

# CS 446 MJT — Homework 4

*haoyuan9*

Version 2

## Instructions.

- Homework is due **Tuesday, April 2, at 11:59pm**; no late homework accepted.
- Everyone must submit individually at gradescope under **hw4**. (There is no **hw4code**!)
- The “written” submission at **hw4 must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L<sup>A</sup>T<sub>E</sub>X, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw4**, gradescope will ask you to mark out boxes around each of your answers; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.

### 1. VC dimension.

This problem will show that two different classes of predictors have infinite VC dimension.

**Hint:** to prove infinite  $\text{VC}(\mathcal{H}) = \infty$ , it is usually most convenient to show  $\text{VC}(\mathcal{H}) \geq n$  for all  $n$ .

- (a) Let  $\mathcal{F} := \{\mathbf{x} \mapsto 2 \cdot \mathbf{1}[\mathbf{x} \in C] - 1 : C \subseteq \mathbb{R}^d \text{ is convex}\}$  denote the set of all classifiers whose decision boundary is a convex subset of  $\mathbb{R}^d$  for  $d \geq 2$ . Prove  $\text{VC}(\mathcal{F}) = \infty$ .

**Hint:** Consider data examples on the unit sphere  $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ .

- (b) Given  $x \in \mathbb{R}$ , let  $\text{sgn}$  denote the sign of  $x$ :  $\text{sgn}(x) = 1$  if  $x \geq 0$  while  $\text{sgn}(x) = -1$  if  $x < 0$ .  
Let  $\sigma > 0$  be given, and define  $\mathcal{G}_\sigma$  to be the set of (sign of) all RBF classifiers with bandwidth  $\sigma$ , meaning

$$\mathcal{G}_\sigma := \left\{ \mathbf{x} \mapsto \text{sgn} \left( \sum_{i=1}^m \alpha_i \exp \left( -\|\mathbf{x} - \mathbf{x}_i\|^2 / (2\sigma^2) \right) \right) : m \in \mathbb{Z}_{\geq 0}, \mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d, \boldsymbol{\alpha} \in \mathbb{R}^m \right\}.$$

Prove  $\text{VC}(\mathcal{G}_\sigma) = \infty$ .

**Remark:** the sign of 0 is not important: you have the freedom to choose some nice data examples and avoid this case.

**Hint:** remember in hw3 it is proved that if  $\sigma$  is small enough, the RBF kernel SVM is close to the 1-nearest neighbor predictor. In this problem,  $\sigma$  is fixed, but you have the freedom to choose the data examples. If the distance between data examples is large enough, the RBF kernel SVM could still be close to the 1-nearest neighbor predictor. Make sure to have an explicit construction of such a dataset.

**Solution.** (Your solution here.)

- (a) For any  $d$ , choose  $n$  examples on the unit sphere  $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ . Then for any labelings, we can always choose one set of data samples  $A$  and define  $C$  as the minimum convex hull of  $A$ . So  $\mathcal{F}$  can shatter any  $n$  samples on the unit sphere. Hence,  $\text{VC}(\mathcal{F}) \geq n$  for all  $n$ . So we have  $\text{VC}(\mathcal{F}) = \infty$ .
- (b) Choose  $n$  data samples, where  $\|\mathbf{x}_{1i} - \mathbf{x}_{1j}\|^2 = \infty$  for any pair of  $i, j$  where  $i \neq j$ . Then construct  $\mathcal{G}_\sigma$  by making  $m = n$ , and choosing  $n$   $\mathbf{x}_i$  where  $\|\mathbf{x}_{2i} - \mathbf{x}_{1i}\|^2 = c$ .  $c$  is a constant.  
Just as HW3, we can devided the indicator function by  $\exp(-\rho^2/2\sigma^2)$  without changing the sign, where  $\rho = \|\mathbf{x}_{2i} - \mathbf{x}_{1i}\|^2$ . Then the classifier becomes

$$\mathcal{G}_\sigma := \left\{ \mathbf{x} \mapsto \text{sgn} \left( \frac{\sum_{i=1}^m \alpha_i \exp \left( -\|\mathbf{x} - \mathbf{x}_{2i}\|^2 / (2\sigma^2) \right)}{\exp \left( -\rho^2 / 2\sigma^2 \right)} \right) : m = n, \mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d, \boldsymbol{\alpha} \in \mathbb{R}^m \right\}.$$

