

Hypothesis Class and ERM

Apply ERM over a **restricted set** of hypotheses \mathcal{H} = hypothesis class $\rightarrow \mathcal{H} \in \{ \text{"linear models"}, \text{"SVM"}, \text{"NNs"}, \dots \}$

- each $h \in \mathcal{H}$ is a function $h: \mathcal{X} \rightarrow \mathcal{Y}$

$\text{ERM}_{\mathcal{H}}$ learner:

$$\underline{\text{ERM}_{\mathcal{H}} \in \arg \min_{h \in \mathcal{H}} L_S(h)}$$

\rightarrow we could find multiple best hypothesis

Which hypothesis classes \mathcal{H} do not lead to overfitting?

\hookrightarrow which one are good

Finite Hypothesis Classes

Assume \mathcal{H} is a finite class: $|\mathcal{H}| < \infty$

Let h_S be the output of $\text{ERM}_{\mathcal{H}}(S)$, i.e. $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$

→ Assumptions ←

- **Realizability**: there exists $h^* \in \mathcal{H}$ such that $L_{\mathcal{D},f}(h^*) = 0$
- **i.i.d.**: examples in the training set are independently and identically distributed (i.i.d) according to \mathcal{D} , that is $S \sim \mathcal{D}^m$

→ Very strong assumption

Observation: realizability assumption implies that $L_S(h^*) = 0$

Can we *learn* (i.e., find using ERM) h^* ? → no deterministic guarantee

(Simplified) PAC learning

Probably Approximately Correct (PAC) learning

Since the training data comes from \mathcal{D} :

- we can only be approximately correct → close to be correct
- we can only be probably correct → high chance

Parameters:

- accuracy parameter ϵ : we are satisfied with a good h_S :
 $L_{\mathcal{D},f}(h_S) \leq \epsilon$ → I don't look for ϕ cause I'm not sure it even exists
- confidence parameter δ : want h_S to be a good hypothesis
with probability $\geq 1 - \delta$ → how sure I am that the model is good

We want both ϵ and δ small (~ 0)

Theorem

Let \mathcal{H} be a finite hypothesis class. Let $\delta \in (0, 1)$, $\varepsilon \in (0, 1)$, and $m \in \mathbb{N}$ such that

$$\text{if } \Rightarrow m \geq \frac{\log(|\mathcal{H}|/\delta)}{\varepsilon} \quad \rightarrow \text{size of the training set}$$

[Then for any f and any D] for which the realizability assumption holds, with probability $\geq 1 - \delta$ we have that for every ERM hypothesis h_S it holds that

I can apply this to every situation $L_{D,f}(h_S) \leq \varepsilon$.

Note: \log = natural logarithm

With finite hypothesis classes (#)

I can almost always find a good hypothesis if I have enough data

$$L_{D,f}(h_S) \leq \varepsilon$$

$$m \geq \frac{\log(|\mathcal{H}|/\delta)}{\varepsilon}$$

Quantifies how large the dataset should be

The smaller the ε and δ , the more data I'll need

Proof (see book as well, Corollary 2.3)

Let $S|x = \{x_1, x_2, \dots, x_m\}$ be the instances in the training set S .
We want to bound to:

$$D^m(\{S|x : L_{D,f}(h_S) > \epsilon\})$$

We call $H_0 = \{h \in H : L_{D,f}(h_S) > \epsilon\}$ (bad hypotheses) and

$$M = \{S|x : \exists h \in H_0, L_S(h) = 0\} \text{ (misleading samples)}$$

Since the realizability assumption holds: $L_S(h_S) = 0$

$$L_{D,f}(h_S) > \epsilon \text{ only if some } h \in H_0 \text{ has } L_S(h) = 0$$

That is, our training data must be in the set M :

$$\{S|x : L_{D,f}(h_S) > \epsilon\} \subseteq M$$

$$\text{Note that: } M = \bigcup_{h \in H_0} \{S|x : L_S(h) = 0\}$$

Therefore $D^m(\{S|x : L_D(h_S) > \epsilon\}) \leq D^m(H) = D^m(\bigcup_{h \in H_0} \{S|x : L_S(h) = 0\})$

$D^m(\bigcup_{h \in H_0} \{S|x : L_S(h) = 0\}) \leq \sum_{h \in H_0} D^m(\{S|x : L_S(h) = 0\})$

union bound

Now let's fix $h \in H_0 : L_S(h) = 0 \iff \forall i = 1, \dots, m : h(x_i) = f(x_i)$

Therefore: $D^m(\{S|x : L_S(h) = 0\}) = D^m(\{S|x : \forall i = 1, \dots, m : h(x_i) = f(x_i)\})$

$\stackrel{m}{=} \prod D(\{x_i : h(x_i) = f(x_i)\})$

\nearrow identically distributed from D

Consider some $i, 1 \leq i \leq m : D(\{x_i : h(x_i) = f(x_i)\})$

$\stackrel{!}{=} 1 - D(\{x_i : h(x_i) \neq f(x_i)\})$

$\stackrel{!}{=} 1 - \underbrace{D(\{x_i : h(x_i) \neq f(x_i)\})}_{L_D(h) = P_{x \sim D}[h(x) \neq f(x)]}$

$= 1 - L_D(h) \leq 1 - \epsilon$

\rightarrow since $h \in H_0$

\hookrightarrow Taylor expansion

$L \leq e^{-\epsilon}$

Combining this result with the product of the probabilities:

$$D^m(\{S_{1x} : L_S(h) = 0\}) \leq \prod_{i=1}^m e^{-\epsilon} = e^{-m\epsilon}$$

Combining the above with the sum of the probabilities:

$$D^m(\{S_{1x} : L_{0,\delta}(h_S) > \epsilon\}) \leq \sum_{h \in H_\delta} e^{-m\epsilon} = |H_\delta| e^{-m\epsilon} \leq |H| e^{-m\epsilon}$$

Now, given the choice of m , we have

$$D(\{S_{1x} : L_{0,\delta}(h_S) > \epsilon\}) \leq |H| e^{-\epsilon \left(\frac{\log |H|}{\delta} \right)} = \delta$$



PAC Learning

Definition (PAC learnability)

A hypothesis class \mathcal{H} is *PAC learnable* if there exist a function $m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm such that for every $\delta, \epsilon \in (0, 1)$, for every distribution \mathcal{D} over \mathcal{X} , and for every labeling function $f: \mathcal{X} \rightarrow \{0, 1\}$, if the realizability assumption holds with respect to $\mathcal{H}, \mathcal{D}, f$, then when running the learning algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples generate by \mathcal{D} and labeled by f , the algorithm returns a hypothesis h such that, with probability $\geq 1 - \delta$ (over the choice of examples): $L_{\mathcal{D}, f}(h) \leq \epsilon$.

$m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$: *sample complexity* of learning \mathcal{H} .

- $m_{\mathcal{H}}$ is the minimal integer that satisfies the requirements.

Corollary

Every finite hypothesis class is PAC learnable with sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil$. *Algorithm to find a good hypothesis*

\Downarrow
ERM