

# 主动学习理论概要

*Steve Hanneke's "Theory of Active Learning"*

*Chapter 2, 4, 5*

孟庆鑫

## 1 基本定义和符号

样本空间  $\mathcal{X}$  配  $\sigma$ -代数  $\mathcal{B}_{\mathcal{X}}$  使  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  为 Borel 空间,  $\mathcal{Y} = \{-1, 1\}$  为标签空间, 对应  $\mathcal{X} \times \mathcal{Y}$  上配  $\sigma$ -代数  $\mathcal{B} = \mathcal{B}_{\mathcal{X}} \otimes 2^{\mathcal{Y}}$ , 并配概率测度  $\mathcal{P}_{XY}$ , 记  $\mathcal{P}$  为  $\mathcal{P}_{XY}$  对  $\mathcal{X}$  的边缘分布, 以下均考虑可测情况.

$$\eta(x) \triangleq \mathbb{P}(Y = +1|X = x), \forall x \in \mathcal{X}, (X, Y) \sim \mathcal{P}_{XY},$$

$$h \text{ 的错误率 } \text{er}(h) \triangleq \mathcal{P}_{XY}((x, y) : h(x) \neq y), \forall h : \mathcal{X} \rightarrow \mathcal{Y},$$

$\mathcal{Z} \triangleq \{(X_i, Y_i)\}_{i=1}^{\infty} \stackrel{iid}{\sim} \mathcal{P}_{XY}$ , 前  $m$  截断  $\mathcal{Z}_m \triangleq \{(X_i, Y_i)\}_{i=1}^m$ , 未标记数据集  $\mathcal{Z}_{\mathcal{X}} \triangleq \{X_i\}_{i=1}^{\infty}$ , 并假设  $\mathcal{Z}_{\mathcal{X}}$  可被完整访问, 在主动学习算法中, 给定  $n$ , 请求观测值  $Y_{i_1}$ , 请求观测值  $Y_{i_2}$ , 依此 ..., 不超过  $n$  次请求, 算法停止并返回  $\hat{h}$ . 我们感兴趣的是算法  $\mathcal{A} : \mathcal{Z}_n \rightarrow \hat{h}$  作为  $n$  的函数, 其行为特征.

**定义 1**  $\forall \mathcal{A}$  为主动学习算法, 如果  $\forall \varepsilon > 0, \forall \delta \in [0, 1], \forall \mathcal{P}_{XY}, \forall n \geq \Lambda(\varepsilon, \delta, \mathcal{P}_{XY})$ ,  $n \in \mathbb{N}$ ,  $(\mathcal{A}, n)$  产生了  $\hat{h}$ ,  $\text{er}(\hat{h}) \leq \varepsilon$  的概率至少为  $1 - \delta$ , 则算法  $\mathcal{A}$  达到了标签复杂性  $\Lambda$ .

噪声率  $\nu \triangleq \inf_{h \in \mathbb{C}} \text{er}(h)$ ,  $\mathbb{C}$  是确定的分类器族, 我们特别关注  $\Lambda(\nu + \varepsilon, \delta, \mathcal{P}_{XY})$ . 记  $f^* = \arg \inf_{h \in \mathbb{C}} \text{er}(h)$ , 有  $f^* \in \bar{\mathbb{C}}$  ( $\mathbb{C}$  的闭包),  $\nu = \text{er}(f^*)$ ,

$$\forall \text{ 分类器族 } \mathcal{H}, \forall \varepsilon \in [0, 1], \varepsilon\text{-极小集 } \mathcal{H}(\varepsilon) \triangleq \{h \in \mathcal{H} : \text{er}(h) - \inf_{g \in \mathcal{H}} \text{er}(g) \leq \varepsilon\},$$

$$\forall \text{ 中心 } h \text{ 的 } \varepsilon\text{-球 } B_{\mathcal{H}, \mathcal{P}}(h, \varepsilon) \triangleq \{g \in \mathcal{H} : \mathcal{P}(x : g(x) \neq h(x)) \leq \varepsilon\},$$

$$\mathcal{H} \text{ 的半径 } R(\mathcal{H}) \triangleq \sup_{h \in \mathcal{H}} \mathcal{P}(x : h(x) \neq f^*(x)) \stackrel{\text{显然}}{=} \inf \{\varepsilon : \mathcal{H} = B_{\mathcal{H}, \mathcal{P}}(f^*, \varepsilon)\},$$

分歧域  $\mathcal{D}(\mathcal{H}) \triangleq \{x \in \mathcal{X} : \exists h, g \in \mathcal{H} \text{ s.t. } h(x) \neq g(x)\}$ , 一种构造算法的策略是: 从  $\mathcal{D}(\mathcal{H})$  里抽样本, 请求观测值, 再反过来更新  $\mathcal{H}$ , 分歧率  $\Delta(\mathcal{H}) \triangleq \mathcal{P}(\mathcal{D}(\mathcal{H}))$ .

