

机器学习理论研究导引

作业三

陈晟 MG21330006

2022 年 4 月 27 日

作业提交注意事项

- (1) 本次作业提交截止时间为 **2022/05/03 23:59:59**, 截止时间后不再接收作业, 本次作业记零分;
- (2) 作业提交方式: 使用此 LaTeX 模板书写解答, 只需提交编译生成的 pdf 文件, 将 pdf 文件上传到以下 ftp 服务器的指定位置:
地址: sftp://210.28.132.67:22, 用户名: mlt2022, 密码: mltspring2022@nju
文件夹位置: /C:/Users/mlt2022/hw_submissions/hw3_submission/ ;
- (3) pdf 文件命名方式: 学号-姓名-作业号-v 版本号, 例 MG1900000-张三-3-v1; 如果需要更改已提交的解答, 请在截止时间之前提交新版本的解答, 并将版本号加一;
- (4) 未按照要求提交作业, 或 pdf 命名方式不正确, 将会被扣除部分作业分数.

1 [50pts] Generalization Bound Based-on VC Dimension

在书中，为了推导基于 VC 维的泛化界，先推导出了基于增长函数的泛化界，再利用 VC 维和增长函数之间的关系（定理 3.1）完成证明。在本题中，我们将从基于 Rademacher 复杂度的泛化界出发，推导基于 VC 维的泛化界。若假设空间 \mathcal{H} 的 VC 维为 d ，

(1) [30pts] 试证明对特征空间上的任一分布 \mathcal{D} 和任一正整数 m 都有

$$\mathfrak{R}_m(\mathcal{H}) \leq \sqrt{\frac{2d \ln\left(\frac{em}{d}\right)}{m}}.$$

(2) [20pts] 试利用基于 Rademacher 复杂度的泛化界，结合 (1) 中的结果，推导基于 VC 维的泛化界。

Proof.

(1)

VC 维: $VCdim(\mathcal{H}) = \max\{m : \Pi_{\mathcal{H}}(m) = 2^m\}$

Let \mathcal{H} be a hypothesis set with $VCdim(\mathcal{H})$. Then for all $m \geq d$,

$$\begin{aligned} \Pi_{\mathcal{H}}(m) &\leq \sum_{i=0}^d \binom{m}{i} \leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \\ &\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} = \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i \\ &= \left(\frac{m}{d}\right)^d \left(1 + \frac{d}{m}\right)^m \leq \left(\frac{m}{d}\right)^d e^d \\ \therefore \Pi_{\mathcal{H}}(m) &\leq \left(\frac{em}{d}\right)^d = O(m^d) \end{aligned}$$

对于 $D = \{x_1, x_2, \dots, x_m\}$, $\mathcal{H}|_D$ 为假设空间 \mathcal{H} 在 D 上的限制。由于 $h \in \mathcal{H}$ 的值域为 $\{-1, +1\}$, 可知 $\mathcal{H}|_D$ 中的元素为模长 \sqrt{m} 的向量。因此由 ppt 中定理可得:

$$\begin{aligned} \therefore \mathcal{R}_m(\mathcal{H}) &= E_D[E_{\sigma}[\sup_{u \in \mathcal{H}|_D} \frac{1}{m} \sum_{i=1}^m \sigma_i u_i]] \\ &\leq E_D\left[\frac{\sqrt{m} \sqrt{2 \ln |\mathcal{H}|_D}}{m}\right] \\ \text{由 } |\mathcal{H}|_D| &\leq \Pi_{\mathcal{H}}(m) \text{ 有} \\ \mathcal{R}_m(\mathcal{H}) &\leq E_D\left[\frac{\sqrt{m} \sqrt{2 \ln \Pi_{\mathcal{H}}(m)}}{m}\right] = \sqrt{\frac{2 \ln \Pi_{\mathcal{H}}(m)}{m}} \\ \text{带入 } \Pi_{\mathcal{H}}(m) &\leq \left(\frac{em}{d}\right)^d \text{ 即证: } \mathfrak{R}_m(\mathcal{H}) \leq \sqrt{\frac{2d \ln\left(\frac{em}{d}\right)}{m}}. \end{aligned}$$

(2)

