

STATISTICAL LEARNING THEORY

Summer Term 2024

Resource: Shai Shalev-Schwarz

Understanding Machine Learning

Online Notes: <http://tkwitt.org> → Teaching
(All my handwritten notes will be available there!)

OTHER COURSE NAMES:

- 1) Advanced machine learning
- 2) Machine learning

MOTIVATION

A motivating example:

	Weight (in g)	Color ($\in \{0, 1\}$)	Tasty
Papaya 1	800	0.1	0
⋮	⋮	⋮	⋮
Papaya N	1200	0.9	1

0... Not tasty
1... Tasty

Let's call the information on this table our **training data**. Based on that we aim to find

$$h: \mathbb{R} \times \mathbb{R} \rightarrow \{0, 1\}$$

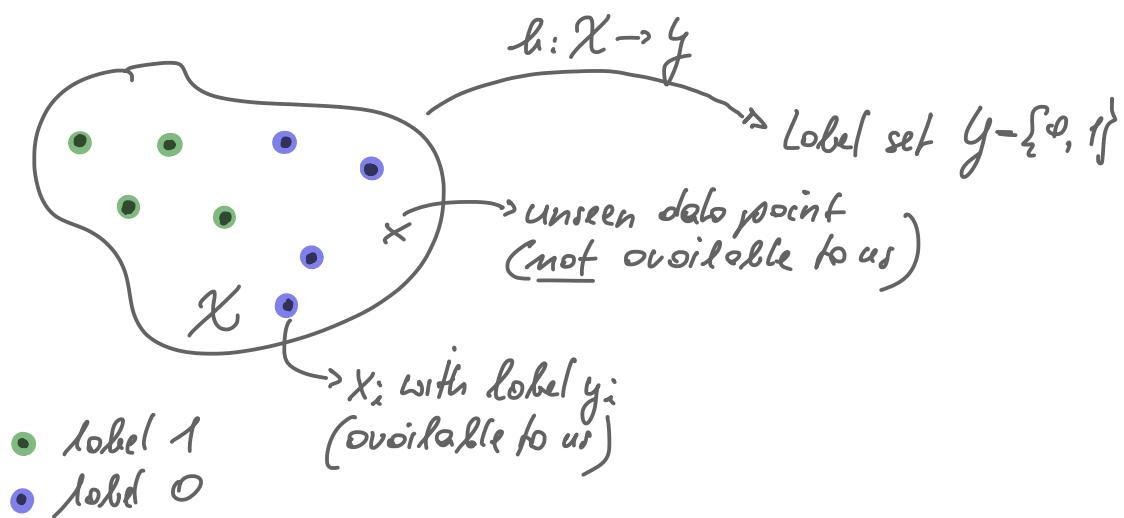
i.e., a function that will take a two-dim. vector as input (weight and color) and output a prediction of whether a papaya is tasty (1) or not (0). We call such a function a **hypothesis**.

More formally, the data available to us comes in the form

$$((x_1, y_1), \dots, (x_N, y_N)) = S$$

with $x_i \in \mathbb{R}^2$ and $y \in \{0, 1\}$. \hookrightarrow Label set $Y = \{0, 1\}$
 \hookrightarrow Domain $X = \mathbb{R}^2$

What do we mean by "learning"



We say, a **learner** receives S and outputs h !
(some algorithm)

Two assumptions we will make initially:

(1) All the x_i 's are drawn

identically and independently (iid)

from some (unknown) distribution D over the domain \mathcal{X} .

(2) The x_i 's are labeled by some (unknown) function

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

which we call the true labeling function, meaning

$$S = ((\underbrace{x_1, f(x_1)}_{y_1}), \dots, (\underbrace{x_N, f(x_N)}_{y_N}))$$

What do we care about?

We care about

$$A = \{x \in \mathcal{X} : h(x) + f(x)\}$$

That is, all the points x in our domain \mathcal{X} , where the hypothesis h differs from the true labeling function f .