Since  $\|\mathbf{x}_{1i} - \mathbf{x}_{1j}\|^2 = \infty$ , we have  $\|\mathbf{x}_{2i} - \mathbf{x}_{1j}\|^2 = \infty$ . Hence

$$\mathbf{x}_{1i} \mapsto \text{sgn}(\alpha_i)$$

Hence, for any labelings for a data sample of size  $n$ , we can choose  $\alpha_i$  the same sign as  $x_{1i}$ . So  $\text{VC}(\mathcal{G}_\sigma) \geq n$  and  $\text{VC}(\mathcal{G}_\sigma) = \infty$ .

## 2. Rademacher complexity of linear predictors.

Let examples  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  be given with  $\|\mathbf{x}_i\| \leq R$ , along with linear functions  $\{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} : \|\mathbf{w}\| \leq W\}$ . The goal in this problem is to show  $\text{Rad}(\mathcal{F}) \leq RW/\sqrt{n}$ .

- (a) For a fixed sign vector  $\varepsilon \in \{-1, +1\}^n$ , define  $\mathbf{x}_\varepsilon := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i$ . Show

$$\max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) \leq W \|\mathbf{x}_\varepsilon\|.$$

**Hint:** Cauchy-Schwarz!

- (b) Show  $\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 \leq R^2/n$ .  
(c) Now combine the pieces to show  $\text{Rad}(\mathcal{F}) \leq RW/\sqrt{n}$ .

**Hint:** one missing piece is to write  $\|\cdot\| = \sqrt{\|\cdot\|^2}$  and use Jensen's inequality.

**Solution.** (*Your solution here.*)

- (a)

$$\begin{aligned} \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) &= \mathbf{w}_{\text{optim}}^\top \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \\ &= \mathbf{w}_{\text{optim}}^\top \mathbf{x}_\varepsilon \\ &\leq \|\mathbf{w}_{\text{optim}}\| \|\mathbf{x}_\varepsilon\| \\ &\leq W \|\mathbf{x}_\varepsilon\| \end{aligned}$$

- (b)

$$\begin{aligned} \mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 &= \mathbb{E}_\varepsilon \sum_j \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{ij} \varepsilon_i \right)^2 \\ &= \frac{1}{n^2} \sum_j \mathbb{E}_\varepsilon \left( \sum_{i=1}^n \mathbf{x}_{ij} \varepsilon_i \right)^2 \\ &= \frac{1}{n^2} \sum_j \left( \text{Var} \left( \sum_{i=1}^n \mathbf{x}_{ij} \varepsilon_i \right) + \mathbb{E}_\varepsilon^2 \left( \sum_{i=1}^n \mathbf{x}_{ij} \varepsilon_i \right) \right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \sum_j \mathbf{x}_{ij}^2 \text{Var}(\varepsilon_i) \\ &= \frac{1}{n^2} \sum_{i=1}^n \|\mathbf{x}_i\|^2 \\ &\leq \frac{1}{n^2} n R^2 \\ &= \frac{R^2}{n} \end{aligned}$$

(c)

$$\begin{aligned}\text{Rad}(\mathcal{F}) &= \mathbb{E}_\varepsilon \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) \\ &\leq \mathbb{E}_\varepsilon W \|\mathbf{x}_\varepsilon\| \\ &= W \sqrt{(\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|)^2} \\ &\leq W \sqrt{\mathbb{E}_\varepsilon (\|\mathbf{x}_\varepsilon\|^2)} \\ &= W \sqrt{\frac{R^2}{n}} \\ &= \frac{WR}{\sqrt{n}}\end{aligned}$$

### 3. Generalization bounds for a few linear predictors.

In this problem, it is always assumed that for any  $(\mathbf{x}, y)$  sampled from the distribution,  $\|\mathbf{x}\| \leq R$  and  $y \in \{-1, +1\}$ .

