

# CS 446 MJT — Homework 4

*junzhew3*  
*haoxuansun8*

2019/04/02

## 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) Assume all the data examples are on the unit sphere  $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ . The number of data

examples are  $\mathbf{n}$  ( $\mathbf{n}$  can be any integer which greater than 1). The  $\mathbf{m}$  data examples are

labeled as +1, the rest  $(\mathbf{n} - \mathbf{m})$  data examples are labeled as -1.

If  $\mathbf{m} = 0$ : we can set  $\mathbf{C} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq 0.5\}$ , so all the points with -1 label are

outside the  $\mathbf{C}$ .

if  $\mathbf{m} = 1$ : we can find a subspace  $\mathbf{C}$  which is tangent to the unit sphere  $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$  at the point which is labeled as +1.

If  $2 \leq \mathbf{m} \leq \mathbf{n}$ : Let us connect all the points with +1 label using subset of  $\mathbb{R}^{d-1}$  one by one to form a subset of  $\mathbb{R}^d$ . Since this subset is a Polyhedra, this subset is a convex subset.

So  $\mathbf{C}$  is equal to this subset. So all the points with +1 label are in the  $\mathbf{C}$ , all the points

with -1 label are outside the  $\mathbf{C}$ .

So  $\text{VC}(\mathcal{F}) \geq n$  for all  $n$ .

SO  $\text{VC}(\mathcal{F}) = \infty$ .

(b) So I will make all data examples  $((\mathbf{X}_i, \mathbf{Y}_i))_{i=1}^n$  be on a line. The distance between data examples is large enough to make sure that the RBF kernel SVM is close to the 1-nearest neighbor predictor. So in this situation, every data examples are support vectors.

1. Set  $m$  is equal to the number of support vectors, so  $m = n$ .
2. Set  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is equal to the support vectors, so  $(\mathbf{x}_1, \dots, \mathbf{x}_n) = (\mathbf{X}_i)_{i=1}^n$
3. Set  $(\alpha_i)_{i=1}^n$  is equal to  $(\mathbf{Y}_i)_{i=1}^n$ .

For each point  $\mathbf{X}_i$ , let  $\rho := \min_{j \in n} \|\mathbf{X}_i - \mathbf{x}_j\|_2$ ,  $T := \{j \in n : \|\mathbf{X}_i - \mathbf{x}_j\| = \rho\}$

$$\text{sgn} \left( \sum_{j=1}^n \alpha_j \exp \left( -\|\mathbf{X}_i - \mathbf{x}_j\|^2 / (2\sigma^2) \right) \right) = \text{sgn} \left( \sum_{j \in T} \alpha_j \exp \left( -\rho^2 / 2\sigma^2 \right) \right). \quad (1)$$

So for each point  $\mathbf{X}_i$ ,  $\rho := 0$  when  $\mathbf{X}_i = \mathbf{x}_i$ , so:

$$\text{sgn} \left( \sum_{j \in T} \hat{\alpha}_j \exp \left( -\rho^2 / 2\sigma^2 \right) \right) = \text{sgn}(\alpha_i) = \text{sgn}(\mathbf{Y}_i). \quad (2)$$

So  $\mathbf{Y}_i(\text{predicted}) = \mathbf{Y}_i(\text{True})$ . Each data example is correctly labeled.

So  $\text{VC}(\mathcal{G}_\sigma) \geq n$  for all  $n$ .

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

## 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) Using Cauchy-Schwarz inequality:

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{w}^\top \mathbf{x}_i = \mathbf{w}^\top \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i \mathbf{x}_i \right) \leq \|\mathbf{w}\| \|\mathbf{x}_\varepsilon\| \leq W \|\mathbf{x}_\varepsilon\| \quad (3)$$

$$\text{So : } \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) \leq W \|\mathbf{x}_\varepsilon\| \quad (4)$$

- (b) Since each value of  $\varepsilon_i$  has a 1/2 probability of being -1 and 1/2 probability of being +1.

$$\text{So } \mathbb{E}(\varepsilon_i) = 0, D(\varepsilon_i) = \mathbb{E}(\varepsilon_i^2) - \mathbb{E}(\varepsilon_i)^2 = 1,$$

$$\text{for any } \mathbf{a}, \mathbf{b} \in [1, n], \mathbf{a} \neq \mathbf{b}, \text{Cov}(\varepsilon_a, \varepsilon_b) = \mathbb{E}[\varepsilon_a - \mathbb{E}(\varepsilon_a)][\varepsilon_b - \mathbb{E}(\varepsilon_b)] = \mathbb{E}(\varepsilon_a \varepsilon_b)$$

$$\text{Since } \varepsilon_a \text{ and } \varepsilon_b \text{ are independent, so } \mathbb{E}(\varepsilon_a \varepsilon_b) = \mathbb{E}(\varepsilon_a) \mathbb{E}(\varepsilon_b) = 0, \text{ so } \text{Cov}(\varepsilon_a, \varepsilon_b) = 0$$