BACKGROUND

(a really short overview
on required concepts)

In the following, $\mathcal{P}(X)$ will denote the power set of X .

Remember that we care about

$$A = \left\{ x \in X : h(x) + f(x) \right\}$$

(\hookrightarrow domain)

In case of $X = \mathbb{R}^n$ and $A \in \mathcal{P}(X)$, we already have a problem.

Measure Problem: find $\mu: \mathcal{P}(\mathbb{R}^n) \rightarrow [0, \infty]$ with properties:

1. $A_i \in \mathcal{P}(\mathbb{R}^n), i \in \mathbb{N}$, pairwise disjoint, we want

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i)$$

2. if $A, B \in \mathcal{P}(\mathbb{R}^n)$ are congruent, we want

$$\mu(A) = \mu(B)$$

3. we want

$$\mu([0, 1]^n) = 1$$

The measure problem is unsolvable for all $n \in \mathbb{N}$!

Solution: We will constrain ourselves to elements of some σ -Algebra over the domain X .

Def. (σ -Algebra): Let S be a non-empty set. A family of sets $\mathcal{F} \subset P(S)$ is called a σ -Algebra over S , if

1. $S \in \mathcal{F}$

2. from $A \in \mathcal{F}$, it follows that $A^c = S \setminus A \in \mathcal{F}$

3. from $A_i \in \mathcal{F}$ with $i \in \mathbb{N}$, it follows that

$$\bigcup_{i=1}^{\infty} A_i \in \mathcal{F} \quad (\text{i.e. for any countable collection})$$

Example(s) (smallest σ -Alg. that contains A):

$$\{\emptyset, A, A^c, S\}$$

$$S = \{q, b, c, d\} \text{ with } \sigma\text{-Algebra } \left\{ \emptyset, \{q, b\}, \{c, d\}, \{q, b, c, d\} \right\}$$

Def. (generator): Given $\mathcal{E} \subset P(S)$ a family of sets and Σ denoting the set of all σ -Algebras that contain \mathcal{E} , we call

$$\sigma(\mathcal{E}) = \bigcap_{\mathcal{F} \in \Sigma} \mathcal{F}$$

the σ -Alg. generated by \mathcal{E} . Further, if it holds that for a σ -Alg. \mathcal{A}

$$\sigma(\mathcal{E}) = \mathcal{A}$$

we call \mathcal{E} the generator of \mathcal{A} .

Example : $\mathcal{E} = \{\{\{1\}\}\}$, $\mathcal{S} = \{1, 2, 3\}$

$$\begin{aligned}\mathcal{O}(\mathcal{E}) &= \mathcal{O}(\{\{\{1\}\}\}) \\ &= \{\emptyset, \{\{1\}\}, \{\{1, 2\}\}, \{\{1, 2, 3\}\}\}\end{aligned}$$

Def. (Topological Space) : A topological space is a tuple (X, τ) with X a set and τ a collection of subsets of X s.t.

1. $\emptyset \in \tau$, $X \in \tau$

2. closed under union, i.e., $\{\overline{T_i}\}_{i \in I} \subseteq \tau \Rightarrow \bigcup_{i \in I} \overline{T_i} \in \tau$

3. closed under finite intersection, i.e.,

$$\{\overline{T_i}\}_{i=1}^n \subseteq \tau \rightarrow \bigcap_{i=1}^n \overline{T_i} \in \tau$$

Remark: The elements of τ are called open sets.

