

## CS 596: Theoretical Machine Learning

Lecturer: Pranjali Awasthi  
Scribe: Pengxiang Wu

Lecture # 8  
Oct 13, 2016

In the PAC learning model, our goal is to produce a hypothesis  $h$  such that  $error(h) \leq \epsilon$ , with probability  $\geq 1 - \delta$ . This is a *strong* learning model, since  $\epsilon$  and  $\delta$  could be arbitrarily small. Now suppose we are given a set of *weak* algorithms which are good but not perfect, is it possible to combine them to produce a strong and perfect algorithm? This process is called *boosting*, which is the main topic of this lecture.

Intuitively, if the weak algorithms complement each other, we could fuse their respective advantages to achieve a perfect algorithm. More specifically, we will prove the following theorem

**Theorem 1 (Boosting).** *Given function class  $H$  and the target function  $h^* : \mathcal{X} \mapsto \{-1, +1\}$ , if for some  $\gamma > 0$  and any distribution  $D$  over  $\mathcal{X}$ ,  $\exists h \in H$  such that  $err_D(h) \leq \frac{1}{2} - \gamma$ , then: (1) for  $\forall \epsilon > 0$  and  $\forall D$  over  $\mathcal{X}$ ,  $\exists g \in MAJ_k(H)$  such that  $err_D(g) \leq \epsilon$ , where  $MAJ_k(H) = \{\text{sgn}(\sum_{i=1}^k \alpha_i h_i(x)); h_i \in H\}$ , and  $k = O(\frac{1}{\gamma^2} \log(\frac{1}{\epsilon}))$ ; (2) such  $g$  can be found efficiently given access to a weak learning algorithm.*

Here  $err_D(g)$  represents true error with respect to the distribution, i.e.,  $err_D(g) = Pr_{x \sim D}[g(x) \neq h^*(x)]$ . The improved function  $g$  can be found using AdaBoost, as we will introduce later.

Theorem 1 implies two facts. First, weak learning and PAC learning are equivalent. Of course, a weak learner is not the same thing as a strong learner; by contrast, the error probability of a weak learner is required only to be slightly smaller than half. Here what we mean by “equivalent” is that, a problem can be weak-learned if and only if it can be strong-learned.

**Intuitive explanation:** Suppose there are three hypotheses  $h_1, h_2, h_3 \in H$  that make independent errors probability of  $p$ . Then the error rate of their majority vote would be  $p^3 + 3p^2(1 - p)$ . If  $h_1, h_2, h_3$  are produced by a weak learning algorithm (i.e.,  $p < \frac{1}{2}$ ), we have  $p^3 + 3p^2(1 - p) < p$ , which indicates an improved performance. One can then bootstrap this process to get down to any arbitrary error rate of  $\epsilon$ .

Now we give a formal proof of the first part of theorem 1.

*Proof.* We prove the theorem by setting up an experts problem and using the guarantee of the randomized multiplicative weights (RMW) algorithm. Let  $X = \{x_1, x_2, \dots, x_N\}$  be the set of experts (i.e., examples). On day  $t$ , the algorithm picks one expert  $x_t \in X$ , and the adversary returns a hypothesis  $h_t \in H$  accordingly. We define the loss as  $1 - h_t(x_t)h^*(x_t)$ . Thus when  $h_t$  predicts correctly, the loss would be 0, otherwise 1. Let's assume that we are using RMW to pick a distribution over experts at each stage. For a given input probability distribution  $p_t$ , the strategy taken by adversary is to find  $h_t \in H$  with maximum loss. Now given arbitrary hypothesis  $h$ , the expected loss of RMW at time  $t$  would be  $E[\text{loss}(h)] = \sum_{x \in X} p_t(x)(1 - h(x)h^*(x))$ . Since, there always exists a weak learner for any distribution  $p_t$ , the adversary can always make sure that the expected loss of RMW at each step is at least  $1 - 2\gamma$ . We also know that on any sequence of length  $T$ , RMW guarantees that  $E[\text{loss}(RMW)] \leq \text{loss of best expert} + \sqrt{\frac{\log N}{T}}$ , we have:

$$1 - 2\gamma \leq E[\text{loss}(RMW)] \leq \text{loss of best expert} + \gamma \quad [\text{when } T \text{ is large}], \quad (1)$$

where the second inequality comes from the fact that when  $T$  is sufficiently large,  $\sqrt{\frac{\log N}{T}}$  will be smaller than  $\gamma$  (since  $\gamma$  is a small number). Thus it holds that:

$$\text{loss of best expert} \geq 1 + \gamma,$$

which means:

$$\text{loss of any expert} \geq 1 + \gamma. \quad (2)$$

Hence we have the following inequality:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T (1 + h_t(x)h^*(x)) \geq 1 + \gamma, \forall x \\ \implies & \frac{1}{T} \sum_{t=1}^T h_t(x)h^*(x) \geq \gamma, \forall x \\ \implies & \left[ \frac{1}{T} \sum_{t=1}^T h_t(x) \right] h^*(x) \geq \gamma, \forall x. \end{aligned} \quad (3)$$