So assume  $\mathbf{x}_i \in \mathbb{R}^r$ :

$$\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 = \mathbb{E}_\varepsilon (x_{\varepsilon_1}^2 + \dots + x_{\varepsilon_r}^2) = \mathbb{E}_\varepsilon [(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n)^2 + \dots + (x_{1r}\varepsilon_1 + \dots + x_{nr}\varepsilon_n)^2]/n^2 \quad (5)$$

$$= (\mathbb{E}_\varepsilon [(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n)^2] + \dots + [\mathbb{E}_\varepsilon (x_{1r}\varepsilon_1 + \dots + x_{nr}\varepsilon_n)^2])/n^2 \quad (6)$$

$$\mathbb{E}_\varepsilon [(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n)^2] = \mathbb{E}_\varepsilon (x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n)^2 + D(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n) \quad (7)$$

$$= (x_{11}\mathbb{E}_\varepsilon(\varepsilon_1) + \dots + x_{n1}\mathbb{E}_\varepsilon(\varepsilon_n))^2 + D(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n) = D(x_{11}\varepsilon_1 + \dots + x_{n1}\varepsilon_n) \quad (8)$$

$$= D(x_{11}\varepsilon_1) + \dots + D(x_{n1}\varepsilon_n) + \text{Cov}(x_{11}\varepsilon_1, x_{21}\varepsilon_2) + \dots + \text{Cov}(x_{11}\varepsilon_1, x_{n1}\varepsilon_n) + \text{Cov}(x_{21}\varepsilon_2, x_{31}\varepsilon_3) \quad (9)$$

$$+ \dots + \text{Cov}(x_{21}\varepsilon_2, x_{n1}\varepsilon_n) + \dots + \text{Cov}(x_{(n-1)1}\varepsilon_{n-1}, x_{n1}\varepsilon_n) \quad (10)$$

$$= D(x_{11}\varepsilon_1) + \dots + D(x_{n1}\varepsilon_n) = x_{11}^2 D(\varepsilon_1) + \dots + x_{n1}^2 D(\varepsilon_n) = x_{11}^2 + \dots + x_{n1}^2 \quad (11)$$

So in the same way:

$$\mathbb{E}_\varepsilon[(x_{12}\epsilon_1 + \cdots + x_{n2}\epsilon_n)^2] = x_{12}^2 + \cdots + x_{n2}^2, \cdots, \mathbb{E}_\varepsilon[(x_{1r}\epsilon_1 + \cdots + x_{nr}\epsilon_n)^2] = x_{1r}^2 + \cdots + x_{nr}^2$$

$$\begin{aligned} \text{So: } \mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 &= [(x_{11}^2 + \cdots + x_{1r}^2) + (x_{21}^2 + \cdots + x_{2r}^2) + \cdots + (x_{n1}^2 + \cdots + x_{nr}^2)]/n^2 \\ &= (\|\mathbf{x}_1\|^2 + \|\mathbf{x}_2\|^2 + \cdots + \|\mathbf{x}_n\|^2)/n^2 \end{aligned}$$

Since  $\|\mathbf{x}_i\| \leq R$ , so  $\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 \leq nR^2/n^2 = R^2/n$ , so it is proved.

$$(c) \text{ Rad}(\mathcal{F}) = \mathbb{E}_\varepsilon \left( \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(\mathbf{x}_i) \right).$$

Using the conclusion in (a), we can get:  $\text{Rad}(\mathcal{F}) \leq W \mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|$

Since  $\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 \leq R^2/n$ , so  $(\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|)^2 \leq \mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\|^2 \leq R^2/n$ , so  $\mathbb{E}_\varepsilon \|\mathbf{x}_\varepsilon\| \leq R/\sqrt{n}$

So  $\text{Rad}(\mathcal{F}) \leq RW/\sqrt{n}$

### 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. Since  $\|\mathbf{w}\| \leq \sqrt{2/\lambda}$ , using the conclusion from the Problem 2, we can get:

$$\text{Rad}(\mathcal{F}) \leq R \sqrt{2/n\lambda} \quad (12)$$

2. There exists  $\rho \geq 0$  so that for any  $f, g \in \mathcal{F}$

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| \leq \rho |f(\mathbf{x}) - g(\mathbf{x})| \quad (13)$$

Here,  $\ell(f(\mathbf{x}), y) = [1 - \mathbf{w}_1^\top \mathbf{x} y]_+$ ,  $\ell(g(\mathbf{x}), y) = [1 - \mathbf{w}_2^\top \mathbf{x} y]_+$ .

