

Passive
Learner

Theorem 10. Suppose that we have a finite set of hypotheses \mathcal{H} (i.e., $|\mathcal{H}| < \infty$) and $\hat{h}_n = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$. Also, assume that the data is separable (i.e., the perfect classifier h^* with no error exists). For any $\epsilon, \delta \in (0, 1)$, we have $\Pr(R(\hat{h}_n) > \epsilon) \leq \delta$ whenever $n \geq \epsilon^{-1} \log(\frac{|\mathcal{H}|}{\delta})$. In other words, for any $\epsilon, \delta \in (0, 1)$, with probability $1 - \delta$, we have $R(\hat{h}_n) \leq \frac{\log(\frac{|\mathcal{H}|}{\delta})}{n}$

$$(x_i, y_i) \in \mathcal{X} \times \{0, 1\} \quad (x_i, y_i) \stackrel{i.i.d.}{\sim} \mathcal{D} \quad i=1, 2, \dots, n$$

Assume $h^* = \arg \min_h R(h)$

and $R(h^*) = 0.$

$$R(h) = \mathbb{E}_{(x, y) \sim \mathcal{D}} [\mathbb{1}\{h(x) \neq y\}]$$

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$$

Active Learner for streaming/sampling-oracle setting.

Separable setting
Assume $\min_h R(h) = 0.$

Algorithm 1 CAL

- 1: Initialize: $Z_0 = \emptyset, V_0 = \mathcal{H}$
 - 2: **for** $t = 1, 2, \dots, n$ **do**
 - 3: Nature reveals unlabeled data point x_t
 - 4: **if** $x_t \in DIS(V_{t-1})$ **then**
 - 5: Query y_t , and set $Z_t = Z_{t-1} \cup (x_t, y_t)$
 - 6: **else**
 - 7: $Z_t = Z_{t-1}$
 - 8: **end if**
 - 9: $V_t = \{h \in \mathcal{H} : h(x_i) = y_i \ \forall (x_i, y_i) \in Z_t\}$ ← "Version space"
 - 10: **end for**
 - 11: **return** any $h \in V_n$
-

Definition 6. For some hypothesis class \mathcal{H} and subset $V \subset \mathcal{H}$ where for each $h \in \mathcal{H}, h : \mathcal{X} \rightarrow \{0, 1\}$, the region of disagreement is defined as

$$DIS(V) = \{x \in \mathcal{X} : \exists h, h' \in V \text{ s.t. } h(x) \neq h'(x)\}$$

which is the set of unlabeled examples x for which there are hypotheses in V that disagree on how to label x .

$$h^* \in V_t \quad \forall t$$

Algorithm 1 CAL

- 1: Initialize: $Z_0 = \emptyset, V_0 = \mathcal{H}$
 - 2: **for** $t = 1, 2, \dots, n$ **do**
 - 3: Nature reveals unlabeled data point x_t
 - 4: **if** $x_t \in DIS(V_{t-1})$ **then**
 - 5: Query y_t , and set $Z_t = Z_{t-1} \cup (x_t, y_t)$
 - 6: **else**
 - 7: $Z_t = Z_{t-1}$ // $y_t = \checkmark h(x_t) = h(x_t) \quad \forall h \in V_{t-1}$
 - 8: **end if**
 - 9: $V_t = \{h \in \mathcal{H} : h(x_i) = y_i \quad \forall (x_i, y_i) \in Z_t\}$
 - 10: **end for**
 - 11: **return** any $h \in V_n$
-

Algorithm 2 Efficient CAL

- 1: Initialize: $Z_0 = \emptyset$
 - 2: **for** $t = 1, 2, \dots, n$ **do**
 - 3: Nature reveals unlabeled data point x_t
 - 4: **if** for $\hat{y} \in \{0, 1\} \exists h_{\hat{y}} \in \mathcal{H} : h_{\hat{y}}(x_s) = y_s, \forall (x_s, y_s) \in Z_{t-1} \cup (x_t, \hat{y})$ **then**
 - 5: Query y_t , and set $Z_t = Z_{t-1} \cup (x_t, y_t)$
 - 6: **else**
 - 7: $Z_t = Z_{t-1}$
 - 8: **end if**
 - 9: **end for**
 - 10: **return** $\arg \min_{h \in \mathcal{H}} \sum_{(x,y) \in Z_t} \mathbf{1}\{h(x) \neq y\}$.
-

$$B(h, r) = \{h' \in \mathcal{H} : \mathbb{E}_x [\mathbf{1}\{h(x) \neq h'(x)\}] \leq r\}$$

Definition 7. The disagreement coefficient of $h \in \mathcal{H}$ with respect to a hypothesis class \mathcal{H} and distribution \mathcal{D}_X is defined as

$$\theta_h^* = \sup_r \frac{\mathbb{P}_{X \sim \mathcal{D}_X}(X \in DIS(B(h, r)))}{r}$$

