

---

# Theoretical Machine Learning

## Lecture 5

Instructor: Haipeng Luo

---

### 1 Erasing the loss and Finite Class Results

Recall that in the last lecture, following the roadmap for statistical learning, we derived the first two steps of relaxation for the value of online learning and arrived at an upper bound in terms of the sequential Rademacher complexity of  $\ell(\mathcal{F})$ :

$$\mathcal{V}^{\text{seq}}(\mathcal{F}, n) \leq \sup_{\mathcal{P} \in \Delta(\mathcal{Z}^n)} \mathbb{E}_{z_{1:n} \sim \mathcal{P}} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n (\mathbb{E}_{z'_t \sim \mathcal{P}(\cdot | z_{1:t-1})} [\ell(f, z'_t)] - \ell(f, z_t)) \right] \leq 2\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})).$$

Next, we further relate it to the sequential Rademacher complexity of  $\mathcal{F}$  when  $\mathcal{F}$  is a function class, and also derive a bound for the finite case.

#### 1.1 Erasing the loss

For many problems (such as online supervised learning), we have  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  and  $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ . Similarly to statistical learning, below we show that in this case we can simply ignore the loss when discussing the learnability of the problem.

**Lemma 1.** *For a binary classification problem with  $\mathcal{Z} = \mathcal{X} \times \{-1, +1\}$ ,  $\mathcal{F} \subset \{-1, +1\}^{\mathcal{X}}$ , and 0-1 loss, one has for any  $\mathcal{Z}$ -valued tree  $(\mathbf{x}, \mathbf{y})$ , there exists another  $\mathcal{X}$ -valued tree  $\mathbf{x}'$  such that*

$$\widehat{\mathcal{R}}^{\text{seq}}(\ell(\mathcal{F}); (\mathbf{x}, \mathbf{y})) = \frac{1}{2} \widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}; \mathbf{x}').$$

Therefore we have  $\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})) \leq \frac{1}{2} \mathcal{R}^{\text{seq}}(\mathcal{F})$ .

**Lemma 2 (Contraction).** *For a regression problem with  $\mathcal{Z} = \mathcal{X} \times \mathbb{R}$  and loss  $\ell(f, (x, y)) = \ell'(f(x), y)$  for some loss  $\ell'(y', y)$  that is  $G$ -Lipschitz in the first parameter, one has*

$$\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})) \leq G \mathcal{R}^{\text{seq}}(\mathcal{F}) \times \mathcal{O}(\ln^{3/2} n).$$

These lemmas are analogues of those in Lecture 2 for statistical learning, with the following differences. For [Lemma 1](#), the statistical learning analogue is  $\widehat{\mathcal{R}}^{\text{iid}}(\ell(\mathcal{F}); (x_{1:n}, y_{1:n})) = \frac{1}{2} \widehat{\mathcal{R}}^{\text{iid}}(\mathcal{F}; x_{1:n})$  for any sequence  $(x_{1:n}, y_{1:n})$ , while for online learning we have moved from a tree  $(\mathbf{x}, \mathbf{y})$  to some other tree  $\mathbf{x}'$  (the reason will be clearly shown in the proof). Nevertheless, note that this does not affect the final conclusion  $\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})) \leq \frac{1}{2} \mathcal{R}^{\text{seq}}(\mathcal{F})$ , similar to  $\mathcal{R}^{\text{iid}}(\ell(\mathcal{F})) = \frac{1}{2} \mathcal{R}^{\text{iid}}(\mathcal{F})$ . For [Lemma 2](#), the same subtlety exists, and in addition, we lose a factor of  $\mathcal{O}(\ln^{3/2} n)$  compared to the statistical learning analogue. It is not clear if this extra factor is necessary or not.

We omit the proof for [Lemma 2](#) and prove [Lemma 1](#) below.