Example: $X = \{1, 2, 3, 4\}$

$$\tau = \{\emptyset, \{1, 2, 3, 4\}, \{2\}, \{1, 2\}, \{2, 3\}, \{1, 2, 3\}\}$$

Def. (BOREL σ-Alg.): Given a topological space (S, \mathcal{Q}) with \mathcal{Q} denoting the system of open sets, then we call

$$\sigma(\mathcal{Q}) =: \mathcal{B}(S)$$

the BOREL σ-Algebra over S. Its elements are called BOREL sets. For $S = \mathbb{R}^n$, we write

$$\mathcal{B}^n := \mathcal{B}(\mathbb{R}^n)$$

Importantly, each of the following systems of sets are generators for \mathcal{B}^n :

① $\{U \subset \mathbb{R}^n : U \text{ open}\}$

② $\{A \subset \mathbb{R}^n : A \text{ closed}\}$

③ $\{[q, b] : q, b \in \mathbb{R}^n \text{ with } q \leq b\}$

④ $\{[-\infty, c] : c \in \mathbb{R}^n\}$

meantif
 $[q, b] = [q_1, b_1] \times \dots \times [q_n, b_n]$
 for $q = (q_1, \dots, q_n)$
 $b = (b_1, \dots, b_n)$

Convention: We extend \mathbb{R} by symbols " $+\infty$ ", " $-\infty$ " as

$$\overline{\mathbb{R}} = \mathbb{R} \cup \{+\infty, -\infty\}$$

$$\overline{\mathcal{B}} = \sigma(\mathcal{B} \cup \{-\infty\} \cup \{+\infty\})$$

Def. (Measurable space): if \mathcal{F} is a σ -Algebra over S , we call (S, \mathcal{F}) a measurable space.

Example: $(\mathbb{R}^n, \mathcal{B}^n)$

Def. (Measurable function): Given (S_1, \mathcal{F}_1) and (S_2, \mathcal{F}_2) measurable spaces, we call

$$f: S_1 \rightarrow S_2$$

a measurable function if

$$\forall E \in \mathcal{F}_2 : f^{-1}(E) \subset \mathcal{F}_1$$

Example: $\mathbb{1}_A : S \rightarrow \mathbb{R}$ // characteristic function

$$S \ni \omega \mapsto \mathbb{1}_A(\omega) = \begin{cases} 1, & \text{if } \omega \in A \\ 0, & \text{else} \end{cases}$$

Consider $\mathbb{1}_A$ as a function mapping to \mathbb{R} :

$$\text{if } B \leq 0, \text{ then } \{\omega : \mathbb{1}_A(\omega) < B\} = \emptyset$$

$$\text{if } B > 1, \text{ then } \{\omega : \mathbb{1}_A(\omega) < B\} = S$$

$$\text{if } 0 < B \leq 1, \text{ then } \{\omega : \mathbb{1}_A(\omega) < B\} = S \setminus A = A^c$$

$\Rightarrow \mathbb{1}_A$ is measurable if $A \in \mathcal{F}$.

↑
all
measurable

Def. (Measure): let (S, \mathcal{F}) be a measurable space.
A function $\mu: \mathcal{F} \rightarrow \overline{\mathbb{R}}$ is called a measure if the
following conditions hold:

1. $\mu(\emptyset) = 0$

2. $\mu(A) \geq 0$ for all $A \in \mathcal{F}$

3. for every sequence $(A_n)_{n \in \mathbb{N}}$ of disjoint
sets from \mathcal{F} , we have

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i) \quad // \text{σ-Additivity}$$

Example (counting measure): $\mu_{\text{count}}: \mathcal{F} \rightarrow \overline{\mathbb{R}}$

$$A \mapsto \begin{cases} |A|, & \text{if } A \text{ is finite} \\ \infty, & \text{else} \end{cases}$$

Def. (Measure space): if (S, \mathcal{F}) is a measurable space
and $\mu: \mathcal{F} \rightarrow \overline{\mathbb{R}}$ is a measure, we call
 (S, \mathcal{F}, μ)
a measure space.