Let  $S = \frac{1}{T} \sum_{t=1}^T h_t(x)$ , and  $q_i = \frac{\text{the number of times } h_i \in H \text{ appears in } S}{T}$ , then it holds that  $\frac{1}{T} \sum_{t=1}^T h_t(x) = \sum_{i=1}^{|H|} q_i h_i(x)$ . Thus  $\forall x \in X$ ,  $[\sum_{i=1}^{|H|} q_i h_i(x)]h^*(x) \geq \gamma$ . Notice that when  $h^*(x) = 1$ ,  $\sum_{i=1}^{|H|} q_i h_i(x) \geq \gamma$ ; similarly,  $\sum_{i=1}^{|H|} q_i h_i(x) \leq -\gamma$  when  $h^*(x) = -1$ . Next we sample  $k$  functions  $h_1, h_2, \dots, h_k$  from distribution  $\vec{q} = \{q_1, q_2, \dots, q_{|H|}\}$ . According to Hoeffding's inequality, we have:

$$\Pr \left( \left| \frac{1}{k} \sum_{i=1}^k h_i(x) - \sum_{i=1}^{|H|} q_i h_i(x) \right| \geq \gamma \right) \leq 2e^{-2k\gamma^2}, \quad (4)$$

where  $\sum_{i=1}^{|H|} q_i h_i(x)$  represents the true average of functions. Now fix distribution  $D$  over  $\mathcal{X}$ , and let  $g = \text{sgn}(\frac{1}{k} \sum_{i=1}^k h_i(x))$ , we get:

$$\begin{aligned} \text{err}_D(g) &= \sum_x P(x) \Pr_g(g \text{ is incorrect in } x | x) \\ &= \sum_x P(x) \Pr_g \left( \left[ \frac{1}{k} \sum_{i=1}^k h_i(x) \right] h^*(x) < 0 \right) \\ &\leq \sum_x P(x) \cdot 2e^{-2k\gamma^2}, \end{aligned} \quad (5)$$

where  $P(x)$  is the probability of picking  $x$ . Since we want  $\text{err}(g) \leq \epsilon$ , it follows that  $\sum_x P(x) \cdot 2e^{-2k\gamma^2} \leq \epsilon$ . Therefore,  $k$  should satisfy  $k \geq \frac{1}{2\gamma^2} \log \frac{2}{\epsilon}$ , i.e.,  $k = O(\frac{1}{\gamma^2} \log(\frac{1}{\epsilon}))$ .  $\square$

Next we prove that strong learner  $g$  can be found efficiently given an algorithm to find a weak learner for every distribution. Let  $X$  be the set of experts, which could be infinite or finite. Now we sample  $m$  i.i.d. training examples  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  from  $D$ . Let  $U_m$  be the uniform distribution over these examples. In the following we will show how to get  $g = \text{sgn}(\frac{1}{k} \sum_{i=1}^k h_i(x)) \in G = \text{MAJ}_k(H)$  such that  $\widehat{\text{err}}(g) = \frac{1}{m} \sum_{i=1}^m I(g(x_i) \neq y_i) \leq \epsilon$ .

From VC theory this would be enough to get small true error as well. Let  $d = VCdim(G)$ , then we have:

$$err(g) \leq \epsilon + \sqrt{\frac{d \log(\frac{1}{\delta})}{m}}, \text{ with probability } \geq 1 - \delta. \quad (6)$$

Therefore, when  $m$  goes large,  $err(g)$  approaches  $\epsilon$ . It is worth mentioning that, if  $d' = VCdim(H)$ , then  $VCdim(G) = O(kd' \log(kd'))$ .

The way we get  $g$  is to run RMW as in the previous case, but only over experts defined by training examples. By the previous analysis, we would be efficiently able to find a  $g$  that has low training error. The algorithm that results is known as Adaboost:

---

**Algorithm 1** AdaBoost

---

- 1: **Input:**  $m$  i.i.d. training examples from uniform distribution  $D$
  - 2: **Initialize:**  $P_1 = (\frac{1}{m}, \dots, \frac{1}{m})$
  - 3: **for**  $t = 1, 2, \dots, T$  **do**
  - 4: Get  $h_t$  which is a weak learner from  $P_t$
  - 5: Let  $\epsilon_t = err(\hat{h}_t) =$  error of  $h_t$  on the training set
  - 6:  $P_{t+1}(i) \propto P_t(i) \exp(-\alpha_t)$ , if  $y_i = h_t(x_i)$
  - 7:  $P_{t+1}(i) \propto P_t(i) \exp(\alpha_t)$ , if  $y_i \neq h_t(x_i)$
  - 8: **end for**
  - 9: **Output:**  $g = \text{sgn}(\sum_{t=1}^T \alpha_t h_t(x))$ , where  $\alpha_t = \frac{1}{2} \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$
- 

**Claim 1.** After  $T$  rounds, the training error of  $g$  is  $\leq e^{-2T\gamma^2}$  in  $T$  rounds. After  $\frac{2}{\gamma^2} \log(\frac{1}{\epsilon})$  rounds, the training error would be  $\leq \epsilon$ .

The above claim states that AdaBoost can find the strong learner  $g$  efficiently. However, it is worth mentioning that, if there exists noise in the training data, AdaBoost could be quite bad.

## Additional Reading

Chapter 10 from the book by Shai Shalev-Shwartz and Shai Ben-David. See online copy: <http://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning/understanding-machine-learning-theory-algorithms.pdf>