

Learning theory

UM1 ch. 3

$$S = \{ (\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_m, y_m) \}$$

\mathcal{F} function class, loss function l

ERM

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \underbrace{\frac{1}{m} \sum_{i=1}^m l(f(\underline{x}_i), y_i)}_{\hat{R}(f)}$$

memorization prediction:

$$f_{\text{mem}}(x) = \begin{cases} y_i & \text{if } \exists (\underline{x}_i, y_i) \in S \\ & x = x_i \\ 0 & \text{o.w.} \end{cases}$$

This leads to a 0 ERM,
but not a good predictor.

Generalization:

$$\text{True risk} \quad R(f) = \mathbb{E}_{(\underline{x}, y) \sim \mathcal{D}} [l(f(\underline{x}), y)]$$

Loss over \mathcal{D} = loss over unseen

Example

$$R(\hat{f}) = \mathbb{E}_{(\underline{x}, y) \sim \mathcal{D}} [l(\hat{f}(\underline{x}), y)]$$

$$R(\hat{f}) = \underbrace{R(\hat{f}) - \hat{R}(\hat{f})}_{\text{Generalization gap.}} + \hat{R}(\hat{f})$$

Assume a realizable case, i.e.,

$$\exists f_* \in \mathcal{F} \text{ s.t. } R(f_*) = 0$$

we have a perfect ERM

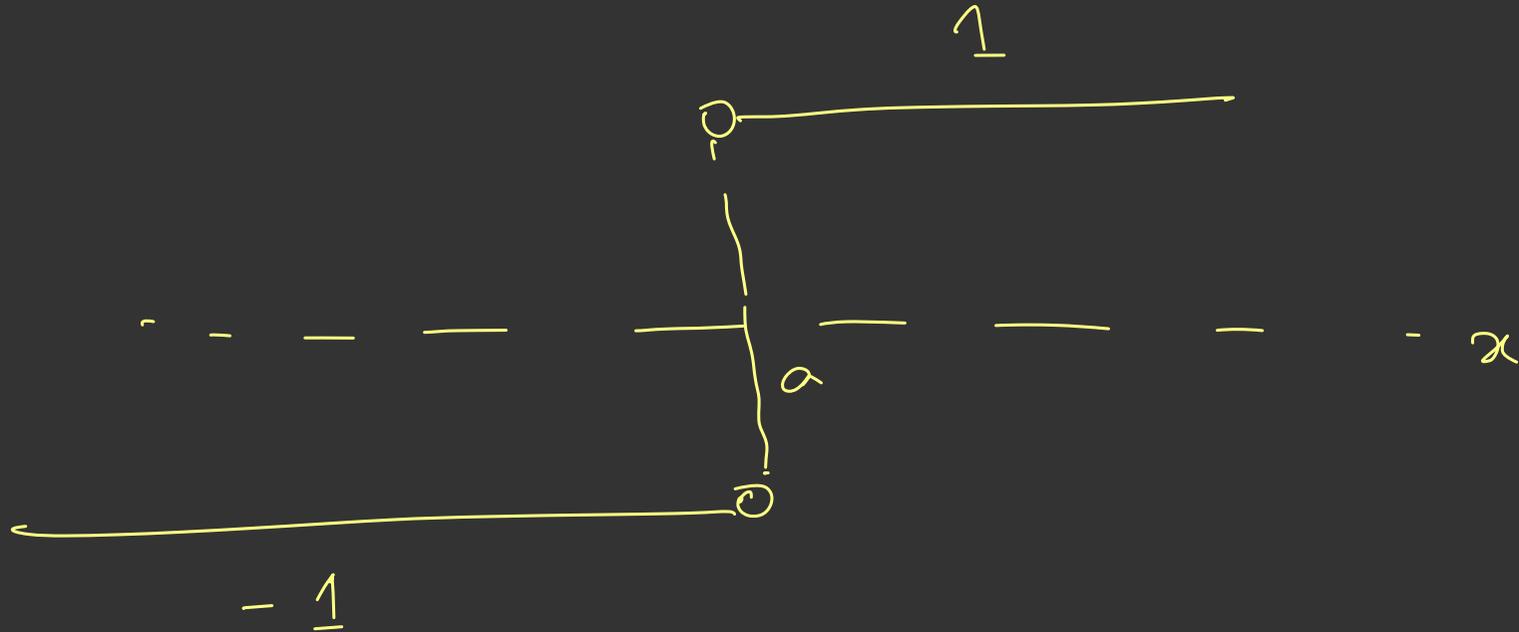
$$\hat{R}(\hat{f}) = 0$$

Q1. Can we find f_* ?