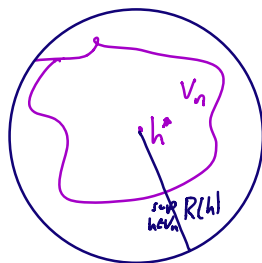
$$\theta^*(\epsilon) = \sup_{r > \epsilon} \dots \leq \frac{1}{\epsilon}$$

Theorem 11. Let $h^* = \arg \min_{h \in \mathcal{H}} R(h)$ and assume $R(h^*) = 0$. Suppose n iid labeled examples $\{(x_i, y_i)\}_{i=1}^n$ are drawn from \mathcal{D} and $V_n = \{h \in \mathcal{H} : h(x_i) = y_i \quad \forall i \in [n]\}$. If we request λ additional labels only when the samples lie in the disagreement region $DIS(V_n)$, where $\lambda = \frac{2\theta_{h^*} \log(|\mathcal{H}|/\delta)}{\delta}$, then, with probability greater than $1 - \delta$ we have $\sup_{h \in V_{n+\lambda}} R(h) \leq \sup_{h \in V_n} \frac{1}{2} R(h)$.

$$\frac{\mathbb{P}_x(X \in DIS(V_n))}{\sup_{h \in V_n} R(h)} \leq \frac{\mathbb{P}(X \in DIS(B(h^*, \sup_{h \in V_n} R(h))))}{\sup_{h \in V_n} R(h)}$$

$$\leq \theta_{h^*}^* \implies \mathbb{P}_x(X \in DIS(V_n)) \leq \theta_{h^*}^* \cdot \sup_{h \in V_n} R(h)$$

$$\mathbb{E}_x [\mathbf{1}\{h(x) \neq h^*(x)\}] = \mathbb{E}_x [\mathbf{1}\{h(x) \neq Y\}] = R(h)$$



$$\sup_{h \in V_{n+\lambda}} R(h) = \sup_{h \in V_{n+\lambda}} P(Y \neq h(x))$$

$$= \sup_{h \in V_{n+\lambda}} P(Y \neq h(x) | x \in \text{DIS}(U_n)) P(x \in \text{DIS}(U_n))$$

$$+ \underbrace{P(Y \neq h(x) | x \notin \text{DIS}(U_n))}_{=0} P(x \notin \text{DIS}(U_n))$$

$\tilde{Y} = h^*(x), h, h^* \in V_n \Rightarrow h(x) = h^*(x)$
 $x \notin \text{DIS}(U_n)$

$$= \sup_{h \in V_{n+\lambda}} \underbrace{P(Y \neq h(x) | x \in \text{DIS}(U_n))}_{\leq \frac{\log(12\lambda/d)}{\lambda}} \underbrace{P(x \in \text{DIS}(U_n))}_{\leq \theta_{h^*}^* \cdot \sup_{h \in V_n} R(h)}$$

$$\leq \frac{\log(12\lambda/d)}{\lambda} \leq \theta_{h^*}^* \cdot \sup_{h \in V_n} R(h)$$

$$\leq \frac{1}{2} \sup_{h \in V_n} R(h). \quad \square$$

Conclude that every λ labels we halve

the risk. If we take

$n = k\lambda$ labels total then

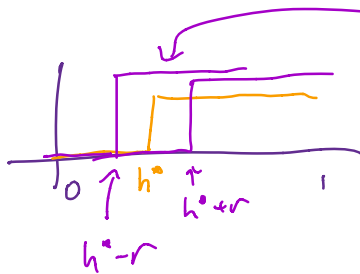
$$\text{w.p. } 1 - k\delta, \quad R(\hat{h}) \leq 2^{-k} = 2^{-n/\lambda}$$

$$\lesssim \exp\left(-\frac{n}{2\theta_{h^*}^* \log(12\lambda/d)}\right)$$

Equivalently, $R(h) \leq \varepsilon$ after $\theta^* \log(12\lambda/d) \log(V/\varepsilon)$ labels.

} cases where $\theta^* = \Theta(1)$.

Ex. \mathcal{H} are thresholds on unit interval (\mathcal{D}_x is uniform(cad))



