## STATISTICAL MACHINE LEARNING (WS2025/26) SEMINAR: VC DIMENSION

**Assignment 1.** Let  $\mathcal{H} = \{h_1, \dots, h_H\}$  be a finite hypothesis class. Prove that the Uniform Law of Large Numbers (ULLN) holds for  $\mathcal{H}$ . Specifically, show that there exists a function  $m_{\mathrm{ul}}^{\mathcal{H}}:(0,1)\times(0,1)\to\mathbb{N}$  such that for any i.i.d. sample  $T_m=\big((x_i,y_i)\in\mathcal{X}\times\mathcal{Y}\mid i=1,2,\ldots,n\}$  $1,\ldots,m$ ) of size  $m\geq m_{\rm ul}^{\mathcal{H}}(\varepsilon,\delta)$  drawn from the distribution p(x,y), the following holds:

$$\mathbb{P}\left(\max_{h\in\mathcal{H}}\left|\hat{R}(T_m,h)-R(p,h)\right|>\varepsilon\right)\leq\delta\,,$$

where

$$R(p,h) = \mathbb{E}_{(x,y)\sim p} \big[\ell(y,h(x))\big] \quad \text{and} \quad \hat{R}(T_m,h) = \frac{1}{m} \sum_{i=1}^m \ell\big(y_i,h(x_i)\big)$$

denote the true and empirical risks, respectively. Assume the 0-1 loss function  $\ell(y,y')=[y\neq 0]$ y'.

## **Solution 1.**

$$\mathbb{P}\Big(\max_{h\in\mathcal{H}}\big|R(p,h)-\hat{R}(T_m,h)\big|\geq\varepsilon\Big)\quad \stackrel{(1)}{=}\quad \mathbb{P}\left(\begin{array}{c} \big|R(p,h^1)-\hat{R}(T_m,h^1)\big|\geq\varepsilon\quad\text{or}\\ \big|R(p,h^2)-\hat{R}(T_m,h^2)\big|\geq\varepsilon\quad\text{or}\\ \vdots\\ \big|R(p,h^H)-\hat{R}(T_m,h^H)\big|\geq\varepsilon\end{array}\right)$$

$$\stackrel{(2)}{\leq}\quad \sum_{h\in\mathcal{H}}\mathbb{P}\Big(\big|R(p,h)-\hat{R}(T_m,h)\big|\geq\varepsilon\Big)$$

$$\stackrel{(3)}{\leq}\quad 2\,|\mathcal{H}|\,e^{-2m\varepsilon^2}$$

- (1)  $a \ge \varepsilon$  or  $b \ge \varepsilon \iff \max\{a, b\} \ge \varepsilon$ (2) Union bound:  $\mathbb{P}(A_1 \text{ or } A_2 \text{ or } \cdots \text{ or } A_n) \le \sum_{i=1}^n \mathbb{P}(A_i)$
- (3) Hoeffding inequality:  $\mathbb{P}(|R(h) R_{\mathcal{T}^m}(h)| \geq \varepsilon) \leq 2e^{-2m\varepsilon^2}$

By setting  $2|\mathcal{H}|e^{-2m\varepsilon^2} = \delta$  and solving for m, we get

$$m_{\mathrm{ul}}^{\mathcal{H}}(\varepsilon, \delta) = \frac{1}{2\varepsilon^2} \log \left( \frac{2|\mathcal{H}|}{\delta} \right) \qquad \Rightarrow \qquad \mathbb{P}\left( \max_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right| \ge \varepsilon \right) \le \delta$$

**Assignment 2.** Let us consider the class of linear classifiers mapping  $x \in \mathbb{R}^d$  to  $\{-1, +1\}$ , that is

$$\mathcal{H} = \left\{ h(\boldsymbol{x}; \boldsymbol{w}, b) = \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b) \mid (\boldsymbol{w}, b) \in (\mathbb{R}^d \times \mathbb{R}) \right\}.$$

Show that the VC dimension of  $\mathcal{H}$  is d+1.

