

# Learnability of a Finite Hypothesis Class in Learning Problems with Noisy Training-Set Labels

Behrad Moniri

University of Pennsylvania

bemoniri@seas.upenn.edu

## Abstract

In this note, the PAC-Learnability of a finite hypothesis class is proved for a learning problem in which the training set labels are flipped with a certain probability.

Let  $\mathcal{H}$  be a finite hypothesis class,  $h^* \in \mathcal{H}$  be the target concept and  $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$  be the training set. The training set labels are corrupted with noise, i.e., the label  $y_i$  is reported as  $h^*(\mathbf{x}_i) \oplus \zeta$  in which  $\zeta \sim \text{Bern}(\eta)$ . We will prove that the algorithm which chooses the function  $h_S$  with the least number of errors on training points, has a low generalization error with high probability.

**Definition 1.** For each  $h \in \mathcal{H}$  and sample set  $S$ , Define  $d(h)$  as the probability that the label  $y$  of a training sample  $\mathbf{x}$  is not equal to  $h(\mathbf{x})$ . Also define  $\hat{d}(h)$  as the empirical loss of  $h$  on the training set,  $\hat{d}(h) = \frac{\sum_{i=1}^m \mathbf{1}[y_i \neq h(\mathbf{x}_i)]}{m}$ .

**Remark 2.** Trivially we have  $d(h^*) = \eta$ .

We will now prove a very useful lemma linking  $d(h)$  and the generalization error of  $h$ .

**Lemma 3.** Let  $\epsilon > 0$  be an arbitrary given number. If  $L_D(h) > \epsilon$ , then we have  $d(h) - d(h^*) \geq \epsilon'$  in which  $\epsilon' = \epsilon(1 - 2\eta)$ .

*Proof.* The label of a training set is flipped with probability  $\eta$ . Hence, we have

$$d(h) = \eta \mathbb{P}[h^*(\mathbf{x}) = h(\mathbf{x})] + (1 - \eta) \mathbb{P}[h^*(\mathbf{x}) \neq h(\mathbf{x})] \quad (1)$$

$$= \eta(1 - L_D(h)) + (1 - \eta)L_D(h) \quad (2)$$

$$= \eta + (1 - 2\eta)L_D(h). \quad (3)$$

For a given  $h$ , if  $L_D(h) > \epsilon$ , we have  $d(h) - \eta = d(h) - d(h^*) > (1 - 2\eta)\epsilon = \epsilon'$ .  $\square$

We will use the well known Hoeffding's inequality (Theorem (4)) in our proof.

**Theorem 4.** Let  $X_1, X_2, \dots, X_m$  be independent random variables. Also assume that each  $X_i$  takes values in  $[a_i, b_i]$  with probability 1. For any  $\epsilon > 0$ , the following inequalities hold for  $S_m = \sum_{i=1}^m X_i$ :

$$\mathbb{P}[S_m - \mathbb{E}[S_m] \geq \epsilon] \leq e^{-2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2}, \quad (4)$$

$$\mathbb{P}[S_m - \mathbb{E}[S_m] \leq -\epsilon] \leq e^{-2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2}. \quad (5)$$

*Proof.* See Appendix (D) of [MRT19].  $\square$

We are now ready to state and prove our main claim.

**Theorem 5.** In the noisy label setting with noise parameters  $\eta$  and  $\eta$  with any distribution  $D$ , let  $\epsilon, \delta > 0$  be arbitrary given numbers. Given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log(|\mathcal{H}|) + \log\left(\frac{2}{\delta}\right) \right)$  training samples, for any  $h \in \mathcal{H}$  with  $L_D(h) > \epsilon$ , we have

$$\hat{d}(h) - \hat{d}(h^*) \geq 0 \quad (6)$$

with probability at least  $\delta$ .