Q2. Can we find a good predictor
for all datasets?

Example - Threshold

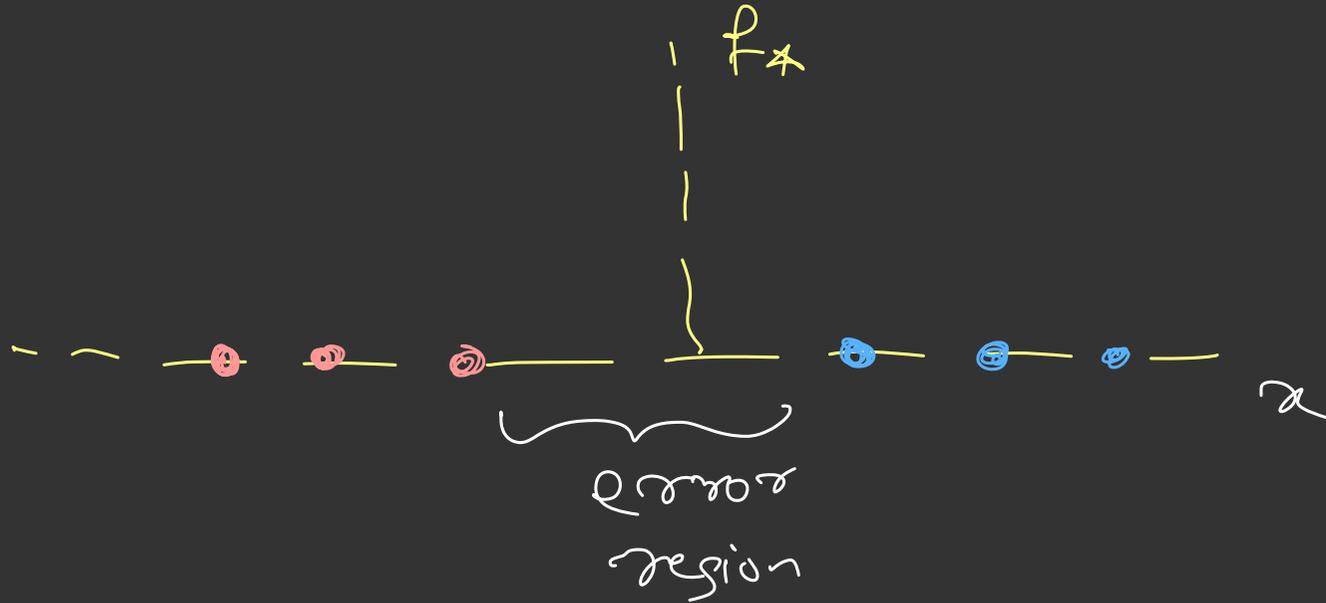
$$f_a(x) = \begin{cases} +1 & \text{if } x \geq a \\ -1 & \text{o.w.} \end{cases}$$



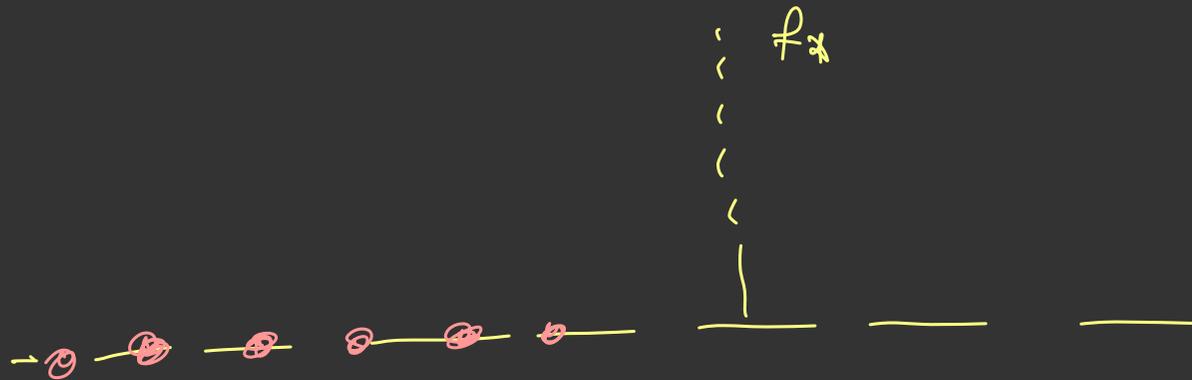
Zero sink

Data is uniformly distributed on the line

Draw 1:



Draw 2:



$$P_0(\text{not seeing } \bullet) = (1-P)^n$$

Probably approximately correct (PAC):