*Proof of Lemma 1.* By definition we have

$$\begin{aligned}\widehat{\mathcal{R}}^{\text{seq}}(\ell(\mathcal{F}); (\mathbf{x}, \mathbf{y})) &= \frac{1}{n} \mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n \epsilon_t \mathbf{1} \{f(\mathbf{x}_t(\epsilon)) \neq \mathbf{y}_t(\epsilon)\} \right] \\ &= \frac{1}{n} \mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n \epsilon_t \frac{1 - \mathbf{y}_t(\epsilon) f(\mathbf{x}_t(\epsilon))}{2} \right] \\ &= \frac{1}{2n} \mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n -\epsilon_t \mathbf{y}_t(\epsilon) f(\mathbf{x}_t(\epsilon)) \right].\end{aligned}$$

We now claim that the random variables  $s_t = -\epsilon_t \mathbf{y}_t(\epsilon)$  for  $t = 1, \dots, n$  are in fact also  $n$  i.i.d. Rademacher random variables, or equivalently, the mapping  $\epsilon \rightarrow \mathbf{s} = (s_1, \dots, s_n)$  is a bijection between  $\{-1, +1\}^n$  and itself. Indeed, this is clear by constructing the inverse mapping  $\mathbf{s} \rightarrow \epsilon$  defined by  $\epsilon_t = -s_t \mathbf{y}_t(\epsilon_{1:t-1})$  (note that  $\epsilon_{1:t-1}$  can be further expressed in terms of  $\mathbf{s}$  recursively).

Based on this fact, we can construct a tree  $\mathbf{x}'$  such that  $\mathbf{x}_t(\epsilon) = \mathbf{x}'_t(\mathbf{s})$  for any  $\epsilon$  and  $t$  (note that the tree is well defined due to the bijection), and thus

$$\widehat{\mathcal{R}}^{\text{seq}}(\ell(\mathcal{F}); (\mathbf{x}, \mathbf{y})) = \frac{1}{2n} \mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n s_t f(\mathbf{x}'_t(\mathbf{s})) \right] = \frac{1}{2} \widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}; \mathbf{x}').$$

Taking sup over  $(\mathbf{x}, \mathbf{y})$  on both sides further proves  $\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})) \leq \frac{1}{2} \mathcal{R}^{\text{seq}}(\mathcal{F})$ .  $\square$

We remark that the tree  $\mathbf{x}'$  is constructed by permuting the paths of  $\mathbf{x}$  according to  $\mathbf{y}$  in some complicated way. As an illustration, consider  $\mathbf{y}$  being the tree with  $+1$  in all nodes. Then it is not hard to see that  $\mathbf{x}'$  is exactly the mirror reflection of  $\mathbf{x}$ . As another example, if  $\mathbf{y}$  has  $+1$  in the root and  $-1$  everywhere else, then  $\mathbf{x}'$  is obtained by swapping the left and right subtrees of the root of  $\mathbf{x}$ .

## 1.2 Finite class

From now on we will focus on bounding  $\mathcal{R}^{\text{seq}}(\mathcal{F})$  for some function class  $\mathcal{F}$ , starting with a finite class. The key is to apply maximal inequality again, restated below for convenience.

**Lemma 3** (Maximal Inequality). *Suppose  $\{U_f\}_{f \in \mathcal{F}}$  is a finite collection of  $\sigma$ -sub-Gaussian random variables. Then we have*

$$\mathbb{E} \left[ \max_{f \in \mathcal{F}} U_f \right] \leq \sigma \sqrt{2 \ln |\mathcal{F}|}.$$

The main result is stated below.

**Theorem 1.** *Let  $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$  be a finite class. We have for any  $\mathcal{X}$ -valued tree  $\mathbf{x}$ ,*

$$\widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}; \mathbf{x}) \leq \frac{1}{n} \sqrt{2 \left( \max_{f \in \mathcal{F}} \max_{\epsilon} \sum_{t=1}^n f^2(\mathbf{x}_t(\epsilon)) \right) \ln |\mathcal{F}|}.$$

*Consequently, if  $\mathcal{Y} \subset [-C, C]$  for some  $C > 0$ , then  $\mathcal{R}^{\text{seq}}(\mathcal{F}) \leq C \sqrt{\frac{2 \ln |\mathcal{F}|}{n}}$ .*

*Proof.* Note that  $\widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}; \mathbf{x}) = \frac{1}{n} \mathbb{E} [\max_{f \in \mathcal{F}} U_f]$  where  $U_f = \sum_{t=1}^n \epsilon_t f(\mathbf{x}_t(\epsilon))$ . Below we show that  $U_f$  is  $\sigma$ -sub-Gaussian with  $\sigma = \max_{f \in \mathcal{F}} \max_{\epsilon} \sqrt{\sum_{t=1}^n f^2(\mathbf{x}_t(\epsilon))}$ , so applying maximal inequality then finishes the proof.