*Proof.* Let  $h \in \mathcal{H}$  be any hypothesis with  $L_D(h) > \epsilon$ . We can decompose  $\hat{d}(h) - \hat{d}(h^*)$  as follows:

$$\hat{d}(h) - \hat{d}(h^*) = [\hat{d}(h) - d(h)] + [d(h) - d(h^*)] + [d(h^*) - \hat{d}(h^*)]. \quad (7)$$

Based on lemma (3), we know that for any  $h$  with  $L_D(h) > \epsilon$ , we have  $d(h) - d(h^*) > \epsilon'$ . We will now show that given enough samples, each of the remaining terms are greater than  $-\frac{\epsilon'}{2}$  with high probability.

1. First, we will show that given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log\left(\frac{2}{\delta}\right) \right)$  samples, we have  $\hat{d}(h^*) - d(h^*) \leq \frac{\epsilon'}{2}$  with probability at least  $1 - \delta/2$ . This is a direct consequence of the Hoeffding's inequality. Note that  $E[\hat{d}(h^*)] = d(h^*)$  and we have

$$\hat{d}(h) = \frac{\sum_{i=1}^m \mathbf{1}[y_i \neq h(\mathbf{x}_i)]}{m} = \sum_{i=1}^m X_i, \quad (8)$$

in which  $X_i \in [0, \frac{1}{m}]$  with probability 1 and  $X_i$ s are jointly independent. The Hoeffding's inequality yields

$$\mathbb{P}[\hat{d}(h^*) - d(h^*) \geq \frac{\epsilon'}{2}] \leq \exp\left[-\frac{m\epsilon'^2}{2}\right] \leq \frac{\delta}{2}. \quad (9)$$

Where the last inequality follows from  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log\left(\frac{2}{\delta}\right) \right)$ .

2. Second, we will prove a *Uniform Convergence* property for  $\mathcal{H}$ . We will prove that given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log(|\mathcal{H}|) + \log\left(\frac{2}{\delta}\right) \right)$  samples, the event

$$\forall h \in \mathcal{H}, d(h) - \hat{d}(h) \leq \frac{\epsilon'}{2} \quad (10)$$

occurs with probability at least  $1 - \frac{\delta}{2}$ .

To prove it, we can write the following chain of inequalities:

$$\mathbb{P}\left[\exists h \in \mathcal{H}, d(h) - \hat{d}(h) \geq \frac{\epsilon'}{2}\right] = \mathbb{P}\left[\bigcup_{h \in \mathcal{H}} \{d(h) - \hat{d}(h) \geq \frac{\epsilon'}{2}\}\right] \quad (11)$$

$$\leq \sum_{h \in \mathcal{H}} \mathbb{P}\left[d(h) - \hat{d}(h) \geq \frac{\epsilon'}{2}\right] \quad (12)$$

$$\leq \sum_{h \in \mathcal{H}} \exp\left[-\frac{m\epsilon'^2}{2}\right] = |\mathcal{H}| \exp\left[-\frac{m\epsilon'^2}{2}\right] \quad (13)$$

Given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log\left(\frac{2}{\delta}\right) + \log(|\mathcal{H}|) \right)$ , the right hand side is less than or equal to  $\frac{\delta}{2}$ .

Thus, given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log\left(\frac{2}{\delta}\right) + \log(|\mathcal{H}|) \right)$ , the first term and third terms of (7) are both greater than  $-\frac{\epsilon'}{2}$  with probability at least  $1 - \frac{\delta}{2}$ . Hence, with probability at least  $1 - \delta$ ,

$$\hat{d}(h) - \hat{d}(h^*) = [\hat{d}(h) - d(h)] + [d(h) - d(h^*)] + [d(h^*) - \hat{d}(h^*)] \quad (14)$$

$$\geq -\frac{\epsilon'}{2} + \epsilon' - \frac{\epsilon'}{2} = 0, \quad (15)$$

which proves the theorem.  $\square$

This theorem states that given  $m \geq \frac{2}{\epsilon^2(1-2\eta)^2} \left( \log\left(\frac{2}{\delta}\right) + \log(|\mathcal{H}|) \right)$  samples, the probability of the algorithm  $h_S = \arg \min_{h \in \mathcal{H}} \hat{d}(h)$  having  $L_D(h_S) \leq \epsilon$  is at least  $1 - \delta$ .

## References

- [MRT19] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. The MIT Press, second edition, 2019.