

LECTURE 11: ONLINE LEARNING

Instructor: Aditya Bhaskara Scribe: Maks Cegielski-Johnson

CS 5966/6966: Theory of Machine Learning

January 1st, 2017

Abstract

In this lecture, we considered online learning in the case of not having a consistent hypothesis in the hypothesis class. We also compared the regret bound between deterministic and randomized algorithms.

1 INTRODUCTION

As we have previously seen, in the online learning model: points arrive and the learner predicts the label, with the true label being revealed later. The goal of online learning is to compete with the best single hypothesis “in hindsight” after T rounds of learning, where we view each hypothesis as an *expert*. We define *regret* as the difference between “the number of mistakes the learner made” and “the number of mistakes the best expert made”.

2 FINITE SET OF HYPOTHESES

2.1 THEOREM. *If there exists hypothesis $h \in \mathcal{H}$ that makes no mistakes, then we can ensure that the learning algorithm A makes at most $\log |\mathcal{H}|$ mistakes overall.*

From this theorem, the idea for an algorithm is to select what the majority of the “non-eliminated” experts say (and there always is an expert that is never eliminated). If the algorithm makes a mistake, then that means we can eliminate at least half of the hypotheses, which also made mistakes.

3 INFINITE SET OF HYPOTHESES

Suppose we know there is a *consistent* hypothesis, can we ensure that the learning algorithm only makes $O(1)$ mistakes? In the last lecture we saw the Littlestone dimension, and we saw that

$$\text{L-Dim} \leq \text{VC-Dim}$$

4 NON-CONSISTENT FINITE SET OF HYPOTHESES

What if there is no perfect expert? Can we ensure that

Turns out this is not possible

$$\# \text{-mistakes}(Alg) < \# \text{-mistakes}(h_{\text{best}}) + \text{small}, \quad \text{where “small” is } \log n$$

This is known as *regret*.

We have previously seen the idea of eliminating hypotheses that made a mistake. What if we keep track of how many mistakes each hypothesis has made? At each step, each hypothesis h_i is assigned some weight w_i , where w_i is small if h_i made many mistakes. In other words, reweight hypotheses rather than eliminating.

5 MULTIPLICATIVE WEIGHT ALGORITHM

Suppose we have $\mathcal{H} = \{h_1, h_2, \dots, h_N\}$ and initial weights $\{w_1^{(1)}, w_1^{(2)}, \dots, w_1^{(N)}\}$ with each $w_1^{(i)} = 1$ for each h_i . And, suppose that if a hypothesis h_i makes a mistake, then $w_2^{(i)} = w_1^{(i)} \cdot (1 - \eta)$ for some $\eta \in (0, 1)$. This means that after t steps, h_i has weight $w_t^{(i)} = (1 - \eta)^{\#\text{mistakes}(h_i)}$.

This gives us the following algorithm:

- At each step t , keep track of $w_t^{(i)}$ for each hypothesis h_i
- Prediction at time t is

$$\begin{cases} 0 & \text{if } \sum_i \text{who predict 0 at } t \cdot w_t^{(i)} \geq \sum_j \text{who predict 1 at } t \cdot w_t^{(j)} \\ 1 & \text{otherwise.} \end{cases}$$

6 ANALYSIS OF MULTIPLICATIVE WEIGHT ALGORITHM

What happens every time we make a mistake? Let's define the "sum of weights" (also known as *potential*) as

This plays the role of the number of hypotheses which are not eliminated

$$(1) \quad \Phi_t = \sum_i w_t^{(i)} \quad \text{at step } t$$

And let's call $\Phi_t^{(0)}$ the weights of hypothesis that predict 0 at time t , and similarly $\Phi_t^{(1)}$ the weights of the hypothesis that predict 1 at step t .