Indeed, with  $U_{f,\tau} = \sum_{t=1}^{\tau} \epsilon_t f(\mathbf{x}_t(\epsilon))$  we have for any  $\lambda > 0$ ,

$$\begin{aligned}\mathbb{E} [\exp(\lambda U_{f,n})] &= \mathbb{E} [\exp(\lambda U_{f,n-1}) \mathbb{E} [\exp(\lambda \epsilon_n f(\mathbf{x}_n(\epsilon))) \mid \epsilon_{1:n-1}]] \\ &\leq \mathbb{E} [\exp(\lambda U_{f,n-1}) \exp(\frac{1}{2} \lambda^2 f^2(\mathbf{x}_n(\epsilon)))]\end{aligned}$$

where the inequality is by the fact that  $\epsilon_n f(\mathbf{x}_n(\epsilon))$  is  $|f(\mathbf{x}_n(\epsilon))|$ -sub-Gaussian. Continuing to peel the last term of  $U_{f,n-1}$  in the same way, we arrive at

$$\mathbb{E} [\exp(\lambda U_{f,n-2}) \mathbb{E} [\exp(\lambda \epsilon_{n-1} f(\mathbf{x}_{n-1}(\epsilon))) \exp(\frac{1}{2} \lambda^2 f^2(\mathbf{x}_n(\epsilon))) \mid \epsilon_{1:n-2}]],$$

but note that the term  $\exp(\frac{1}{2}\lambda^2 f^2(\mathbf{x}_n(\epsilon)))$  also involves the randomness of  $\epsilon_{n-1}$ , so we cannot directly proceed in the same way. Instead, we bound it by the worst case:

$$\begin{aligned} & \mathbb{E} \left[ \exp(\lambda U_{f,n-2}) \mathbb{E} [\exp(\lambda \epsilon_{n-1} f(\mathbf{x}_{n-1}(\epsilon))) \mid \epsilon_{1:n-2}] \max_{\epsilon_{n-1}} \exp(\frac{1}{2}\lambda^2 f^2(\mathbf{x}_n(\epsilon))) \right] \\ & \leq \mathbb{E} \left[ \exp(\lambda U_{f,n-2}) \max_{\epsilon_{n-1}} \exp(\frac{1}{2}\lambda^2 f^2(\mathbf{x}_{n-1}(\epsilon)) + \frac{1}{2}\lambda^2 f^2(\mathbf{x}_n(\epsilon))) \right] \\ & \hspace{15em} (\epsilon_{n-1} f(\mathbf{x}_{n-1}(\epsilon)) \text{ is } |f(\mathbf{x}_{n-1}(\epsilon))|\text{-sub-Gaussian)} \\ & \leq \mathbb{E} \left[ \exp(\lambda U_{f,n-2}) \max_{\epsilon_{n-2}, \epsilon_{n-1}} \exp(\frac{1}{2}\lambda^2 f^2(\mathbf{x}_{n-1}(\epsilon)) + \frac{1}{2}\lambda^2 f^2(\mathbf{x}_n(\epsilon))) \right] \end{aligned}$$

Continuing in the same fashion, we arrive at

$$\mathbb{E} [\exp(\lambda U_{f,n})] \leq \max_{\epsilon} \exp\left(\frac{\lambda^2}{2} \sum_{t=1}^n f^2(\mathbf{x}_t(\epsilon))\right) \leq \exp(\lambda^2 \sigma^2 / 2),$$

which shows that  $U_f$  is  $\sigma$ -sub-Gaussian.  $\square$

This shows that any finite class with bounded value is online learnable, and will play a key role in following development with infinite classes.