Some elementary properties : Let (S, \mathcal{F}, μ) be a measure space. Also, let $A, B \in \mathcal{F}$ and $A_n \in \mathcal{F}$ with $n \in \mathbb{N}$. Then, it holds that

1. if A and B are disjoint, then

$$\mu(A \cup B) = \mu(A) + \mu(B)$$

2. if $A \subset B$ and $\mu(A) < \infty$, then

$$\mu(\underbrace{B \setminus A}_{\text{set difference}}) = \mu(B) - \mu(A)$$

3. if $A \subseteq B$, then

$$\mu(A) \leq \mu(B)$$

4. $\mu\left(\bigcup_{n=1}^{\infty} A_n\right) \leq \sum_{n=1}^{\infty} \mu(A_n)$ // called σ -Additivity



Def. (Probability space): Given (S, \mathcal{F}, D) a measure space and $D(S) = 1$,

then we call (S, \mathcal{F}, D) a probability space, and D a probability measure.

Remarks: S ... interpreted as the set of outcomes of a random experiment

\mathcal{F} ... interpreted as the set of events to which we want to assign probabilities to.

D ... the prob. measure assigns to each event $E \in \mathcal{F}$ a probability $D(E) \in [0, 1]$

Another remark: Based on $(S_1, \mathcal{F}_1, \mu)$ a measure space, a measurable space (S_2, \mathcal{F}_2) and a measurable function $f: S_1 \rightarrow S_2$, we can easily construct a new measure ν :

$$\nu_f: \mathcal{F}_2 \rightarrow [0, \infty]$$

$$B \mapsto \nu_f(B) = \mu(f^{-1}(B))$$

pre-image of B under f

We call ν_f the "push-forward" measure of μ under f .

Def. (Random Variable): Given $(S_1, \mathcal{F}_1, \mathbb{P})$ a prob. space and (S_2, \mathcal{F}_2) a measurable space, we call a measurable function

$$X: S_1 \rightarrow S_2$$

a random variable. If $S_2 = \mathbb{R}$, we say X is a real random var.

Also, the push-forward measure \mathbb{P}_X on S_2 is a prob. measure, since

$$\mathbb{P}_X(S_2) = \mathbb{P}\left(\underbrace{X^{-1}(S_2)}_{S_1}\right) = \mathbb{P}(S_1) = 1$$

We call \mathbb{P}_X the distribution of X .

CONVENTION(s) on notation:

X ... random variable

$$\{X \in A\} = \{\omega \in S : X(\omega) \in A\}$$

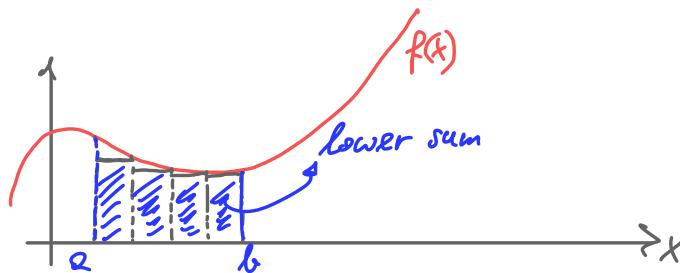
$$\{X = c\} = \{\omega \in S : X(\omega) = c\}$$

$$\mathbb{P}_X(A) = \mathbb{P}\left(\{\omega \in S : X(\omega) \in A\}\right)$$

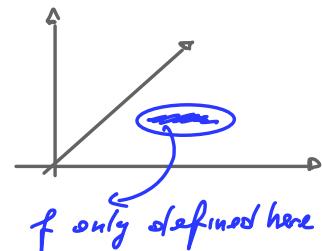
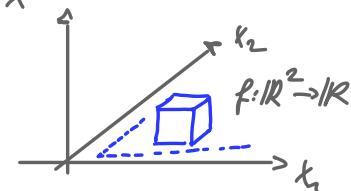
$$\mathbb{P}_X(\{c\}) = \mathbb{P}\left(\{\omega \in S : X(\omega) = c\}\right)$$