**Solution 2.** The proof has two steps:

- (1) Show that  $d_{VC}(\mathcal{H}) \geq d+1$ . We construct a set of d+1 input points and show that they can be shattered.
- (2) Show that  $d_{VC}(\mathcal{H}) < d+2$ . We show that no set of d+2 input points can be shattered.
- 1) Lower bound:  $d_{VC}(\mathcal{H}) \geq d+1$ . We construct d+1 points in  $\mathbb{R}^d$  that can be shattered by linear classifiers. Let the first d points be the standard basis vectors:

$$\boldsymbol{x}_i = [0, \dots, \underbrace{1}_{i\text{-th coordinate}}, \dots, 0], \quad \forall i \in 1, \dots, d.$$

Let the last point be the origin,

$$x_{d+1} = 0.$$

Fix an arbitrary sequence of labels  $(y_1, y_2, \dots, y_{d+1}) \in \{-1, +1\}^{d+1}$ . We construct the parameters of the linear classifier as

$$\mathbf{w} = [y_1, y_2, \dots, y_d], \qquad b = \frac{1}{2}y_{d+1}.$$

Now verify that the classifier

$$h(\boldsymbol{x}) = \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b)$$

correctly predicts all d+1 labels. For  $i \in 1, \ldots, d$ ,

$$h(\boldsymbol{x}_i) = \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b) = \operatorname{sign}(y_i + \frac{1}{2}y_{d+1}) = y_i.$$

For the last point,

$$h(\boldsymbol{x}_{d+1}) = \operatorname{sign}(\langle \boldsymbol{w}, \boldsymbol{x}_{d+1} \rangle + b) = \operatorname{sign}(\frac{1}{2}y_{d+1}) = y_{d+1}.$$

Thus, the constructed classifier realizes any labeling of these d+1 points, proving that the set is shattered. Therefore,

$$d_{\text{VC}}(\mathcal{H}) \ge d + 1.$$

**2)** Upper bound:  $d_{VC}(\mathcal{H}) < d+2$ . We now show that no set of d+2 points in  $\mathbb{R}^d$  can be shattered by  $\mathcal{H}$ . Let  $\{x_1,\ldots,x_{d+2}\}\subset\mathbb{R}^d$  be arbitrary. We lift these points into  $\mathbb{R}^{d+1}$  by appending a constant coordinate:

$$z_i = [x_i, 1], \quad \forall i \in 1, \dots, d+2.$$

We also represent affine classifiers as homogeneous ones:

$$\hat{h}(\boldsymbol{z}) = \operatorname{sign}(\langle \hat{\boldsymbol{w}}, \boldsymbol{z} \rangle), \quad \text{where } \hat{\boldsymbol{w}} = [\boldsymbol{w}, b].$$

Clearly,  $h(x_i) = \hat{h}(z_i)$  for all i. Since we have d+2 points in  $\mathbb{R}^{d+1}$ , they must be linearly dependent. Hence, there exist coefficients  $(a_1, \ldots, a_{d+2})$ , not all zero, such that

$$\sum_{i=1}^{d+2} a_i oldsymbol{z}_i = oldsymbol{0}.$$

Without loss of generality, assume  $a_1 \neq 0$  and express

$$oldsymbol{z}_1 = \sum_{i=2}^{d+2} a_i' oldsymbol{z}_i, \quad ext{where } a_i' = -rac{a_i}{a_1}.$$

Now define labels for the last d+1 points as

$$y_i = \operatorname{sign}(a_i'), \quad \forall i \in \{2, \dots, d+2\}.$$

Assume that there exists a classifier  $\hat{h}$  correctly predicting these labels, i.e.

$$\hat{h}(\boldsymbol{z}_i) = y_i \iff \operatorname{sign}(\langle \hat{\boldsymbol{w}}, \boldsymbol{z}_i \rangle) = \operatorname{sign}(a_i'), \quad \forall i \in 2, \dots, d+2.$$

Then, for all  $i \geq 2$ , we have  $a'_i \langle \hat{\mathbf{w}}, \mathbf{z}_i \rangle \geq 0$ , and for at least one i, this product is strictly positive. Now, consider the prediction for  $\mathbf{z}_1$ :

$$\langle \hat{m{w}}, m{z}_1 
angle = \sum_{i=2}^{d+2} a_i' \langle \hat{m{w}}, m{z}_i 
angle.$$