(1) If  $(1 - \mathbf{w}_1^\top \mathbf{x} y) > 0, (1 - \mathbf{w}_2^\top \mathbf{x} y) > 0$ :

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = |\mathbf{w}_2^\top \mathbf{x} y - \mathbf{w}_1^\top \mathbf{x} y| = |\mathbf{w}_2^\top \mathbf{x} - \mathbf{w}_1^\top \mathbf{x}| = |f(\mathbf{x}) - g(\mathbf{x})| \quad (14)$$

(2) If  $(1 - \mathbf{w}_1^\top \mathbf{x} y) \leq 0, (1 - \mathbf{w}_2^\top \mathbf{x} y) \leq 0$ :

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = 0 \leq |f(\mathbf{x}) - g(\mathbf{x})| \quad (15)$$

(3) If  $(1 - \mathbf{w}_1^\top \mathbf{x} y) \leq 0, (1 - \mathbf{w}_2^\top \mathbf{x} y) > 0, \mathbf{w}_1^\top \mathbf{x} y \geq 1, \mathbf{w}_2^\top \mathbf{x} y < 1$ :

If  $y = 1$ :  $\mathbf{w}_1^\top \mathbf{x} \geq 1, \mathbf{w}_2^\top \mathbf{x} < 1$

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = |1 - \mathbf{w}_2^\top \mathbf{x} y| = 1 - \mathbf{w}_2^\top \mathbf{x} \leq \mathbf{w}_1^\top \mathbf{x} - \mathbf{w}_2^\top \mathbf{x} = |f(\mathbf{x}) - g(\mathbf{x})| \quad (16)$$

If  $y = -1$ :  $\mathbf{w}_1^\top \mathbf{x} \leq -1, \mathbf{w}_2^\top \mathbf{x} > -1$

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = |1 - \mathbf{w}_2^\top \mathbf{x} y| = 1 + \mathbf{w}_2^\top \mathbf{x} \leq \mathbf{w}_2^\top \mathbf{x} - \mathbf{w}_1^\top \mathbf{x} = |f(\mathbf{x}) - g(\mathbf{x})| \quad (17)$$

(4) If  $(1 - \mathbf{w}_1^\top \mathbf{x} y) > 0, (1 - \mathbf{w}_2^\top \mathbf{x} y) \leq 0, \mathbf{w}_1^\top \mathbf{x} y < 1, \mathbf{w}_2^\top \mathbf{x} y \geq 1$ :

If  $y = 1$ :  $\mathbf{w}_1^\top \mathbf{x} < 1, \mathbf{w}_2^\top \mathbf{x} \geq 1$

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = |1 - \mathbf{w}_1^\top \mathbf{x} y| = 1 - \mathbf{w}_1^\top \mathbf{x} \leq \mathbf{w}_2^\top \mathbf{x} - \mathbf{w}_1^\top \mathbf{x} = |f(\mathbf{x}) - g(\mathbf{x})| \quad (18)$$

If  $y = -1$ :  $\mathbf{w}_1^\top \mathbf{x} > -1, \mathbf{w}_2^\top \mathbf{x} \leq -1$

$$|\ell(f(\mathbf{x}), y) - \ell(g(\mathbf{x}), y)| = |1 - \mathbf{w}_1^\top \mathbf{x} y| = 1 + \mathbf{w}_1^\top \mathbf{x} \leq \mathbf{w}_1^\top \mathbf{x} - \mathbf{w}_2^\top \mathbf{x} = |f(\mathbf{x}) - g(\mathbf{x})| \quad (19)$$

So  $\rho = 1$  here.

3. There exists  $[a, b]$  so that  $\ell(f(\mathbf{x}), y) \in [a, b]$  for any  $f \in \mathcal{F}$ .

Here  $\ell(f(\mathbf{x}), y) = [1 - \mathbf{w}^\top \mathbf{x}y]_+ \geq 0$ , so  $a=0$  here.

In order to get the maximun of the  $\ell(f(\mathbf{x}), y)$ , ( $1 - \mathbf{w}^\top \mathbf{x}y > 0$ ), so now:

$$\ell(f(\mathbf{x}), y) = [1 - \mathbf{w}^\top \mathbf{x}y]_+ = 1 - \mathbf{w}^\top \mathbf{x}y = 1 - \|\mathbf{x}\| \|\mathbf{w}\| \cos(\theta) y \leq 1 + \|\mathbf{x}\| \|\mathbf{w}\| = 1 + R\sqrt{2/\lambda} \quad (20)$$

So  $b = \left(1 + R\sqrt{2/\lambda}\right)$  when  $y \cos(\theta) = -1$ ,  $\|\mathbf{x}\| = R$ ,  $\|\mathbf{w}\| = \sqrt{2/\lambda}$ .

So  $b - a = \left(1 + R\sqrt{2/\lambda}\right)$

4. With probability  $\geq 1 - \delta$ , every  $f \in \mathcal{F}$  satisfies:

$$\mathcal{R}_\ell(f) \leq \widehat{\mathcal{R}}_\ell(f) + 2\rho \text{Rad}(\mathcal{F}) + 3(b - a) \sqrt{\frac{\ln(2/\delta)}{2n}} \quad (21)$$

Here  $\rho = 1$ ,  $b - a = \left(1 + R\sqrt{2/\lambda}\right)$  and  $\text{Rad}(\mathcal{F}) \leq R\sqrt{2/n\lambda}$ .

$$\text{So: } \mathcal{R}(\hat{f}) \leq \widehat{\mathcal{R}}(\hat{f}) + R\sqrt{\frac{8}{\lambda n}} + 3\left(1 + R\sqrt{2/\lambda}\right) \sqrt{\frac{\ln(2/\delta)}{2n}}. \quad (22)$$