定义 2  $\forall r_0 > 0, \forall h$  对  $\mathbb{C}$  在  $\mathcal{P}$  下的分歧率为

$$\theta_h(r_0) = \sup_{r > r_0} \frac{\Delta(B(h, r))}{r}.$$

特殊的, 当  $h = f^*$  时记为  $\theta(r_0)$ , 称为分类器族  $\mathbb{C}$  在  $\mathcal{R}_Y$  下的分歧率.

情况 1  $\exists a \in [1, \infty), \alpha \in [0, 1], \forall h \in \mathbb{C}$ ,

$$\mathcal{P}(x : h(x) \neq f^*(x)) \leq a(\text{er}(h) - \text{er}(f^*))^\alpha.$$

例 1  $\mathcal{X} = [0, 1], \mathbb{C} = \{1_{[z, 1]}^\pm : z \in (0, 1)\}$ . 注意到  $\mathbb{C}$  同构于  $z \in (0, 1)$ , 问题等价于从  $(0, 1)$  上查找分段点, 可采用类似二分查找策略. 相较于被动学习, 可以带来指数级的提升.

## 2 标签复杂度下界

给定分布  $\mathcal{P}$ ,  $\mathbb{C}$  的  $\varepsilon$ -覆盖数  $\mathcal{N}(\varepsilon, \mathcal{P}) \triangleq \min\{|\mathcal{H}| : \bigcup_{h \in \mathcal{H}} B(h, \varepsilon) = \mathbb{C}\}$ .

定理 1  $\forall \mathcal{P}, \forall \Lambda$  (主动学习算法所达到的),  $\forall \varepsilon > 0, \exists \mathcal{R}_{XY}$  与  $\mathcal{P}$  相容, 且

$$\Lambda(\varepsilon, \delta, \mathcal{R}_{XY}) \geq \lceil \log_2((1 - \delta)\mathcal{N}(2\varepsilon, \mathcal{P})) \rceil.$$

考虑噪声情况的下界, 有下面的定理.

定理 2 任给  $\gamma \in (0, 1), \delta \in (1, \frac{1}{4}), n \in \mathbb{N}$ , 令  $p_0 = \frac{1}{2} - \frac{\gamma}{2}, p_1 = \frac{1}{2} + \frac{\gamma}{2}$ . 给定  $\hat{t} : \{0, 1\}^n \rightarrow \{0, 1\}$  (可能是随机的). 如果

$$n < 2 \left\lfloor \frac{1 - \gamma^2}{2\gamma^2} \ln \left( \frac{1}{8\delta(1 - 2\delta)} \right) \right\rfloor,$$

则对  $t \sim \text{Bernoulli}(\frac{1}{2}), B_1|t, \dots, B_n|t \stackrel{iid}{\sim} \text{Bernoulli}(p_t), B_1, \dots, B_n$  与  $\hat{t}$  独立, 则  $\hat{t}(B_1, \dots, B_n) \neq t$  的概率大于  $\delta$ .

定理 3 存在全局常量  $q \in \mathbb{R}^+$  使得: 如果  $|\mathbb{C}| \geq 3$ , 则  $\forall \Lambda, \forall v \in (1, \frac{1}{2})$ , 以及充分小的  $\varepsilon, \delta > 0, \exists \mathcal{R}_{XY}, \text{er}(f^*) = v$ , 使得

$$\Lambda(v + \varepsilon, \delta, \mathcal{R}_{XY}) \geq q \left( \frac{v^2}{\varepsilon^2} \right) \left( d + \text{Log} \left( \frac{1}{\delta} \right) \right),$$

其中  $\text{Log}(x) \triangleq \max\{\ln(x), 1\}, x \geq 0$ . 进一步,  $\forall a \in [4, \infty), \alpha \in (0, 1]$ , 以及充分小的  $\varepsilon, \delta > 0, \exists \mathcal{R}_{XY}$  满足情况 1, 则有

$$\Lambda(v + \varepsilon, \delta, \mathcal{R}_{XY}) \geq qa^2 \left( \frac{1}{\varepsilon} \right)^{2-2\alpha} \left( d + \text{Log} \left( \frac{1}{\delta} \right) \right).$$