设 \mathcal{H} 为一族从 $\{-1, +1\}$ 中取值且 VC 维等于 d 的函数。那么，对于任意的 $\delta > 0$ ，下面的不等式至少有 $1 - \delta$ 的概率对假设 $h \in \mathcal{H}$ 成立

$$\mathcal{R}_m(h) \leq \hat{\mathcal{R}}_m(h) + \sqrt{\frac{2d \log \frac{em}{d}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

因此，这种情况下的泛化界满足 $\mathcal{R}_m(h) \leq \hat{\mathcal{R}}_m(h) + \mathcal{O}\left(\sqrt{\frac{\log\left(\frac{em}{d}\right)}{\frac{m}{d}}}\right)$

它说明了 $\frac{m}{d}$ 对于泛化的重要性。这个推论同样符合奥卡姆剃刀原则

(由 Sauer 引理，我们还能得到以下的 VC，不通过 Rademacher 复杂度得到: $\mathcal{R}_m(h) \leq \hat{\mathcal{R}}_m(h) + \sqrt{\frac{8d \log \frac{2em}{d} + 8 \log \frac{4}{\delta}}{m}}$)

2 [50pts] Generalization Bound and Covering Numbers

设 \mathcal{H} 为一假设空间, 假设的定义域为 \mathcal{X} , 值域为 $\mathcal{Y} \subset \mathbb{R}$. $\forall \epsilon > 0$, 可如下定义覆盖数 (Covering Number):

$$\mathcal{N}(\mathcal{H}, \epsilon) = \min \{k \in \mathbb{N} | \exists \{h_1, \dots, h_k\} \subset \mathcal{H}, \text{ s.t. } \forall h \in \mathcal{H}, \exists i \in [k], \|h - h_i\|_\infty \leq \epsilon\}, \quad (2.1)$$

其中 $\|h - h_i\|_\infty = \max_{x \in \mathcal{X}} |h(x) - h_i(x)|$. 覆盖数可以衡量一个假设空间的复杂度: 覆盖数越大, 意味着这一假设空间越复杂. 本题利用覆盖数证明了平方损失下的一个泛化界, 该结论也可以说明覆盖数可以衡量假设空间的复杂度. 令 \mathcal{D} 为 $\mathcal{X} \times \mathcal{Y}$ 上的一个分布, 且有标记样本根据这一分布采样得到. 定义 $h \in \mathcal{H}$ 的泛化误差为

$$R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h(x) - y)^2]. \quad (2.2)$$

训练集 $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ 上的经验误差为

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2. \quad (2.3)$$

设 \mathcal{H} 是有界的, 即 $\exists M > 0$, 使得 $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h \in \mathcal{H}, |h(x) - y| \leq M$.

(1) [15pts] 令 $L_S = R(h) - \hat{R}_S(h)$, 试证明 $\forall h_1, h_2 \in \mathcal{H}, S$,

$$|L_S(h_1) - L_S(h_2)| \leq 4M \|h_1 - h_2\|_\infty. \quad (2.4)$$

(2) [15pts] 设 \mathcal{H} 可以被 k 个子集 $\mathcal{B}_1, \dots, \mathcal{B}_k$ 覆盖, 即 $\mathcal{H} = \mathcal{B}_1 \cup \dots \cup \mathcal{B}_k$. 试证明 $\forall \epsilon > 0$,

$$\Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon \right] \leq \sum_{i=1}^k \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |L_S(h)| \geq \epsilon \right]. \quad (2.5)$$

(3) [20pts] 令 $k = \mathcal{N}(\mathcal{H}, \frac{\epsilon}{8M})$, $\mathcal{B}_1, \dots, \mathcal{B}_k$ 为 \mathcal{H} 的覆盖, 其中 $\forall i \in [k]$, \mathcal{B}_i 的圆心为 h_i , 半径为 $\frac{\epsilon}{8M}$. 试证明

$$\Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |L_S(h)| \geq \epsilon \right] \leq \Pr_{S \sim \mathcal{D}^m} \left[|L_S(h_i)| \geq \frac{\epsilon}{2} \right]. \quad (2.6)$$

并利用 Hoeffding 不等式证明

$$\Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \leq \mathcal{N}\left(\mathcal{H}, \frac{\epsilon}{8M}\right) 2e^{-\frac{m\epsilon^2}{2M^4}}. \quad (2.7)$$

Proof.