## 2 Online Binary Classification

Next we move on to discuss the learnability of infinite classes, starting from binary classification with 0-1 loss. Recall that for statistical learning, we made a key observation that even if  $\mathcal{F}$  is infinite, what really matters is the projection  $\mathcal{F}|_{\mathbf{x}_{1:n}}$ , which is always finite. Similarly, for online learning we also have

$$\widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}, \mathbf{x}) = \frac{1}{n} \mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n \epsilon_t f(\mathbf{x}_t(\epsilon)) \right] = \frac{1}{n} \mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{v} \in V} \sum_{t=1}^n \epsilon_t \mathbf{v}_t(\epsilon) \right] \leq \sqrt{\frac{2 \ln |V|}{n}} \quad (1)$$

where  $V = \mathcal{F}|_{\mathbf{x}} = \{(f \circ \mathbf{x}_1, \dots, f \circ \mathbf{x}_n) \mid f \in \mathcal{F}\}$  is the projection of  $\mathcal{F}$  onto tree  $\mathbf{x}$ , which is a set of  $\{-1, +1\}$ -valued trees. Note that  $\mathcal{F}|_{\mathbf{x}}$  is always finite, so we have yet again moved from an infinite class to a finite class. However, how large can  $|\mathcal{F}|_{\mathbf{x}}$  be? Since a tree of depth  $n$  has  $2^n - 1$  nodes, the cardinality of  $\mathcal{F}|_{\mathbf{x}}$  can be as bad as  $2^{2^n - 1}$ , leading to a very trivial bound. On the other hand, recall that in statistical learning, for a set of  $n$  samples  $\mathbf{x}_{1:n}$ ,  $|\mathcal{F}|_{\mathbf{x}_{1:n}}$  can only be at most  $2^n$ .

Since both  $2^{2^n - 1}$  and  $2^n$  are trivial bounds anyway, maybe we should just hope that  $|\mathcal{F}|_{\mathbf{x}}$  is small for common problems with a class  $\mathcal{F}$  that is not too complex? This is unfortunately not true, since  $|\mathcal{F}|_{\mathbf{x}}$  can be way too large even for a very simple class. To see this, consider the following class defined over  $\mathcal{X} = \mathbb{R}$ :

$$\mathcal{F} = \left\{ f_{\theta}(x) = \begin{cases} +1, & \text{if } x = \theta \\ -1, & \text{else} \end{cases} \mid \theta \in \mathbb{R} \right\}. \quad (2)$$

This class is intuitively simple since each classifier  $f_{\theta}$  in the class is predicting +1 for one and only one specific input  $\theta$ . Indeed, it is clear that this class cannot even shatter a set of size two, and thus  $\text{VCdim}(\mathcal{F}) = 1$ , which means it is (easily) learnable in the statistical learning setting.

However, it is easy to construct a tree such that  $|\mathcal{F}|_{\mathbf{x}} = 2^{n-1}$ , which again makes the bound in Equation (1) trivial. To show this, simply let  $\mathbf{x}$  have distinct values in all the leaves. Then  $\mathcal{F}|_{\mathbf{x}}$  at least contains  $2^{n-1}$  different trees, each of which has a different leaf with value +1.

So does this mean that  $|\mathcal{F}|_{\mathbf{x}}$  is not the right complexity measure, or is this simple class really not online learnable? It would be very unfortunate if even a class as simple as this is not online learnable. Fortunately, it turns out that this is not the case and the projection is really not the right concept to consider. To see how to fix this, note that the projection is really a set  $V$  of  $\{-1, +1\}$ -valued trees, such that

$$\forall f \in \mathcal{F}, \exists \mathbf{v} \in V, \text{ s.t. } \forall \epsilon \in \{-1, +1\}^n, f(\mathbf{x}_t(\epsilon)) = \mathbf{v}_t(\epsilon) \text{ holds for all } t = 1, \dots, n.$$

