6)

Definition of training error

Training Error Proof: It's all about the Normalizers

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [y_i g(x_i) \leq 0] \quad (5)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \exp \{-y_i g(x_i)\} \quad (6)$$

$$(7)$$

$\mathbb{1} [u \leq 0] \leq \exp -u$ is true for all real u .

Training Error Proof: It's all about the Normalizers

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [y_i g(x_i) \leq 0] \quad (5)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \exp \{-y_i g(x_i)\} \quad (6)$$

$$(7)$$

Final distribution $D_{t+1}(i)$

$$D_{t+1}(i) = \frac{\exp \{-y_i \sum_{s=1}^t \alpha_s h_s(x_i)\}}{m \prod_{s=1}^t Z_s} \quad (8)$$

Training Error Proof: It's all about the Normalizers

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [y_i g(x_i) \leq 0] \quad (5)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \exp \{-y_i g(x_i)\} \quad (6)$$

$$= \frac{1}{m} \sum_{i=1}^m \left[m \prod_{t=1}^T Z_t \right] D_{T+1}(i) \quad (7)$$

$$(8)$$

m 's cancel, D is a distribution

Training Error Proof: It's all about the Normalizers

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [y_i g(x_i) \leq 0] \quad (5)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \exp \{-y_i g(x_i)\} \quad (6)$$

$$= \frac{1}{m} \sum_{i=1}^m \left[m \prod_{t=1}^T Z_t \right] D_{T+1}(i) \quad (7)$$

$$= \prod_{t=1}^T Z_t \quad (8)$$

Training Error Proof: Weak Learner Errors

Single Weak Learner

$$Z_t = \sum_{i=1}^m D_t(i) \exp\{-\alpha_t y_i h_t(x_i)\} \quad (9)$$

$$= \quad (10)$$

$$= \quad (11)$$

$$= \quad (12)$$

Training Error Proof: Weak Learner Errors

Single Weak Learner

$$Z_t = \sum_{i=1}^m D_t(i) \exp\{-\alpha_t y_i h_t(x_i)\} \quad (9)$$

$$= \sum_{i:\text{right}} D_t(i) \exp\{-\alpha_t\} + \sum_{i:\text{wrong}} D_t(i) \exp\{\alpha_t\} \quad (10)$$

$$= \quad (11)$$

$$= \quad (12)$$

Training Error Proof: Weak Learner Errors

Single Weak Learner

$$Z_t = \sum_{i=1}^m D_t(i) \exp\{-\alpha_t y_i h_t(x_i)\} \quad (9)$$

$$= \sum_{i:\text{right}} D_t(i) \exp\{-\alpha_t\} + \sum_{i:\text{wrong}} D_t(i) \exp\{\alpha_t\} \quad (10)$$

$$= (1 - \epsilon_t) \exp\{-\alpha_t\} + \epsilon_t \exp\{\alpha_t\} \quad (11)$$

$$= \quad (12)$$

Training Error Proof: Weak Learner Errors

Single Weak Learner

$$Z_t = \sum_{i=1}^m D_t(i) \exp\{-\alpha_t y_i h_t(x_i)\} \quad (9)$$

$$= \sum_{i:\text{right}} D_t(i) \exp\{-\alpha_t\} + \sum_{i:\text{wrong}} D_t(i) \exp\{\alpha_t\} \quad (10)$$

$$= (1 - \epsilon_t) \exp\{-\alpha_t\} + \epsilon_t \exp\{\alpha_t\} \quad (11)$$

$$= (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} + \epsilon_t \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \quad (12)$$

Training Error Proof: Weak Learner Errors

Single Weak Learner

$$Z_t = (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} + \epsilon_t \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}} \quad (9)$$

Normalization Product

$$\prod_{t=1}^T Z_t = \prod_{t=1}^T 2 \sqrt{\epsilon_t (1 - \epsilon_t)} = \sqrt{1 - 4 \left(\frac{1}{2} - \epsilon_t \right)^2} \quad (10)$$

$$(11)$$

Training Error Proof: Weak Learner Errors

Normalization Product

$$\prod_{t=1}^T Z_t = \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)} = \sqrt{1-4\left(\frac{1}{2}-\epsilon_t\right)^2} \quad (9)$$

$$\leq \prod_{t=1}^T \exp\left\{-2\left(\frac{1}{2}-\epsilon_t\right)^2\right\} \quad (10)$$

$$(11)$$

Training Error Proof: Weak Learner Errors

Normalization Product

$$\prod_{t=1}^T Z_t = \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)} = \sqrt{1-4\left(\frac{1}{2}-\epsilon_t\right)^2} \quad (9)$$

$$\leq \prod_{t=1}^T \exp\left\{-2\left(\frac{1}{2}-\epsilon_t\right)^2\right\} \quad (10)$$

$$= \exp\left\{-2\sum_{t=1}^T \left(\frac{1}{2}-\epsilon_t\right)^2\right\} \quad (11)$$

Generalization

VC Dimension

$$\leq 2(d + 1)(T + 1) \lg[(T + 1)e]$$

Margin-based Analysis

AdaBoost maximizes a linear program maximizes an L_1 margin, and the weak learnability assumption requires data to be linearly separable with margin 2γ