

第 3 次作业

赵丰 2017310711

November 18, 2017

- 
- **Acknowledgments:** This coursework refers to textbook.
  - **Collaborators:**
    - 2.2 was solved with the help from  
[https://en.wikipedia.org/wiki/Sauer%E2%80%933Shelah\\_lemma](https://en.wikipedia.org/wiki/Sauer%E2%80%933Shelah_lemma).
    - 2.3 was solved with the help from  
<http://www.cs.nyu.edu/~mohri/ml07/sol2.pdf>
    - Discussion with the teaching assistant helped me finishing 3(b).
- 

I use `enumerate` to generate answers for each question:

1. 1.1. (a)

$$\hat{\mathcal{R}}_S(cG + d) = \mathbb{E}_{\sigma} \left[ \sup_{g \in G} \frac{\sigma \cdot (cg_S + d)}{m} \right] \quad (1)$$

因为  $\forall g \in G$ ,

$$\mathbb{E}_{\sigma} \left[ \frac{d}{m} \sum_{i=1}^m \sigma_i \right] = 0 \quad (2)$$

所以

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in G} \frac{\sigma \cdot (cg_S + d)}{m} \right] = |c| \mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_{i=1}^m \text{sign}(c) \sigma_i g(z_i) \right] \quad (3)$$

因为  $\text{sign}(c)\sigma_i$  与  $\sigma_i$  具有相同的分布, 所以

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in G} \frac{\sigma \cdot (cg_S + d)}{m} \right] = |c| \hat{\mathcal{R}}_S(G) \quad (4)$$

从而推出:

$$\hat{\mathcal{R}}_S(cG + d) = |c| \hat{\mathcal{R}}_S(G) \quad (5)$$

(b)

$$\hat{\mathcal{R}}_S(\text{conv}(G)) = \mathbb{E}_{\sigma} \left[ \sup_{g_i \in G} \frac{\sigma \cdot (\sum_{i=1}^n \alpha_i g_i)}{m} \right] \quad (6)$$

其中  $\mathbf{g}_i = (g_i(z_1), \dots, g_i(z_m))^T$   
一方面:

$$\begin{aligned} \mathbb{E}_{\sigma} \left[ \sup_{g_i \in G} \frac{\sigma \cdot (\sum_{i=1}^n \alpha_i g_i)}{m} \right] &\leq \sum_{i=1}^n \alpha_i \mathbb{E}_{\sigma} \left[ \sup_{g_i \in G} \frac{\sigma \cdot g_i}{m} \right] \\ &= \sum_{i=1}^n \alpha_i \hat{\mathcal{R}}_S(G) \\ &= \hat{\mathcal{R}}_S(G) \end{aligned}$$

另一方面, 取  $n = 1, \alpha_1 = 1$  有

$$\mathbb{E}_{\sigma} \left[ \sup_{g_i \in G} \frac{\sigma \cdot (\sum_{i=1}^n \alpha_i g_i)}{m} \right] \geq \mathbb{E}_{\sigma} \left[ \sup_{g_1 \in G} \frac{\sigma \cdot (\alpha_1 g_1)}{m} \right] = \hat{\mathcal{R}}_S(G) \quad (7)$$

所以推出

$$\hat{\mathcal{R}}_S(\text{conv}(G)) = \hat{\mathcal{R}}_S(G) \quad (8)$$

(c) 可以证明  $\sup_{a \in A, b \in B} (a + b) = \sup_{a \in A} a + \sup_{b \in B} b$ , 所以

$$\begin{aligned} \hat{\mathcal{R}}_S(G_1 + G_2) &= \mathbb{E}_{\sigma} \left[ \sup_{\substack{g_1 \in G_1 \\ g_2 \in G_2}} \frac{\sigma \cdot (g_1 + g_2)}{m} \right] \\ &= \mathbb{E}_{\sigma} \left[ \sup_{g_1 \in G_1} \frac{\sigma \cdot g_1}{m} \right] + \mathbb{E}_{\sigma} \left[ \sup_{g_2 \in G_2} \frac{\sigma \cdot g_2}{m} \right] \\ &= \hat{\mathcal{R}}_S(G_1) + \hat{\mathcal{R}}_S(G_2) \end{aligned}$$

1.2.

$$\begin{aligned} \hat{\mathcal{R}}_S(G) &= \mathbb{E}_{\sigma} \left[ \sup_{h \in H} \frac{\sum_{i=1}^m \sigma_i 1_{h(x_i) \neq y_i}}{m} \right] \\ &= \mathbb{E}_{\sigma} \left[ \sup_{h \in H} \frac{\sum_{i=1}^m \sigma_i \frac{1 - h(x_i)y_i}{2}}{m} \right] \end{aligned} \quad (9)$$

因为  $|y_i| = 1$ , 由 **P1.1(a)** 的结论

$$\begin{aligned} \mathbb{E}_{\sigma} \left[ \sup_{h \in H} \frac{\sum_{i=1}^m \sigma_i \frac{1 - h(x_i)y_i}{2}}{m} \right] &= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in H} \frac{\sum_{i=1}^m \sigma_i h(x_i)}{m} \right] \\ &= \frac{1}{2} \hat{\mathcal{R}}_{S_x}(H) \end{aligned} \quad (10)$$