However, suppose we have a set  $V$  of  $\{-1, +1\}$ -valued trees such that a similar statement holds but importantly with two quantifiers swapped:

$$\forall f \in \mathcal{F}, \forall \epsilon \in \{-1, +1\}^n, \exists v \in V, \text{ s.t. } f(x_t(\epsilon)) = v_t(\epsilon) \text{ holds for all } t = 1, \dots, n.$$

Then this is in fact already enough for Equation (1) to hold (try to convince yourself)! A set  $V$  with the above property is called a *zero-cover* of  $\mathcal{F}|_{\mathbf{x}}$ , and the zero-covering number  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})$  is defined as the size of the smallest zero-cover. We have thus shown the following:

$$\widehat{\mathcal{R}}^{\text{seq}}(\mathcal{F}, \mathbf{x}) \leq \sqrt{\frac{2 \ln \mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})}{n}}.$$

So how large can  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})$  be then? First of all, this is clearly always not larger than  $|\mathcal{F}|_{\mathbf{x}}$  (since  $\mathcal{F}|_{\mathbf{x}}$  is a zero-cover of itself). Second,  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})$  is in fact always bounded by  $2^n$ . This is because the set of all the possible trees with the same value at each level is always a zero-cover for any class, and there are clearly  $2^n$  such trees (since each level takes one of the two possible values). This is of course still a trivial bound, but it is at least the same trivial bound as the one for a projection in statistical learning, indicating that this might be the right complexity measure.

For a class with specific structures,  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})$  can be much smaller. For example, the simple class defined in Equation (2) has zero-covering number  $n + 1$ , implying that it is online learnable (as we hope). We defer the formal proof to the next subsection, but illustrate with a simpler case when  $\mathbf{x}$  contains no identical value along any path. In this case we only need the following  $n + 1$  trees to cover  $\mathcal{F}|_{\mathbf{x}}$ : a tree with  $-1$  in every node, and for each  $t = 1, \dots, n$ , a tree with  $+1$  for all nodes at level  $t$  and  $-1$  everywhere else.

## 2.1 Combinatorial parameter

Similarly to previous discussions for statistical learning, next we will introduce some combinatorial parameter that provides a reasonable upper bound for  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}})$  and that is easier to bound. First, we say that  $\mathcal{F}$  shatters a tree  $\mathbf{x}$  if for any  $\epsilon \in \{-1, +1\}^n$ , there exists  $f \in \mathcal{F}$  such that  $f(x_t(\epsilon)) = \epsilon_t$  holds for all  $t = 1, \dots, n$ . Now, we define the *Littlestone dimension*  $\text{Ldim}(\mathcal{F})$  of a class  $\mathcal{F}$  as the depth of the largest tree that can be shattered by  $\mathcal{F}$  ( $\text{Ldim}(\mathcal{F})$  is defined as 0 if no tree can be shattered by  $\mathcal{F}$ , and  $\infty$  if for any  $n$  there exists a tree of depth  $n$  that is shattered by  $\mathcal{F}$ ).

Note that one always has  $\text{VCdim}(\mathcal{F}) \leq \text{Ldim}(\mathcal{F})$ , since if  $x_{1:n}$  is shattered by  $\mathcal{F}$ , then a tree  $\mathbf{x}$  of depth  $n$  such that all nodes at level  $t$  have value  $x_t$  for  $t = 1, \dots, n$  is clearly shattered by  $\mathcal{F}$  as well. Using this fact, one sees that  $\text{Ldim}(\mathcal{F}) = 0$  implies  $\text{VCdim}(\mathcal{F}) = 0$  and thus  $\mathcal{F}$  contains only one function.

Similar to VC dimension, to prove  $\text{Ldim}(\mathcal{F}) = d$  we have to 1) construct a tree of depth  $d$  that can be shattered by  $\mathcal{F}$  and 2) prove that no tree of depth  $n + 1$  can be shattered by  $\mathcal{F}$ . As an example, we argue that the simple class defined by Equation (2) has Littlestone dimension exactly 1 (the same as its VC dimension). Clearly it can shatter a tree with depth 1 (in fact, any tree with depth 1). On the other hand, it cannot shatter any tree with depth 2 (just consider  $\epsilon = (+1, +1)$  and  $\epsilon = (+1, -1)$ ).