### 3 基于分歧的主动学习

#### 3.1 CAL

Cohn, Atlas, Ladner 给出了通用主动学习算法可实现的例子:

---

**Algorithm 1 : CAL ( $n$ ):**

---

```

1:  $m \leftarrow 0, Q \leftarrow \emptyset$ 
2: while  $|Q| < n \ \& \ m < 2^n$ , do
3:    $m \leftarrow m + 1$ 
4:   if  $\forall y \in \mathcal{Y}, \exists h \in \mathbb{C} \text{ s.t. } \text{er}_{Q \cup \{(X_m, y)\}}(h) = 0$ , then
5:     request label  $Y_m$ ;  $Q \leftarrow Q \cup \{(X_m, Y_m)\}$ 
6: return any  $\hat{h} \in \mathbb{C} \text{ s.t. } \text{er}_Q(\hat{h}) = 0$ 

```

---

注意第 4 行的逻辑含义是: 不论对  $X_m$  打上何种标签, 分类器族  $\mathbb{C}$  对  $Q$  形成的分布都有分类器可接受, 这意味着  $X_m$  的标签无法断言, 需要请求  $Y_m$ . 算法 1 可做如下等价:

---

**Algorithm 2 : CAL ( $n$ ):**

---

```

1:  $m \leftarrow 0, t \leftarrow 0, V \leftarrow \mathbb{C}$ 
2: while  $t < n \ \& \ m < 2^n$ , do
3:    $m \leftarrow m + 1$ 
4:   if  $X_m \in \mathcal{D}(V)$ , then
5:     request label  $Y_m$ ;  $V \leftarrow \{h \in V : h(X_m) = Y_m\}$ ;  $t \leftarrow t + 1$ 
6: return any  $\hat{h} \in V$ 

```

---

**定理 4** CAL 达到了标签复杂性  $\Lambda$ , 使得对可实现的  $\mathcal{R}_{XY}$ ,  $\forall \varepsilon, \delta > 0$ , 有

$$\Lambda(\varepsilon, \delta, \mathcal{R}_{XY}) \lesssim \theta(\varepsilon) \left( d \log(\theta(\varepsilon)) + \log\left(\frac{\log(1/\varepsilon)}{\delta}\right) \right) \log\left(\frac{1}{\varepsilon}\right).$$

符号  $\lesssim$  在  $u(\varepsilon, \delta) \lesssim v(\varepsilon, \delta)$  里指  $\exists c \in \mathbb{R}^+$  且与  $\mathbb{C}, \mathcal{R}_{XY}$  或其他问题特定变量都无关, 并满足  $u(\varepsilon, \delta) \leq cv(\varepsilon, \delta)$ ,  $\forall \varepsilon, \delta \in (0, 1)$ . 定理 4 与  $\varepsilon$  渐近相关的界为

$$\mathcal{O}\left(\theta(\varepsilon) \log\left(\frac{1}{\varepsilon}\right) \log\left(\theta(\varepsilon) \log\left(\frac{1}{\varepsilon}\right)\right)\right).$$

**定理 5**  $\forall m \in \mathbb{N} \cup \{0\}, \forall r \in (0, 1), \mathbb{E}[\Delta(V_m^*)] \geq (1-r)^m \Delta(B(f^*, r))$ . 进一步, 这意味着  $\forall \varepsilon \in (0, 1)$ ,

$$\mathbb{E}\left[N\left(\left\lceil \frac{1}{\varepsilon} \right\rceil\right)\right] \geq \frac{\theta(\varepsilon)}{2}.$$

**定理 6**  $\forall n \in \mathbb{N}, \forall r \in (0, 1), \mathbb{E}[\Delta(V_{M(n)}^*)] \geq \Delta(B(f^*, r)) - nr$ . 进一步, 这意味着  $\forall n \in \mathbb{N}, \forall \varepsilon \in (0, 1)$ ,

$$n \leq \frac{\theta(\varepsilon)}{2} \implies \mathbb{E}[\Delta(V_{M(n)}^*)] \geq \frac{\Delta(B(f^*, r))}{2}.$$

考虑

$$\begin{aligned} \{q_m \in \{0, 1\}\}_{m=1}^{\infty} : q_m \perp\!\!\!\perp \mathcal{Z} \mid \{(X_i, q_i Y_i)\}_{i=1}^{m-1}, X_m, \\ \{\hat{h}_m\}_{m=0}^{\infty} : \hat{h}_m \perp\!\!\!\perp \mathcal{Z} \mid \{(X_i, q_i Y_i)\}_{i=1}^m, \end{aligned}$$

随机变量  $q_m$  的取值表示是否请求  $Y_m$ , CAL 作为一种选择性抽样算法, 由

$$(\{q_m\}_{m=1}^{\infty}, \{\hat{h}_m\}_{m=0}^{\infty}) : \hat{h}_m \in V_m^*, q_m = \mathbb{1}_{\mathcal{D}(V_{m-1}^*)}(X_m)$$

确定.