因此推出

$$\hat{\mathcal{R}}_S(G) = \frac{1}{2} \hat{\mathcal{R}}_{S_x}(H) \quad (11)$$

2. (a) 对于直线上  $m$  个不同的点, 考虑含  $k (1 \leq k \leq m)$  个正类有  $m + 1 - k$  种分法, 外加全为负类的一种分法, 有

$$\Pi_H(m) = 1 + \sum_{k=1}^m (m + 1 - k) = 1 + \frac{m(m+1)}{2} \quad (12)$$

根据 Sauer 引理

$$\Pi_H(m) \leq \binom{m}{0} + \binom{m}{1} + \binom{m}{2} = 1 + \frac{m(m+1)}{2} \quad (13)$$

因此对于数轴上区间作为假设集的例子, 引理给出的上界是紧的。

- (b) 从  $\{1, 2, \dots, m\}$  中选不超过  $d$  个数共有  $K = \sum_{k=0}^d \binom{m}{k}$  中取法，对应于第  $i$  种取法， $h_i(x)$  是这样一个向量，将取出的数对应的分量位置置 1，其余位置置 0。注意  $h_i$  虽然是从  $\mathcal{X}$  到  $\{0, 1\}^m$  空间的一个映射但却不依赖于  $x \in \mathcal{X}$ 。将这  $K$  个映射构成我们考虑的函数空间  $H$ ，由定义可知  $H$  的增长函数为  $K$ 。当  $m = d$  时  $K = 2^d$ ，由 VC 维的定义可知  $H$  的 VC 维为  $d$ 。所以我们构造出了这样一个假设集合使得 Sauer 引理取到了等号。
- (c) 只需对  $\mathbb{R}^n$  中证明对任意  $n+3$  个点只用一个闭球均不能完全分离两类。反设  $\exists \mathbf{x}_1, \dots, \mathbf{x}_{n+3}$ ，有  $\forall y_i \in \{0, 1\}, 1 \leq i \leq n+3$ ，均  $\exists r, \mathbf{x}^*$ ，使得  $y_i(\|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - r^2) > 0, 1 \leq i \leq n+3$
- 下面考虑  $\mathbb{R}^{n+1}$  中的点集  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{n+3}$ ，其中  $\tilde{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i \\ \|\mathbf{x}_i\|_2^2 \end{pmatrix}$ ，对于这  $n+3$  个  $\mathbb{R}^{n+1}$  中的点，由前面假设：
- $\forall y_i \in \{0, 1\}, 1 \leq i \leq n+3$ ，构造  $\mathbf{a} = \begin{pmatrix} -2\mathbf{x}^* \\ 1 \end{pmatrix}, b = \|\mathbf{x}^*\|_2^2 - r^2$ ，则有  $y_i(\mathbf{a} \cdot \tilde{\mathbf{x}}_i + b) > 0, 1 \leq i \leq n+3$ ，即用超平面在  $n+1$  维空间中对任意  $n+3$  个点都可以完全分离。这与超平面的 VC 维等于考虑的空间维数加 2 矛盾。
- 因此用  $\mathbb{R}^n$  中的闭球分类 VC 维至多为  $n+2$ 。

3. (a) 对于给定的样本，由  $\hat{h}$  定义我们有

$$\mathbb{E}[\widehat{\text{err}}(\hat{h})] \leq \mathbb{E}[\widehat{\text{err}}(h^*)] = \text{err}(h^*) \quad (14)$$

其中最后一式利用了对于给定的假设，经验误差的均值等于泛化误差。另一方面，

$$\mathbb{E}[\text{err}(\hat{h})] \geq \mathbb{E}[\text{err}(h^*)] = \text{err}(h^*) \quad (15)$$

- (b) 设  $\tilde{S}$  与  $S$  只有一个样本点不同，设其指标为  $(x'_i, y'_i)$ ，记函数

$$\Phi(S) = \widehat{\text{err}}(\hat{h}) - \mathbb{E}[\widehat{\text{err}}(\hat{h})] \quad (16)$$

$$\Phi(S) - \Phi(S') = \min_{h \in \mathcal{H}, S} \widehat{\text{err}}(h) - \min_{h \in \mathcal{H}, S'} \widehat{\text{err}}(h) \quad (17)$$

$$\geq \min_{h \in \mathcal{H}} \frac{1}{m} (1_{h(x_i) \neq y_i} - 1_{h(x'_i) \neq y'_i}) \quad (18)$$

$$\geq \frac{-1}{m} \quad (19)$$

同理

$$\Phi(S') - \Phi(S) \geq \frac{-1}{m} \quad (20)$$

从而有

$$|\Phi(S) - \Phi(S')| \leq \frac{1}{m} \quad (21)$$

由 McDiarmid 不等式，有  $1 - \delta$  的概率使得：

$$|\Phi(S) - \Phi(S')| \leq \sqrt{\frac{\ln(2/\delta)}{2m}} \quad (22)$$

由于  $\delta$  是小量,  $\frac{\ln(2)}{m}$  相比  $\frac{\ln(1/\delta)}{m}$  为小量, 因此

$$|\Phi(S) - \Phi(S')| = O\left(\frac{\sqrt{\ln(1/\delta)}}{m}\right) \quad (23)$$

常数为  $\frac{\sqrt{2}}{2}$

- (c) (a) 中的结果说明, 极小化泛化误差得到的最优估计器的统计平均性能要优于极小化经验误差的估计器 (第二个不等式), 但针对两种情况下的最优估计器, 泛化误差要大于经验误差的平均 (第一个不等式)。(b) 中的结果说明, 当样本足够多时, 经验误差以大概率集中在经验误差的平均值附近。