A very short primer on the LEBESGUE INTEGRAL

Key idea of Riemann integral: $f: \mathbb{R} \rightarrow \mathbb{R}$

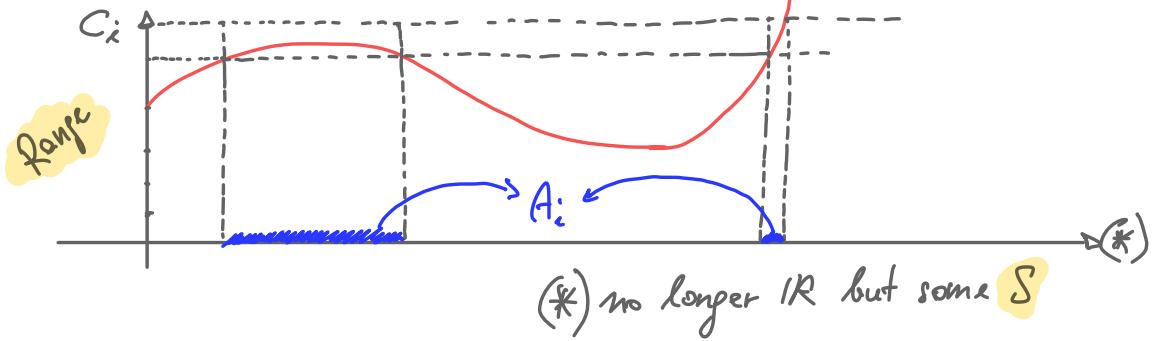


- 1) Difficult to extend to higher dimensions
- 2) Reliance on continuity
- 3) :



Question: How to partition this domain?

Idea of LEBESGUE interpretation: Partition the range of a function $f: S \rightarrow \mathbb{R}$ into intervals to get to an approximation by so called "elementary functions".



We need a way to measure A_i ; If we can do that, we can write

$$c_i \cdot \mu(A_i)$$

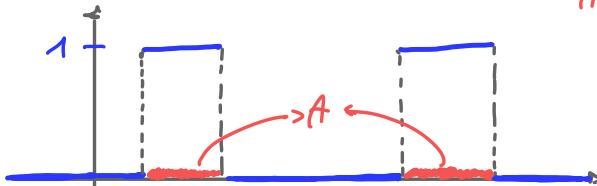
measure

.... this allows us to eventually write

$$\sum_i c_i \cdot \mu(A_i) \quad \text{OR} \quad \int_S f d\mu$$

A little bit more detail: Let's look at the case of "simple" functions. So, we have (S, \mathcal{F}, μ) , $\mu: \mathcal{F} \rightarrow [0, \infty]$

$$\mathbb{1}_A: S \rightarrow \mathbb{R}, A \in \mathcal{F}$$

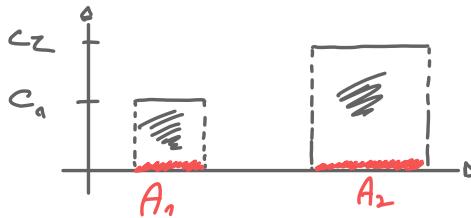


Any reasonable notion of interpretation of $\mathbb{1}_A$ should return $\mu(A)$:

$$\text{Integral} \quad \mathbb{I}(\mathbb{1}_A) = \mu(A)$$

finite here!

Now, a "simple" function can be written as $f(x) = \sum_{i=1}^n c_i \cdot \mathbb{1}_{A_i}(x)$ with $A_i \in \mathcal{F}$ and $c_i \in \mathbb{R}$. Visually:



We can define the integral of f via: $\mathbb{I}(f) = \sum_{i=1}^n c_i \mu(A_i)$. But, since $\mu: \mathcal{F} \rightarrow [0, \infty]$ and $c_i \in \mathbb{R}$, we could get something like

$$10 \cdot \infty - 8 \cdot \infty$$

Solution: $\{f: S \rightarrow \mathbb{R} \mid f \text{ is "simple" and } f \geq 0\} = T^+$. For $f \in T^+$ we have a well-defined representation of f as

$$f(x) = \sum_{i=1}^n c_i \cdot \mathbb{1}_{A_i}(x) \text{ with } c_i \geq 0$$

and we define the LEBESGUE integral (for this class) as:

$$\mathbb{I}(f) := \sum_{i=1}^n c_i \cdot \mu(A_i) \quad \text{or} \quad \int f d\mu \quad \text{or} \quad \int f(x) d\mu(x) \quad x \in S$$