The following theorem is an exact analogue of Sauer's lemma and it provides a bound on the zero-covering number in terms of the Littlestone dimension. Its proof also reveals how to construct a zero-cover recursively.

**Theorem 2.** *Suppose  $\mathbf{x}$  is any  $\mathcal{X}$ -valued tree with depth  $n$  and  $\mathcal{F} \subset \{-1, +1\}^{\mathcal{X}}$  has Littlestone dimension  $d \leq n$ , then*

$$\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}}) \leq \sum_{i=0}^d \binom{n}{i} \leq \left(\frac{en}{d}\right)^d.$$

*Proof.* Let  $g(d, n) = \sum_{i=0}^d \binom{n}{i}$ . We will prove  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}}) \leq g(d, n)$  via induction on the value of  $d + n$  (the second inequality has been proven in Lecture 2). The base case  $d + n = 1$  is trivial since the only configuration is  $d = 0$  and  $n = 1$ , in which case  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}}) = 1$  clearly. Now we assume that the statement holds for any  $n' > d'$  such that  $n' + d' < n + d$ , and prove  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}}) \leq g(d, n)$ . The case when  $d = 0$  is again trivial, so we assume  $d > 0$ .

First, we provide a way to recursively construct a zero-cover for  $\mathcal{F}|_{\mathbf{x}}$ . Depending on the prediction for the root of  $\mathbf{x}$ , we split the class  $\mathcal{F}$  into two subclasses:  $\mathcal{F}_- = \{f \in \mathcal{F} \mid f(x_1) = -1\}$  and

$\mathcal{F}_+ = \{f \in \mathcal{F} \mid f(\mathbf{x}_1) = -1\}$ . Let  $\mathbf{x}^\ell$  and  $\mathbf{x}^r$  be the left and right subtrees of the root of  $\mathbf{x}$  (which have depth  $n - 1$ ), and  $V_+^\ell$  and  $V_+^r$  be the smallest zero-cover of  $\mathcal{F}_+|_{\mathbf{x}^\ell}$  and  $\mathcal{F}_+|_{\mathbf{x}^r}$  respectively. Now we construct a set  $V_+$  in the following way: 1) the root of every tree in  $V_+$  has value  $+1$ , 2) pair the element from  $V_+^\ell$  and  $V_+^r$  to form the left and right subtrees of the root for trees in  $V_+$ , so that each element from  $V_+^\ell$  and  $V_+^r$  appears at least once. It is clear that this can be done such that  $|V_+| = \max\{|V_+^\ell|, |V_+^r|\}$ . Moreover, it is also clear that  $V_+$  is a zero-cover of  $\mathcal{F}_+|_{\mathbf{x}}$ . In the exact same way, we construct  $V_-$  such that  $|V_-| = \max\{|V_-^\ell|, |V_-^r|\}$  and  $V_-$  is a zero-cover of  $\mathcal{F}_-|_{\mathbf{x}}$ . Finally, we have that  $V_- \cup V_+$  is a zero-cover of  $\mathcal{F}|_{\mathbf{x}}$ .

It remains to bound  $|V_-|$  and  $|V_+|$ . The key observation is that it is impossible that  $\mathcal{F}_-$  and  $\mathcal{F}_+$  both have Littlestone dimension  $d$ . Otherwise, there are trees  $\mathbf{x}^-$  and  $\mathbf{x}^+$  of depth  $n$  that can be shattered by  $\mathcal{F}_-$  and  $\mathcal{F}_+$  respectively. By pairing  $\mathbf{x}^-$  and  $\mathbf{x}^+$  as the left and right subtrees of the root  $\mathbf{x}_1$ , we obtain a tree with depth  $n + 1$  that can be shattered by  $\mathcal{F}$ , which is a contradiction to  $\text{Ldim}(\mathcal{F}) = d$ . Without loss of generality, we can thus assume  $\mathcal{F}_-$  has Littlestone dimension at most  $d - 1$ . Using the inductive hypothesis, we thus have