(1)

由于 $R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h(x) - y)^2]$. 训练集 $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ 上的经验误差为

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

$$\therefore L_S = R(h) - \hat{R}_S(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h(x) - y)^2] - \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

$$\forall h_1, h_2 \in \mathcal{H}, S, |L_S(h_1) - L_S(h_2)| =$$

$$|\mathbb{E}_{(x,y) \sim \mathcal{D}} [(h_1(x) - y)^2] - \frac{1}{m} \sum_{i=1}^m (h_1(x_i) - y_i)^2 - (\mathbb{E}_{(x,y) \sim \mathcal{D}} [(h_2(x) - y)^2] - \frac{1}{m} \sum_{i=1}^m (h_2(x_i) - y_i)^2)|$$

$\therefore \forall h_1, h_2 \in \mathcal{H}$

$$\therefore \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h_1(x) - y)^2] = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h_2(x) - y)^2]$$

$$\frac{1}{m} \sum_{i=1}^m (h_2(x_i) - y_i)^2 - \frac{1}{m} \sum_{i=1}^m (h_1(x_i) - y_i)^2 =$$

$$\frac{1}{m} \sum_{i=1}^m ((h_2(x_i) - y_i)^2 - (h_1(x_i) - y_i)^2) =$$

$$\frac{1}{m} \sum_{i=1}^m ((h_2(x_i) + h_1(x_i) - 2y_i)(h_2(x_i) - h_1(x_i)))$$

$\therefore \mathcal{H}$ 是有界的, 即 $\exists M > 0$, 使得 $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h \in \mathcal{H}, |h(x) - y| \leq My$

$$\therefore |(h_1(x_i) + h_2(x_i) - 2y_i)^2| \leq |(2 * \max(h_1(x_i), h_2(x_i)) - 2y_i)^2| \leq 4M$$

并且显然, $|h_1(x_i) - h_2(x_i)| \leq \|h_1 - h_2\|_\infty = \max_{x \in \mathcal{X}} |h_1(x) - h_2(x)|$

$$\therefore |L_S(h_1) - L_S(h_2)| = |\frac{1}{m} \sum_{i=1}^m (h_2(x_i) - y_i)^2 - \frac{1}{m} \sum_{i=1}^m (h_1(x_i) - y_i)^2| =$$

$$|\frac{1}{m} \sum_{i=1}^m ((h_2(x_i) + h_1(x_i) - 2y_i)(h_2(x_i) - h_1(x_i)))| =$$

$$\frac{1}{m} \sum_{i=1}^m |(h_2(x_i) + h_1(x_i) - 2y_i)(h_2(x_i) - h_1(x_i))| \leq \frac{1}{m} \sum_{i=1}^m (4M \|h_1 - h_2\|_\infty) = 4M \|h_1 - h_2\|_\infty$$

即证明 $\forall h_1, h_2 \in \mathcal{H}, S, |L_S(h_1) - L_S(h_2)| \leq 4M \|h_1 - h_2\|_\infty$

(2)

当 $k=1$ 时, 显然有 $\mathcal{H} = B_1$

所以 $\Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] = \sum_{i=1}^1 \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon]$ 题目要求显然成立

假设 $k=n$ 时, 对任意 $\mathcal{H} \cap B_1, B_2 \dots B_n$, 都有 $\Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] = \sum_{i=1}^n \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon]$

成立

则 $k=n+1$ 时, 会新增加一个 B_{n+1} , 由上可知该 B_{n+1} 可以是 $k=n$ 的某一种子集划分情况

下, 某一个子集 B_p 的子集, 因此可以有

$$\sum_{i=1}^n + 1 \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon] = \sum_{i=1}^n \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon] + \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_{n+1}} |L_S(h)| \geq \epsilon]$$

$$\therefore \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] \leq \sum_{i=1}^{n+1} \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon] \text{ 成立}$$

由数学归纳法可知, 对任意 $k(k>0)$, 有

$$\mathcal{H} = B_1 \cup \dots \cup B_k. \forall \epsilon > 0, \text{ 有 } \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] \leq \sum_{i=1}^k \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon].$$