Some properties: (that follow immediately) **IMPORTANT!**

for $f, g \in T^+$ and $\alpha, \beta \geq 0$

(A) $\int(\alpha f + \beta g) d\mu = \alpha \cdot \int f d\mu + \beta \cdot \int g d\mu$ LINEARITY

for $f, g \in T^+$ and $f \leq g$

(B) $\int f d\mu \leq \int g d\mu$ MONOTONICITY

Overall, even though we just looked at "simple" functions, the construction extends to all measurable functions $f: S \rightarrow \mathbb{R}$.

(also (A) and (B) hold for all measurable functions)

Def. (Expected value of a RV): Given a prob. space $(S, \mathcal{F}, \mathbb{P})$ and a \mathbb{P} -integrable real RV $X: S \rightarrow \mathbb{R}$, then

$$\mathbb{E}[X] = \int X d\mathbb{P}$$

is called the expected value of X .

We have:

- If $X \geq 0$, then $\mathbb{E}[X] \geq 0$

Linearity

- $\mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for RV's X and Y
- $\mathbb{E}[\alpha X] = \alpha \mathbb{E}[X]$ for constant α

Def. (MARKOV inequality): Given prob. space $(S, \mathcal{F}, \mathbb{P})$ and a non-negative RV $X (\geq 0)$, we have for $Q > 0$

$$\underbrace{\mathbb{P}[X \geq Q]}_{\mathcal{D}(\{ \omega \in S : X(\omega) \geq Q \})} \leq \frac{\mathbb{E}[X]}{Q}$$

$$\mathcal{D}(\{ \omega \in S : X(\omega) \geq Q \})$$

Proof: Define $\varphi: S \rightarrow \mathbb{R}$ with $\varphi(x) = \begin{cases} Q, & \text{if } X \geq Q \\ 0, & \text{else } (X < Q) \end{cases}$

$\rightarrow 0 \leq \varphi(x) \leq X(x)$. Hence, by monotonicity

we have $\int X d\mathbb{P} \geq \int \varphi d\mathbb{P} = Q \cdot \mathcal{D}(\{ \omega \in S : X(\omega) \geq Q \})$.
 \Rightarrow as $Q > 0$, we have $\frac{1}{Q} \cdot \int X d\mathbb{P} \geq \mathcal{D}(\{ \dots \})$

and since $D(\{\omega \in S : X(\omega) \geq Q\}) = P[X \geq Q]$, we
get $P[X \geq Q] \leq \frac{\int X dD}{Q} = \frac{E(X)}{Q}$ □

CONTINUING ML MATERIAL

1. **Domain** set X ; we call $x \in X$ an instance
2. **Label** set Y , e.g., $Y = \{0, 1\}$
3. **Training data** set $S = ((x_1, y_1), \dots, (x_m, y_m))$ with $(x_i, y_i) \in X \times Y = \mathcal{Z}$
4. A **learner** that receives S and outputs $h: X \rightarrow Y$ which we call a hypothesis.

Assumption: For now, we assume the x_i 's are drawn iid from some probability measure D over the domain X and labeled by some unknown function $f: X \rightarrow Y : y_i = f(x_i)$

We are interested in:

$$D\left(\{x \in X : h(x) \neq f(x)\}\right) = \mathbb{P}_{x \sim D} [h(x) \neq f(x)] = L_{D,f}(h)$$

"Generalization error"

The empirical version of this is

$$\frac{1}{m} \cdot \left| \{i \in [m] : h(x_i) \neq f(x_i)\} \right| = L_S(h)$$

"Empirical error"
(empirical risk)