Leslie Valiant (1984)

Error parameter: ϵ

Confidence parameter: δ

$$P_{\sigma} \left[\text{err}(f) \leq \epsilon \right] \geq 1 - \delta$$

Defn: A function class \mathcal{F} is PAC learnable if there exists an algorithm A and a function $m_{\mathcal{F}} : (0, 1)^2 \rightarrow \mathbb{N}$ with the following property:

For every labelling function $f \in \mathcal{F}$,
for every distribution \mathcal{D} on the instance space X , and for all $\epsilon, \delta \in (0, 1)$, if A has access S of size $m \geq m_{\mathcal{F}}(\epsilon, \delta)$, then with probability $1 - \delta$ (over the

choices of the training set), A outputs
a predictor \hat{f} s.t.

$$\mathbb{P}_x \left[\hat{f}(x) \neq f(x) \right] \leq \epsilon.$$

$x \sim \mathcal{D}$

Finite function classes are PAC learnable

Consider $\mathcal{F} = \{f_1, \dots, f_{|\mathcal{F}|}\}$

Binary classification.

Theorem: Every finite function class
is PAC learnable with sample complexity

$$m_{\mathcal{F}}(\epsilon, \delta) \leq \left\lceil \frac{\log(|\mathcal{F}|/\delta)}{\epsilon} \right\rceil$$

Alternatively, \hat{f} is any ERM obtained
over a training set of size m , with

prob. $1-\delta$

$$R(\hat{f}_S) \leq \frac{\log(|F|/\delta)}{m}$$

Proof:

Define $B = \{ f \in \mathcal{F} : f \text{ is } \epsilon\text{-bad} \}$

- we say \hat{f}_S is ϵ -bad if $R(\hat{f}_S) > \epsilon$

- we need $P_{\mathcal{X}} [\hat{f}_S \text{ is } \epsilon\text{-bad}] \leq \delta$

\hookrightarrow ERM

So

$$P_{\mathcal{X}} [\hat{f}_S \text{ is } \epsilon\text{-bad}] = P_{\mathcal{X}} [\hat{f}_S \text{ is } \epsilon\text{-bad} \wedge \hat{R}(\hat{f}_S) = 0]$$

$$\leq \sum_{\mathcal{S}} \mathbb{P}_{\mathcal{R}} [\exists f \in \mathcal{F} : f \text{ is } \epsilon\text{-bad} \wedge \hat{R}(f) = 0]$$

$$= \sum_{\mathcal{S}} \mathbb{P}_{\mathcal{R}} [\exists f \in \mathcal{B} : \hat{R}(f) = 0]$$

$$\leq \sum_{f \in \mathcal{B}} \sum_{\mathcal{S}} \mathbb{P}_{\mathcal{R}} [\hat{R}(f) = 0]$$

[Union bound]

$$= \sum_{f \in \mathcal{B}} \sum_{\mathcal{S}} \mathbb{P}_{\mathcal{R}} [\forall i \in [m] : f(x_i) = f_*(x_i)]$$

$$= \sum_{f \in \mathcal{B}} \prod_{i=1}^m \mathbb{P}_{\mathcal{R}} [f(x_i) = f_*(x_i)]$$

$$= \sum_{f \in \mathcal{B}} \prod_{i=1}^m [1 - \mathbb{P}_{\mathcal{R}} [f(x_i) \neq f_*(x_i)]]$$

$$= \sum_{f \in B} \prod_{i=1}^m [1 - R(f)]$$

$$\leq \sum_{f \in B} (1 - \epsilon)^m$$

$$= |B| (1 - \epsilon)^m \leq |F| (1 - \epsilon)^m$$

$$\boxed{|B| \subseteq F}$$

$$\leq |F| \exp(-m\epsilon)$$

$$\leq \delta \quad := \delta$$

$$1 - a \leq \exp(-a)$$

$$|F| \exp(-m\epsilon) \leq \delta$$

$$\Rightarrow m \geq \frac{\log(|F|/\delta)}{\epsilon} \quad \square$$

Example:

$$F = \left\{ \underline{\omega} \mid \sigma(\underline{\omega}^T \underline{x}) \right\}, \quad \omega_i \in \{-10, \dots, 0, \dots, 10\}$$

$$|F| = 21^d$$

$$m_F(\epsilon, \delta) \leq \frac{d \log 21 + \log(1/\delta)}{\epsilon}$$

Are infinite classes PAC learnable?