$$|V_+| = \max\{|V_+^\ell|, |V_+^r|\} = \max\{\mathcal{N}_0(\mathcal{F}_+|_{\mathbf{x}^\ell}), \mathcal{N}_0(\mathcal{F}_+|_{\mathbf{x}^r})\} \leq g(d, n - 1),$$

and

$$|V_-| = \max\{|V_-^\ell|, |V_-^r|\} = \max\{\mathcal{N}_0(\mathcal{F}_-|_{\mathbf{x}^\ell}), \mathcal{N}_0(\mathcal{F}_-|_{\mathbf{x}^r})\} \leq g(d - 1, n - 1).$$

Therefore,  $\mathcal{N}_0(\mathcal{F}|_{\mathbf{x}}) \leq |V| = |V_-| + |V_+| \leq g(d - 1, n - 1) + g(d, n - 1) = g(d, n)$  (the last step is proven in Lecture 2). This finishes the proof.  $\square$

We remark that the concept of zero-covering is the key to allow a construction with  $|V_+| = \max\{|V_+^\ell|, |V_+^r|\}$ . If we focus on the projection instead, we would have arrived at something like  $|V_+| = |V_+^\ell| \times |V_+^r|$ . Applying this theorem directly, we conclude that the zero-covering number of the simple class defined by Equation (2) is indeed bounded by  $g(1, n) = n + 1$  (even if the tree contains identical elements in some paths). The proof of Theorem 2 also reveals a recursive way to construct such a cover with size  $n + 1$  (try to construct explicitly a zero-cover for a small tree of depth say 3).

The following bound on the value of the game is a direct corollary based on previous discussions.

**Corollary 1.** *For any class of binary classifier  $\mathcal{F}$  with  $d = \text{Ldim}(\mathcal{F})$ , we have*

$$\mathcal{V}^{\text{seq}}(\mathcal{F}, n) \leq \mathcal{R}^{\text{seq}}(\mathcal{F}) \leq \sqrt{\frac{2d \ln\left(\frac{en}{d}\right)}{n}}.$$

So finite Littlestone dimension is sufficient for online learnability. It turns out that it is also *necessary* for online learnability, indicating that this sequence of upper bounding is tight.

**Theorem 3.** *If  $\text{Ldim}(\mathcal{F}) = \infty$ , then for any algorithm and any  $n$ , there exists an environment such that the expected average regret of this algorithm is at least  $1/2$ .*

*Proof.* Let  $\mathbf{x}$  be a tree of depth  $n$  shattered by  $\mathcal{F}$  (which always exists since  $\text{Ldim}(\mathcal{F}) = \infty$ ), and  $\epsilon_1, \dots, \epsilon_n$  be i.i.d. Rademacher random variables. At time  $t$ , the environment chooses  $(\mathbf{x}_t(\epsilon), \epsilon_t)$ . In such an environment, no matter what the algorithm is (proper or improper), its expected total loss is always exactly  $n/2$ . On the other hand, by the definition of shattering, there is always an  $f \in \mathcal{F}$  with perfect prediction on this dataset, which means that the expected regret is at least  $n/2$ .  $\square$

In fact, in HW2 you will prove an even stronger statement (with  $d = \text{Ldim}(\mathcal{F})$ ):

$$\mathcal{V}^{\text{seq}}(\mathcal{F}, n) \geq \sqrt{\frac{1}{8} \min\left\{\frac{d}{n}, 1\right\}},$$

further showing that the upper bound we obtain is very tight. Closing the gap  $\ln\left(\frac{en}{d}\right)$  between the upper and lower bounds remains open.