### 3.2 噪声情况

注意到取得最佳分类器  $f^*$ , 仍然有  $\nu = \text{er}(f^*)$ , 下面考虑 CAL 带噪声的情况. 一个具体的算法 (Balcan, Beygelzimer, Langford) 是  $A^2$  策略的变体, 算法如下. 记  $\delta_m \triangleq \frac{\delta}{\log_2^2(2m)}$ ,  $\delta \in (0, 1)$ ,  $m \in \mathbb{N}$ .

---

**Algorithm 3 : RobustCAL $_{\delta}(n)$ :**

---

```

1:  $m \leftarrow 0, i \leftarrow 1, Q_i \leftarrow \emptyset$ 
2: while  $|Q_i| < n \ \& \ m < 2^n$ , do
3:    $m \leftarrow m + 1$ 
4:   if  $\forall y \in \mathcal{Y}, \exists h \in \mathbb{C} \text{ s.t. } h(X_m) = y, \forall j < i, \frac{(\text{er}_{Q_j}(h) - \text{er}_j^*)|Q_j|}{U(2^j, \delta_{(2^j)})} \leq 2^j$ , then
5:     request label  $Y_m$ ;  $Q_i \leftarrow Q_i \cup \{(X_m, Y_m)\}$ 
6:   if  $\log_2(m) \in \mathbb{N}$ , then
7:      $\text{er}_j^* \leftarrow \min\{\text{er}_{Q_j}(h) : h \in \mathbb{C}, \forall j < i, (\text{er}_{Q_j}(h) - \text{er}_j^*)|Q_j| \leq U(2^j, \delta_{(2^j)})2^j\}$ ;
8:      $i \leftarrow i + 1; Q_i \leftarrow Q_{i-1}$ 
9: return any  $\hat{h} \in \mathbb{C} \text{ s.t. } \forall j < i, (\text{er}_{Q_j}(h) - \text{er}_j^*)|Q_j| \leq U(2^j, \delta_{(2^j)})2^j$ 

```

---

同 CAL 一样, 算法 3 也有等价形式 (算法 4).

**定理 7**  $\forall \delta \in (0, 1)$ , RobustCAL $_{\delta}$  达到了标签复杂度  $\Lambda$ , 使得  $\forall \mathcal{R}_{XY}, a, \alpha$  满足情况 1,  $\forall \varepsilon \in (0, 1)$ ,

$$\Lambda(\nu + \varepsilon, \delta, \mathcal{R}_{XY}) \lesssim a^2 \theta(a\varepsilon^{\alpha}) \left(\frac{1}{\varepsilon}\right)^{2-2\alpha} \left(d \log(\theta(a\varepsilon^{\alpha})) + \log\left(\frac{\log(a/\varepsilon)}{\delta}\right)\right) \log\left(\frac{1}{\varepsilon}\right).$$

**Algorithm 4 : RobustCAL<sub>δ</sub>(n):**


---

```

1:  $m \leftarrow 0, Q \leftarrow \emptyset, V \leftarrow \mathcal{C}$ 
2: while  $|Q| < n \ \& \ m < 2^n$ , do
3:    $m \leftarrow m + 1$ 
4:   if  $X_m \in \mathcal{D}(V)$ , then
5:     request label  $Y_m$ ;  $Q \leftarrow Q \cup \{(X_m, Y_m)\}$ 
6:   if  $\log_2(m) \in \mathbb{N}$ , then
7:      $V \leftarrow \left\{ h \in V : \left( \text{er}_Q(h) - \min_{g \in V} \text{er}_Q(g) \right) |Q| \leq U(m, \delta_m) m \right\}$ 
8: return any  $\hat{h} \in V$ 