Growth function:

$$\Pi_F(S) = \left| \{ f(x_1), \dots, f(x_m) \} \right|$$

All possible labellings that the training points could have according to F

Growth function: (max. possible labelling over all sets of size m)

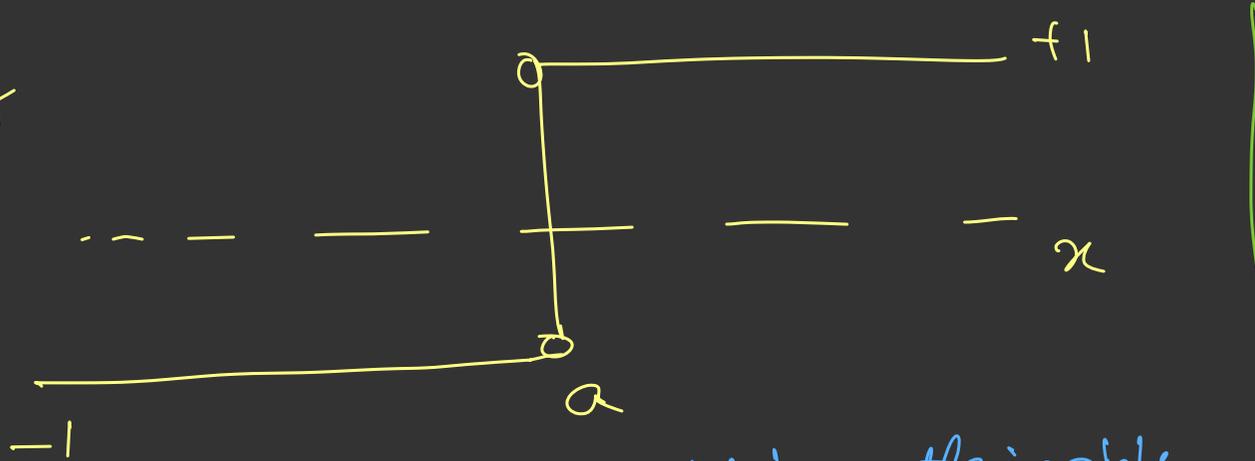
$$\Pi_F(m) := \max_{S: |S|=m} |\Pi_F(S)|$$

Example:

$$f_a(x) = \begin{cases} +1 & \text{if } x \geq a \\ -1 & \text{o.w.} \end{cases}$$