**Summary.** Combining all steps, we have shown for binary classification problems

$$\begin{aligned} \sqrt{\frac{1}{8} \min \left\{ \frac{d}{n}, 1 \right\}} &\leq \mathcal{V}^{\text{seq}}(\mathcal{F}, n) \\ &\leq \sup_{\mathcal{P} \in \Delta(\mathcal{Z}^n)} \mathbb{E}_{z_{1:n} \sim \mathcal{P}} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n (\mathbb{E}_{z'_t \sim \mathcal{P}(\cdot | z_{1:t-1})} [\ell(f, z'_t)] - \ell(f, z_t)) \right] \\ &\leq 2\mathcal{R}^{\text{seq}}(\ell(\mathcal{F})) \leq \mathcal{R}^{\text{seq}}(\mathcal{F}) \leq \sup_{\mathbf{x}} \sqrt{\frac{2 \ln \mathcal{N}_0(\mathcal{F} | \mathbf{x})}{n}} \leq \sqrt{\frac{2d \ln \left( \frac{en}{d} \right)}{n}} \end{aligned}$$

where  $d = \text{Ldim}(\mathcal{F})$ .

## 2.2 Online learning is strictly harder

In Lecture 1, via the online-to-batch conversion we showed that online learning is at least as hard as statistical learning. Is it strictly harder? The example of the simple class defined by Equation (2) does not indicate that because the VC dimension and the Littlestone dimension coincide (both are 1). Instead, let's consider the threshold function class defined over  $\mathcal{X} = \mathbb{R}$ :

$$\mathcal{F} = \left\{ f_{\theta}(x) = \begin{cases} +1 & \text{if } x \leq \theta \\ -1 & \text{else} \end{cases} \mid \theta \in \mathbb{R} \right\}, \quad (3)$$

which has VC dimension exactly 1 as discussed in Lecture 2. It turns out that the Littlestone dimension of this seemingly simple class is infinity!

**Proposition 1.** *The Littlestone dimension of the threshold function class (Equation (3)) is  $\infty$ .*

*Proof.* To see this, consider an infinite  $[0, 1]$ -valued tree  $\mathbf{x}$  with root being  $1/2$ , and the left child and right child of a node with value  $a$  at level  $t$  being  $a - \frac{1}{2^{t+1}}$  and  $a + \frac{1}{2^{t+1}}$  respectively. (This is much easier to interpret if you draw a picture.)

Now for any  $n$  and any path/labeling  $\epsilon$ , let  $\theta_1$  be the last node of this path  $\mathbf{x}_n(\epsilon)$  and  $\theta_2$  be the last node on this path with label  $-\epsilon_n$ . Then the claim is that any value  $\theta$  in between  $\theta_1$  and  $\theta_2$  satisfies:  $f_{\theta}(\mathbf{x}_t(\epsilon)) = \epsilon_t$  for all  $t = 1, \dots, n$ , and thus  $\mathcal{F}$  shatters this tree. Indeed, note that the tree is constructed such that if  $\epsilon_t = +1$ , then every node in the path below level  $t$  has value larger than the node at level  $t$ . Similarly, if  $\epsilon_t = -1$ , then every node in the path below level  $t$  has value smaller than the node at level  $t$ . Therefore, suppose  $\epsilon_n = +1$ , then all the nodes with label  $+1$  on the path must have a value smaller than  $\theta_1$ , and all the nodes with label  $-1$  on the path must have a value larger than  $\theta_2$ , and thus  $f_{\theta}$  predicts all the labels correctly if  $\theta \in [\theta_1, \theta_2]$ . The case for  $\epsilon_n = -1$  is similar.  $\square$

Based on previous discussions, we conclude that this simple class is learnable in the statistical learning setting, but not learnable in the online setting. More generally, it is clear that the class of linear classifiers

$$\mathcal{F} = \{ f_{\theta, b}(x) = \text{sign}(\langle x, \theta \rangle + b) \mid \theta \in \mathbb{R}^d, b \in \mathbb{R} \}$$

also has infinite Littlestone dimension (since it subsumes the threshold class), while having a finite VC dimension  $d + 1$ . This illustrates that online learning is not just as hard as statistical learning, but is in fact *strictly* harder than statistical learning.

So if even learning linear classifiers is impossible, is online learning just too hard to be meaningful? The answer is yes in some sense for online classification, or put differently the 0-1 loss is too hard for online learning beyond finite classes. However, next time we will show that for online regression, many more possibilities open up.