

# Compression Schemes

Alireza Kazemipour

CMPUT 654: Theoretical Foundations of Machine Learning

December 7th, 2023

What? Why? How?

# What? PAC Learning $\longrightarrow$ Setting

- Domain set  $\mathcal{X}$ .
- Probability distribution over  $\mathcal{X}$ .  $\mathcal{D}: \mathcal{X} \rightarrow [0, 1]$
- Label set  $\mathcal{Y}$ .  $\mathcal{Y} = \{0, 1\}$
- True labeling function  $f$ .  $f: \mathcal{X} \rightarrow \mathcal{Y}$
- Training data  $\mathcal{S}^m$ .  $\mathcal{S}^m = ((x_1, y_1), \dots, (x_m, y_m))$
- Concept class  $\mathcal{H}$ .  $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$
- The learner's output  $h \in \mathcal{H}$ .  $h: \mathcal{X} \rightarrow \mathcal{Y}$
- Measure of success
$$L_{\mathcal{D}, f}(h) := \mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq f(x)] := \mathcal{D}(\{x : h(x) \neq f(x)\})$$

# What? PAC Learnability $\longrightarrow$ Finite Concept Classes

- Objective:

$$\mathcal{D}^m(\{\mathcal{S}|_x : L_{D,f}(h_{\mathcal{S}}) > \epsilon\}) \leq \delta$$

- So:

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

- $\delta$ : Confidence Parameter
- $\epsilon$ : Accuracy Parameter

# What? PAC Learnability $\longrightarrow$ Finite VC Dimension

- Objective:

$$\mathcal{D}^m(\{\mathcal{S}|_x : L_{D,f}(h_{\mathcal{S}}) > \epsilon\}) \leq \delta$$

- So:

$$m \geq \max\left(\frac{32d}{\epsilon} \log\left(\frac{16e}{\epsilon}\right), \frac{16}{\epsilon} \log\left(\frac{2}{\delta}\right)\right) \quad (1)$$

$$m \geq \max\left(\frac{8d}{\epsilon} \log\left(\frac{8d}{\epsilon}\right), \frac{4}{\epsilon} \log\left(\frac{2}{\delta}\right)\right) \quad (2)$$

- $d$ :  $VCdim(\mathcal{H})$
- $e$ : Euler's number

# What? PAC Learnability $\longrightarrow$ Compression Schemes

## Claim

*Compression Schemes give weaker conditions for PAC learnability  $\iff$  A lower lower bound on the Sample Complexity*

# Why? PAC Learnability $\longrightarrow$ Compression Schemes

*”Can we do better?”*

# How? PAC Learnability $\longrightarrow$ Compression Schemes

- Kernel  $\kappa$ .  $\kappa : \bigcup_{m=k}^{\infty} \mathcal{S}^m \rightarrow \mathcal{S}^k$  (Compressor)
  - Reconstructor  $\rho$ .  $\rho : \mathcal{S}^k \times \mathcal{X} \rightarrow \mathcal{Y} = \{0, 1\}$  (Decompressor)
- 
- $\forall m \geq k$ ,  $\mathcal{S}^k$  is a subsequence of length  $k$  of  $\mathcal{S}^m$
  - $\forall m, x_i \in \mathcal{S}^m|_x \implies \rho(\kappa(\mathcal{S}^m), x_i) = f(x_i)$



- (Previously) Objective:

$$\mathcal{D}^m(\{\mathcal{S}|_x : L_{D,f}(h_{\mathcal{S}}) > \epsilon\}) \leq \delta$$

- (Now) Objective:

$$\mathcal{D}^m(\{\mathcal{S}|_x : L_{D,f}(\rho(\kappa(\mathcal{S}^m), x)) > \epsilon\}) \leq \delta$$

# How? Compression Schemes $\longrightarrow$ Punchline

Let  $T$  be the collection of all  $k$ -element subsequences of the sequence  $(1, 2, \dots, m)$ . For any  $\bar{t} = (t_1, \dots, t_k) \in T$ :

$$A_{\bar{t}} = \{\mathcal{S}^m : \kappa(\mathcal{S}^m) = \mathcal{S}^k\}$$

$$E_{\bar{t}} = \{\mathcal{S} \in A_{\bar{t}} : P(\{x : \rho(\kappa(\mathcal{S}), x) = f(x)\}) < 1 - \epsilon\}$$

$$U_{\bar{t}} = \{\mathcal{S}^m : P(\{x : \rho(\mathcal{S}^k, x) = f(x)\}) < 1 - \epsilon\}$$

$$B_{\bar{t}} = \{\mathcal{S}^m : \text{mark } \rho(\mathcal{S}^k, x_i) = f(x), \forall x_i \text{ s.t. } i \notin t\}$$

$$E_{\bar{t}} = U_{\bar{t}} \cap A_{\bar{t}} \xrightarrow{A_{\bar{t}} \subseteq B_{\bar{t}}} P(E_{\bar{t}}) \leq P(U_{\bar{t}} \cap A_{\bar{t}}) \leq \binom{m}{k} (1 - \epsilon)^{m-k}$$

■ Previously:

$$m \geq \max\left(\frac{8d}{\epsilon} \log\left(\frac{8d}{\epsilon}\right), \frac{4}{\epsilon} \log\left(\frac{2}{\delta}\right)\right)$$


■ Now:

$$m \geq \max\left(\frac{4k}{\epsilon} \log\left(\frac{4k}{\epsilon}\right) + 2k, \frac{2}{\epsilon} \log\left(\frac{1}{\delta}\right)\right)$$

## Theorem

*If a Concept Class is PAC-Learnable then it has a Compression Scheme of size  $k^1$ .*

---

<sup>1</sup>S. Moran and A. Yehudayoff. Sample compression schemes for vc classes. Journal of the ACM (JACM), 63(3):1–10, 2016. 

- ① Why are there different bounds based on  $VCdim$ ?
- ② Why is the Accuracy Parameter  $\epsilon$  turned into Confidence-like Parameter ?
- ③ Why/When is  $k \leq d$ ?
- ④ Is there an *Information Theoretic* approach for the same purpose?

**Thank You! :)**