$B(h^*, r)$

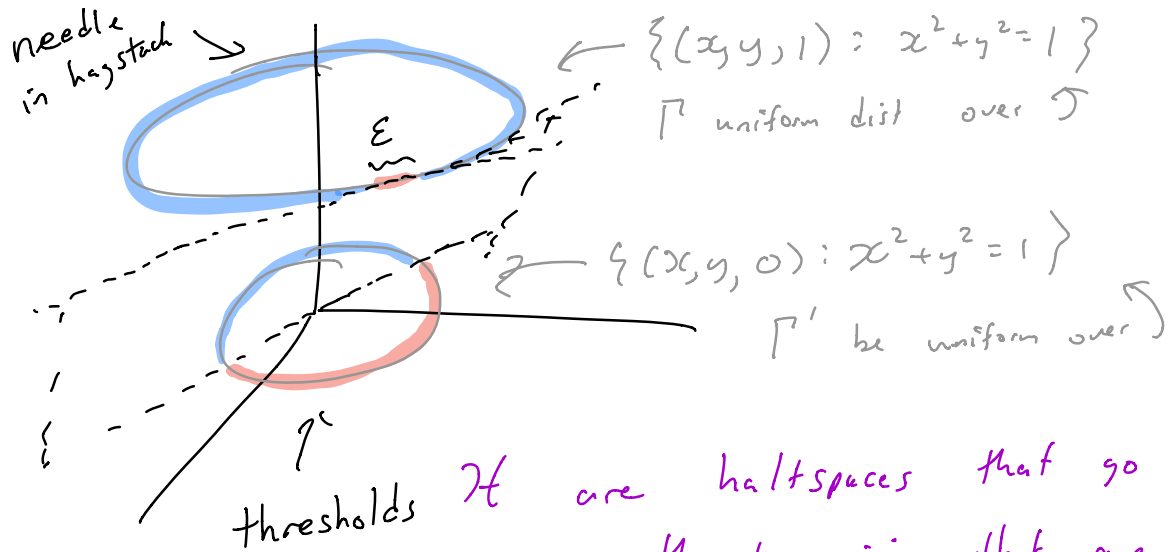
$$\mathbb{P}(X \in \text{DIS}(B(h^*, r))) = 2r$$

$$\Rightarrow \theta^*(\varepsilon) \leq \frac{2r}{r} = 2.$$

If \mathcal{H} halfspaces going through origin

and \mathcal{D}_x is uniform on sphere

$$\text{then } \theta^*(\varepsilon) \leq \sqrt{d}$$



\mathcal{H} are halfspaces that go through origin that are tilted to chop off a tiny bit of Γ .

Assume $h^* \in \mathcal{H}$ w/ $Y = h^*(x)$

$$D_X = (1 - \beta) \Gamma + \beta \Gamma' \quad \text{for } \beta \ll 1$$

If $\beta = 1$ CAL achieves $\log(1/\epsilon)$ labels

If $\beta = 0$ any algorithm requires $1/\epsilon$ labels.

If you run CAL and wait for samples

only from Γ' , requires $\log(1/\epsilon)$ labels

and $\frac{1}{\beta} \log(1/\epsilon)$ unlabeled data

Consider a finite hypothesis space \mathcal{H} and consider any $Q \subset \binom{\mathcal{H}}{2}$ where $(h, h') \in Q$ can be considered an edge connecting any two hypotheses. For any $\hat{y} \in \{0, 1\}$ define $\mathcal{H}_{(x, \hat{y})} = \{h \in \mathcal{H} : h(x) = \hat{y}\}$. We say an example x ρ -splits Q if requesting its label reduces the number of edges by at least a fraction $\rho \in (0, 1)$:

$$\max\{|Q \cap \mathcal{H}_{(x, 0)}|, |Q \cap \mathcal{H}_{(x, 1)}|\} \leq (1 - \rho)|Q|.$$

We are now ready to introduce the splitting index.

Definition 8. Fix any subset $S \subset \mathcal{H}$ and $Q \subset \binom{S}{2}$ such that $\mathbb{P}(h(X) \neq h'(X)) \geq \epsilon, \forall (h, h') \in Q$. Then we say S is (ρ, ϵ, τ) -splittable if $\mathbb{P}(X \text{ splits } Q) \geq \tau$.

Basically, the definition is saying that to reduce the number of pairs of hypotheses that differ by at least ϵ by a fraction at least ρ , requires $1/\tau$ unlabeled data. If \mathcal{H} is finite, and \mathcal{H} is (ρ, ϵ, τ) -splittable, then it is almost immediate that there exists an algorithm that requires $1/(\tau\rho)$ unlabeled data and $1/\rho$ labels to identify an ϵ -good classifier ([\[Dasgupta, 2005b\]](#) suggests one, though it is computationally intractable). What is more important is the reproduced lower bound:

Theorem 12 ([\[Dasgupta, 2005b\]](#)). Fix any hypothesis space \mathcal{H} and distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$. Suppose that for some $\rho \in (0, 1)$, $\epsilon \in (0, 1)$ and some $\tau \in (0, 1/2)$, the set $S \subset \mathcal{H}$ is not (ρ, ϵ, τ) -splittable. Then any active learning strategy that achieves an accuracy of $\epsilon/2$ on all target hypotheses in S must, with probability at least $3/4$ (taken over the random sampling of data), either draw $\geq 1/\tau$ unlabeled samples, or must request $\geq 1/\rho$ labels.