```

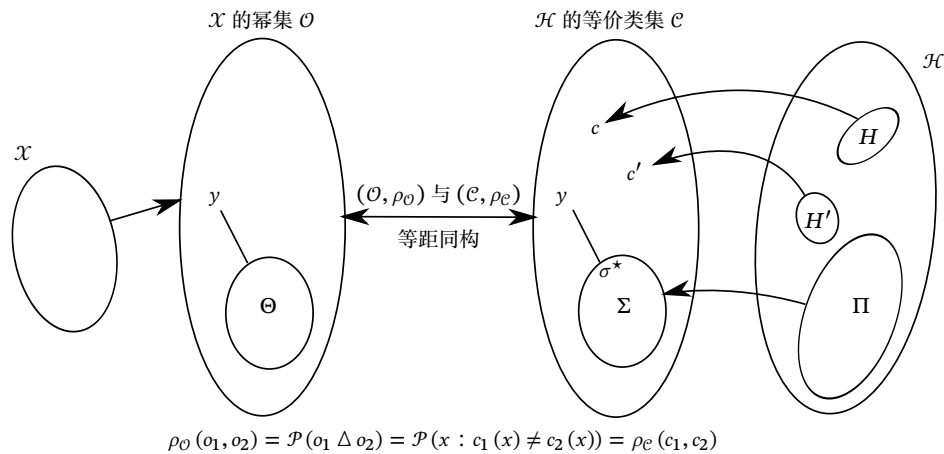
---

进一步,

$$\Lambda(v + \varepsilon, \delta, \mathcal{R}_{XY}) \lesssim \theta(v + \varepsilon) \left( \frac{v^2}{\varepsilon^2} + \text{Log} \left( \frac{1}{\varepsilon} \right) \right) \left( d \text{Log}(\theta(v + \varepsilon)) + \text{Log} \left( \frac{\text{Log}(1/\varepsilon)}{\delta} \right) \right).$$

## 4 注

本质上, 基本定义是在指出分类器集合  $\mathcal{H}$  的拓扑性质. 理想情况, 假设  $\mathcal{H}$  足够大, 大到任何一种  $\mathcal{X}$  的二分类标签情况, 都对应一簇分类器  $H \subset \mathcal{H}$ , 这样, 一方面, 将  $H$  视作等价类集  $\mathcal{C}_{\mathcal{X}, \mathcal{H}}$  的点, 一方面记  $\mathcal{X}$  的幂集  $\mathcal{O}_{\mathcal{X}}$ , 并假设  $\mathcal{X}$  上的测度函数为  $\mathcal{P}$ , 且不存在非空集零测, 这样  $\mathcal{O}_{\mathcal{X}}$  与  $\mathcal{C}_{\mathcal{X}, \mathcal{H}}$  可以构造等距同构. 略去脚标, 如下图.



若  $h \in H$ , 则  $H = \{\forall h' \in \mathcal{H} : h'(x) = h(x), \forall x \in \mathcal{X}\}$ ,  $\mathcal{H} \supset H \mapsto c \in \mathcal{C}$ . 因为是等距同构, 与距离相关的定义都可以对等写出, 如  $\varepsilon$ -闭球.

现在考虑分类器族并不足够大, 不足以打散  $\mathcal{X}$ , 可以用  $\Pi$  来表示实际的分类器族,  $\Pi \subset \mathcal{H}$ , 显然有如图对应的  $\Sigma, \Theta$ , 图中  $y$  是标签, 表示  $\Pi$  没有一个元素可以做到

完全正确的分类的情况, 于是问题等价于求  $y$  对  $\Sigma$  的投影,  $\sigma^* = \arg \inf_{\sigma \in \Sigma} \rho_c(\sigma, y) = \arg \inf_{h \in \Pi} \mathcal{P}(x : h(x) \neq y)$ , 这与噪声率的定义有出入, 噪声率用  $\mathcal{R}_Y$  定义的.

再次审视度量空间, 必须要求不存在非空集零测, 否则  $\rho$  是伪度量, 好在按照度量的理解可以弱化到伪度量上. 注意  $\rho_O(o_1, o_2) = \mathcal{P}(o_1 \Delta o_2) = \mathcal{P}(x : c_1(x) \neq c_2(x)) = \rho_C(c_1, c_2)$  里的对称差  $o_1 \Delta o_2 = o_1 \cup o_2 \setminus o_1 \cap o_2$ , 将其扩展到子集  $O$  上, 有

$$\rho_O(O) = \mathcal{P}\left(\bigcup_{o \in O} o \setminus \bigcap_{o \in O} o\right) = \mathcal{P}(x : \exists c_1, c_2 \in C \text{ s.t. } c_1(x) \neq c_2(x)) = \rho_C(C).$$

这就得到了分歧域和分歧率的定义.