Consider a data set with $x_1 < x_2 < x_3$

$$\Pi_F(3) = 4$$



more generally

$$\Pi_F(m) = m + 1$$

$$\ll 2^m$$

Not attainable

$$-1, -1, -1$$

$$-1, -1, +1$$

$$-1, +1, +1$$

$$+1, +1, +1$$

$$-1, +1, -1$$

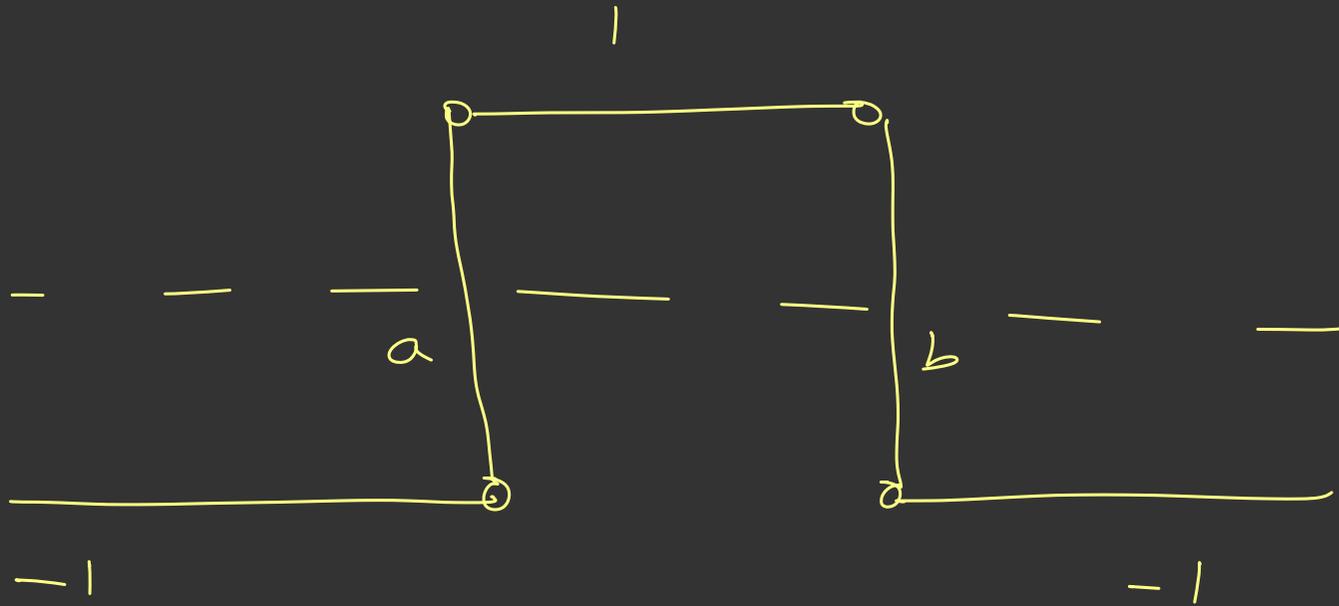
$$1, -1, +1$$

$$1, 1, -1$$

$$1, -1, -1$$

Example:

$$f_{a,b}(x) = \begin{cases} +1 & x \in [a,b) \\ -1 & \text{o.w.} \end{cases}$$



w.l.o.g

$$x_1 < x_2 < x_3$$

only $1, -1, 1$ is not attainable

$$\pi_{\mathbb{F}}(3) = 7$$

$$\pi_{\mathbb{F}}(m) = \frac{m(m+1)}{2} + 1$$

General bound on $\Pi_F(m)$.

VC dimension of F



Vapnik - Chervonenkis (1970)

Def (Shattering):

A set S of inputs m said to be shattered
by a function class F if $|\Pi_F(S)| = 2^{|S|}$.

Def (VC dimension) :

VC dimension of a function class \mathcal{F} or
 $VC(\mathcal{F})$ is the size of the largest set S
that can be shattered by \mathcal{F} .

To show that a function class has $VC(\mathcal{F}) = d$,
we must show that

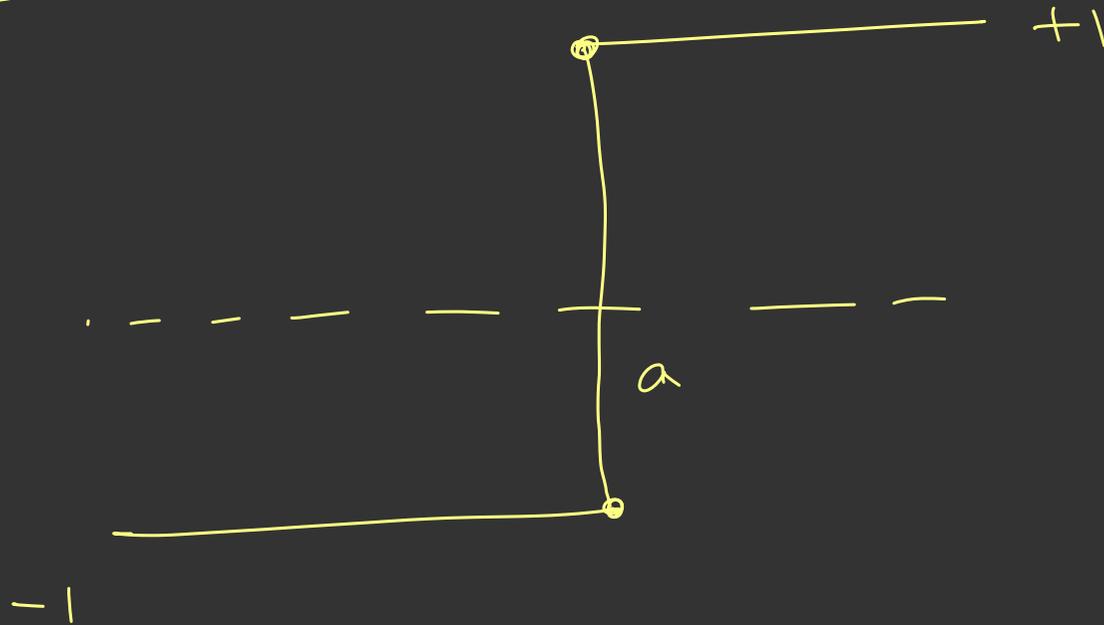
★ There is a set S of d points that is
shattered by \mathcal{F}

★ There is not set S of $d+1$ points that
is shattered by \mathcal{F} .

Examples:

$$f_a(x) = \begin{cases} +1 & \text{if } x \geq a \\ -1 & \text{o.w.} \end{cases}$$

① Threshold:



Any set of 1 can be shattered,
however, no set of size 2 can be shattered.

Thus, $VC(F) = 1$

(2) Intervals:

$$VC(F) = 2$$

(3) Half spaces: ?