Consider the following version of the soft-margin SVM:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \left[1 - \mathbf{w}^\top \mathbf{x}_i y_i\right]_+ = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \widehat{\mathcal{R}}_{\text{hinge}}(\mathbf{w}).$$

Let  $\hat{\mathbf{w}}$  denote the (unique!) optimal solution, and  $\hat{f}(\mathbf{x}) = \hat{\mathbf{w}}^\top \mathbf{x}$ .

Prove that for any regularization level  $\lambda > 0$ , with probability at least  $1 - \delta$ , it holds that

$$\mathcal{R}(\hat{f}) \leq \widehat{\mathcal{R}}(\hat{f}) + R \sqrt{\frac{8}{\lambda n}} + 3 \left(1 + R \sqrt{2/\lambda}\right) \sqrt{\frac{\ln(2/\delta)}{2n}}.$$

**Hint:** use the fact from slide 5/61 of the first ML Theory lecture that  $\|\hat{\mathbf{w}}\| \leq \sqrt{2/\lambda}$ , the linear predictor Rademacher complexity bound from the previous problem, and the Rademacher generalization theorem on slide 57 of the final theory lecture.

**Solution.** (*Your solution here.*)

1. First, let's prove Hinge Loss is 1-Lipschitz.

For any  $\mathbf{w}_1$  and  $\mathbf{w}_2$ , if  $\mathbf{w}_1 > 1$ , then

$$\begin{aligned} |[1 - \mathbf{w}_1^\top xy]_+ - [1 - \mathbf{w}_2^\top xy]_+| &= |[1 - \mathbf{w}_2^\top xy]_+| \\ &\leq |[\mathbf{w}_1^\top xy - \mathbf{w}_2^\top xy]_+| \\ &\leq |\mathbf{w}_1^\top xy - \mathbf{w}_2^\top xy| \end{aligned}$$

Similarly, we can prove that when  $\mathbf{w}_2 > 1$ ,  $|[1 - \mathbf{w}_1^\top xy]_+ - [1 - \mathbf{w}_2^\top xy]_+| \leq |\mathbf{w}_1^\top xy - \mathbf{w}_2^\top xy|$

When  $\mathbf{w}_1 \leq 1, \mathbf{w}_2 \leq 1$

$$\begin{aligned} |[1 - \mathbf{w}_1^\top xy]_+ - [1 - \mathbf{w}_2^\top xy]_+| &= |1 - \mathbf{w}_1^\top xy - 1 + \mathbf{w}_2^\top xy| \\ &= |\mathbf{w}_1^\top xy - \mathbf{w}_2^\top xy| \end{aligned}$$

Hence, Hinge Loss using affine  $\mathbf{f}$  is 1-Lipschitz.

In addition, since  $\|\mathbf{x}\| \leq R$  and  $\|\mathbf{w}\| \leq \sqrt{2/\lambda}$ , we have  $0 \leq [1 - \mathbf{w}^\top xy]_+ \leq 1 + R \sqrt{2/\lambda}$

Hence,

$$\begin{aligned} \mathcal{R}(\hat{f}) &\leq \widehat{\mathcal{R}}(\hat{f}) + \text{Rad}(\mathcal{F}) + (b - a) \sqrt{\frac{\ln(1/\delta)}{n}} \\ &\leq \widehat{\mathcal{R}}(\hat{f}) + R \sqrt{\frac{2}{\lambda n}} + (1 + R \sqrt{2/\lambda}) \sqrt{\frac{\ln(1/\delta)}{n}} \\ &\leq \widehat{\mathcal{R}}(\hat{f}) + R \sqrt{\frac{8}{\lambda n}} + 3 \left(1 + R \sqrt{2/\lambda}\right) \sqrt{\frac{\ln(2/\delta)}{2n}}. \end{aligned}$$