Suppose $\Phi_t^{(0)} > \Phi_t^{(1)}$. Then, if the learning algorithm predicted 0, and it was wrong, we have that

$$\Phi_{t+1} = \Phi_t^{(1)} + (1 - \eta)\Phi_t^{(0)}$$

If we made a mistake, in the worst case, then

$$\Phi_{t+1} \leq \frac{1}{2}\Phi_t + \frac{1}{2}(1 - \eta)\Phi_t = \left(1 - \frac{\eta}{2}\right)\Phi_t$$

We have the $|\mathcal{H}|$ term because at $t = 0$, each $w_i = 1$ for each hypothesis in \mathcal{H}

So after T steps,

$$\Phi_T \leq |\mathcal{H}| \cdot \left(1 - \frac{\eta}{2}\right)^{\#\text{mistakes}(Alg)}$$

Say $\exists h$ that made only a "few" mistakes

Is there a lower bound on Φ_t ?

If there exists a best hypothesis h_i that made $\leq k$ mistakes, then

$$w_T^{(i)} = (1 - \eta)^k \implies \Phi_T \geq (1 - \eta)^k \implies (1 - \eta)^k \leq |\mathcal{H}| \leq \left(1 - \frac{\eta}{2}\right)^m$$

where $m = \#-mistakes(Alg)$. Then,

$$k \log(1 - \eta) \leq m \log(1 - \eta/2) + \log |\mathcal{H}|$$

and so this implies

$$m \leq k \frac{\log(1 - \eta)}{\log(1 - \eta/2)} - \frac{\log |\mathcal{H}|}{\log(1 - \eta/2)}$$

and then we have

$\log(1 - x) \approx -x$ for small x

$$\frac{\log(1 - \eta)}{\log(1 - \eta/2)} \leq (2 + \eta) \quad \text{and} \quad \frac{\log |\mathcal{H}|}{\log(1 - \eta/2)} \approx 2 \frac{\log |\mathcal{H}|}{\eta}$$

6.1 THEOREM. For any η , the weighted majority algorithm achieves the guarantee:

$$\#-mistakes(Alg) \leq (2 + \eta) \cdot \#-mistakes(h_{opt}) + \frac{1}{\eta} \log |\mathcal{H}|$$

Notice that the optimum-mistake-count term has a factor of 2. We are interested in whether we can reduce this value.

7 RANDOMIZED ALGORITHM

7.1 THEOREM. No deterministic algorithm can achieve a number of mistakes $< 2 \cdot \#-mistakes(opt)$.

Proof (Cover Example)

Consider $\mathcal{H} = \{h_1, h_2\}$ where h_1 always predicts 1 and h_2 always predicts 0. Then, suppose we have any deterministic algorithm A and we get learn and predict on T examples. Now, suppose that the true label is always the opposite of what A predicts. Then, it must be the case that

$$\#-mistakes(A) = T \quad \text{and} \quad \#-mistakes(opt) \leq \frac{T}{2}$$

And this is a 2-factor loss \square .

The model we want to consider is one where the adversary (who knows the true label function f) knows what the learning algorithm is, but does not see the coin-tosses. Suppose the algorithm is

$$\begin{cases} 0 & \text{with probability } p \\ 1 & \text{with probability } 1 - p. \end{cases}$$

What happens in the cover example? The claim is that the algorithm will do as well as the optimal hypothesis in this model.

Question. In this setting, can we beat the factor of 2? That is, can we achieve

$$\mathbb{E}[\#-mistakes(A)] \leq (1 + \epsilon) \#-mistakes(opt) + \frac{\log |\mathcal{H}|}{\epsilon}$$

8 WEIGHTED MAJORITY - VERSION 2

A natural idea is to treat the weights as probabilities. Consider the following algorithm:

- Start with $w_0^{(i)} = 1, \forall i$
- At time t :
 - Algorithm picks h_i with probability $\frac{w_t^{(i)}}{\Phi_t}$ and outputs h_i output
 - For each $j \in [|\mathcal{H}|]$,

$$w_{t+1}^{(j)} = \begin{cases} (1 - \eta)w_t^{(j)} & \text{if } j \text{ was wrong} \\ w_t^{(j)} & \text{otherwise.} \end{cases}$$

8.1 THEOREM. *The expected regret of the algorithm satisfies the “right” bound:*

$$\mathbb{E}[\#-mistakes(A)] \leq (1 + 2\eta)\#-mistakes(opt) + \frac{\log |\mathcal{H}|}{\eta}$$