(3)

$$\therefore k = \mathcal{N}(\mathcal{H}, \frac{\epsilon}{8M})$$

$$\therefore k = \mathcal{N}(\mathcal{H}, \frac{\epsilon}{8M}) = \min \{k \in \mathbb{N} | \exists \{h_1, \dots, h_k\} \subset \mathcal{H}, \text{ s.t. } \forall h \in \mathcal{H}, \exists i \in [k], \|h - h_i\|_\infty \leq \frac{\epsilon}{8M}\}$$

由于 $\forall i \in [k], B_i$ 的圆心为 h_i , 半径为 $\frac{\epsilon}{8M}$

$$\therefore \Pr_{S \sim \mathcal{D}^m} [|L_S(h_i)| \geq \frac{\epsilon}{2}] = R(h_i) - \hat{R}_S(h_i) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(h_i(x) - y)^2] - \frac{1}{m} \sum_{j=1}^m (h_i(x_j) - y_j)^2$$

根据 ppt 对引理 4.1 证明过程中的 $\Pr(\sup_{h \in \mathcal{H}} |\hat{E}_D(h) - \hat{E}_{D'}(h)| \geq \frac{1}{2}\epsilon) \geq \frac{1}{2} \Pr(\sup_{h \in \mathcal{H}} |E(h) - \hat{E}_D(h)| > \epsilon)$

其中 D 和 D' 都是从 D 独立同分采样的训练集

$$\Pr_{S \sim \mathcal{D}^m} [|L_S(h_i)| \geq \frac{\epsilon}{2}] \geq \Pr_{B_i \in Q} [E_{S_i \sim D^m} [\mathcal{I}(\sup_{h \in \mathcal{H}} |\hat{E}_{B_i}(h) - \hat{E}_{D'}(h)| \geq \frac{1}{2}\epsilon)]]$$

$$\geq \Pr_{h_i \sim B_i} [\mathcal{I}(|\hat{E}_D(h_0) - E(h_0)| - |\hat{E}_{h_i}(h_0) - E(h_0)| \geq \frac{1}{2}\epsilon)] \geq 1 - \Pr(|\hat{E}_{h_i}(h_0) - E(h_0)| > \frac{1}{2}\epsilon)$$

由 Chebyshev 不等式 (1.21) 可得

$$\Pr(|\hat{E}_{h_i}(h_0) - E(h_0)| > \frac{1}{2}\epsilon) \leq \frac{4(1-E(h_0))E(h_0)}{\epsilon^2 m} \leq \frac{1}{\epsilon^2 m}$$

由于 $\forall i \in [k], B_i$ 的圆心为 h_i , 半径为 $\frac{\epsilon}{8M}$, 于是可得 $\Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon] \leq$

$$\Pr_{S \sim \mathcal{D}^m} [|L_S(h_i)| \geq \frac{\epsilon}{2}]$$

b. 利用 Hoeffding 不等式证明 $\Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \epsilon] \leq \mathcal{N}(\mathcal{H}, \frac{\epsilon}{8M}) 2e^{-\frac{m\epsilon^2}{2M^4}}$

由 (2) $\Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] \leq \sum_{i=1}^k \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in B_i} |L_S(h)| \geq \epsilon]$

$$\begin{aligned}
& \text{且 } \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] = \Pr_{S \sim \mathcal{D}^m} [\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon] \\
& \therefore \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \leq \sum_{i=1}^k \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \\
& \leq \mathcal{N} \left(\mathcal{H}, \frac{\epsilon}{8M} \right) \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \\
& \text{显然, 由 Hoeffding 不等式, 有: } \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \leq 2e^{-\frac{m\epsilon^2}{2M^4}} \\
& \therefore \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \leq \mathcal{N} \left(\mathcal{H}, \frac{\epsilon}{8M} \right) \Pr_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{B}_i} |R(h) - \hat{R}_S(h)| \geq \epsilon \right] \leq \\
& \mathcal{N} \left(\mathcal{H}, \frac{\epsilon}{8M} \right) 2e^{-\frac{m\epsilon^2}{2M^4}}
\end{aligned}$$