Because each term  $a'_i \langle \hat{\boldsymbol{w}}, \boldsymbol{z}_i \rangle \geq 0$ , the sum is non-negative and strictly positive unless all are zero. Hence,

$$\hat{h}(\boldsymbol{z}_1) = \operatorname{sign}(\langle \hat{\boldsymbol{w}}, \boldsymbol{z}_1 \rangle) = +1.$$

But if we assign the label  $y_1 = -1$ , no classifier can realize this labeling, because  $z_1$  is forced into the positive class. Thus, the set of d+2 points cannot be shattered by linear classifiers, implying

$$d_{VC}(\mathcal{H}) < d+2.$$

Assignment 3. Consider a hypothesis space of classifiers

$$\mathcal{H} = \left\{ h(x; a) = \operatorname{sign}(\sin(ax)) \mid a \in \mathbb{R} \right\}.$$

That is, each  $h \in \mathcal{H}$  is determined by a single parameter  $a \in \mathbb{R}$  and it maps real valued input  $x \in \mathbb{R}$  to a set of hidden labels  $\{+1, -1\}$  based on the sign of the score  $\sin(xa)$ . Show that the VC dimension of  $\mathcal{H}$  is infinite.

Hint: Show that for arbitrary set of labels  $\{y^i \in \{+1, -1\} \mid i = 1, ..., m\}$  the inputs  $\{x^i = 10^{-i} \mid i = 1, ..., m\}$  can be predicted correctly by h(x; a) with

$$a = \pi \left( 1 + \frac{1}{2} \sum_{i=1}^{m} (1 - y^i) 10^i \right)$$

**Solution 3.** The plan of the proof is to construct points  $(x_1, x_2, \ldots, x_m)$  such that for arbitrary labels  $(y_1, y_2, \ldots, y_m) \in \{-1, +1\}^m$ , there will be a, which we also construct, such that all the points will be correctly classified by  $h(x) = \sin(ax)$ , i.e. the set of points is shattered. Thus, we conclude that the VC dimension of  $\mathcal{H}$  is  $\infty$ .

(1) Set the m points to be

$$x_i = 10^{-i}, \quad i \in \{1, \dots, m\}.$$

and fix an arbitrary sequence of the labels  $(y_1, y_2, \dots, y_m) \in \{-1, +1\}^m$ .

(2) Let the classifier parameter be:

$$a = \pi + \pi \sum_{i=1}^{m} \hat{y}_i \, 10^j$$

where

$$\hat{y}_i = \frac{1 - y_i}{2} \in \{0, 1\} \ .$$

are auxiliary labels.

(3) We will use the indetity

$$\sin\left(\pi(k+t)\right) = (-1)^k \sin(\pi t),\,$$

which is valid for every integer  $k \in \mathbb{N}$ .

(4) Let us rewrite the value of the discriminant function at  $x_i$ :

$$a x_{i} = \sin \left( \frac{\pi + \pi \sum_{j=1}^{m} \hat{y}_{j} 10^{j}}{10^{i}} \right)$$

$$= \sin \left( \pi \left[ 10^{-i} + \sum_{j=1}^{i-1} \hat{y}_{j} 10^{j-i} + \hat{y}_{i} + \sum_{j=i+1}^{m} \hat{y}_{j} 10^{j-i} \right] \right)$$

$$= \sin \left( \pi \sum_{j=i+1}^{m} \hat{y}_{j} 10^{j-i} + \pi \left( 10^{-i} + \sum_{j=1}^{i-1} \hat{y}_{j} 10^{j-i} + \hat{y}_{i} \right) \right)$$

$$= \sin \left( \pi \left( 10^{-i} + \sum_{j=1}^{i-1} \hat{y}_{j} 10^{j-i} + \hat{y}_{i} \right) \right)$$

$$= \left\{ \sin(\pi \delta) \text{ if } \hat{y}_{i} = +1 \\ \sin(\pi \delta + \pi) \text{ if } \hat{y}_{i} = -1 \right\}$$

$$= \operatorname{sign}(y_{i})$$