Convention: $S|_x = (x_1, \dots, x_m)$

Claim: $\mathbb{E}_{S|x \sim D^m} [L_S(h)] = L_{D,f}(h) \quad (\mathbb{E} \dots \text{expected value})$

$$\mathbb{E}_{S|x \sim D^m} [L_S(h)] = \mathbb{E}_{S|x \sim D^m} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq f(x_i)} \right] // \text{by def.}$$

empirical error

$$= \frac{1}{m} \cdot \sum_{i=1}^m \mathbb{E}_{x_i \sim D} [\mathbb{1}_{h(x_i) \neq f(x_i)}] // \text{by linearity}$$

$$= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{x \sim D} [\mathbb{1}_{h(x) \neq f(x)}] // \text{as all } x_i's \text{ are drawn iid from } D$$

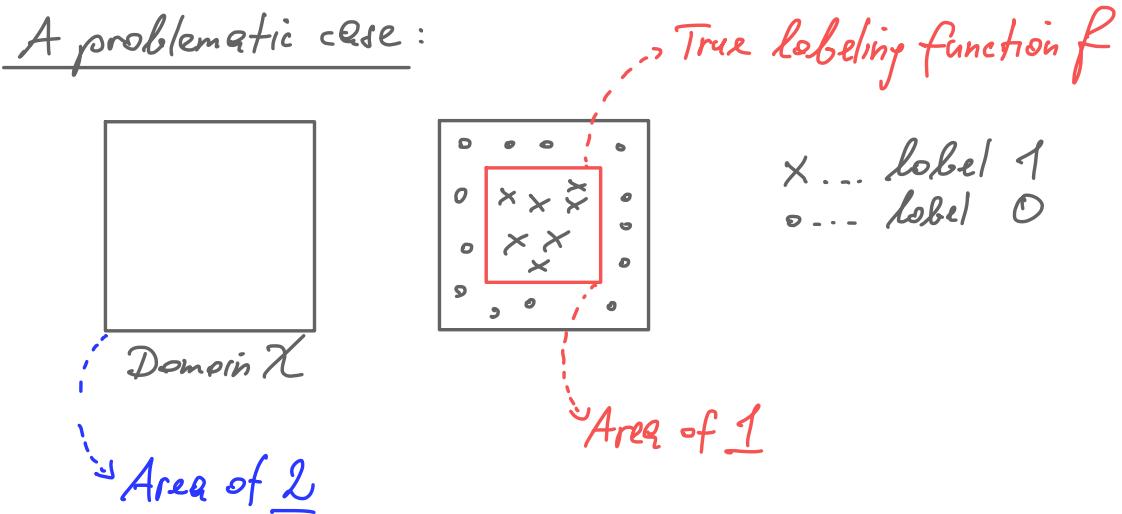
$$= \frac{1}{m} \sum_{i=1}^m \mathbb{P}_{X \sim D} [h(x) \neq f(x)]$$
$$= \cancel{\frac{1}{m} \cdot m} \cdot \mathbb{P}_{X \sim D} [h(x) \neq f(x)]$$

$= L_{D,f}(h)$ generalization error

This establishes the claim.

Our first learning paradigm

Empirical risk minimization (ERM): As we only have access to the training data (S), it's natural to try to select h such that the empirical risk (empirical error) is minimized. We call such a h an empirical risk minimizer.



- Say the distribution on X is uniform
- Say we have an ERM algorithm that returns h_S s.t.

$$h_S(x) = \begin{cases} y_i, & \text{if } \exists i \in \{1, \dots, m\} : x_i = x \\ 0, & \text{else} \end{cases}$$

// i.e., a lookups table

Obviously, h_S is correct on all instances in $S \rightarrow h_S$ is an empirical risk minimizer, meaning

$$L_S(h_S) = 0!$$

But, on unseen instances from our distribution on X (uniform), h_S is only correct 50% of the time, i.e.,

$$L_{D,F}(h_S) = \frac{1}{2} !$$

This is called overfitting.

Hypothesis class (H): We restrict searching for h to H , i.e., a class of functions from $X \rightarrow Y$ and we write

$$\text{ERM}_H(S) \in \underset{h \in H}{\operatorname{arg\,min}} L_S(h)$$

ERM over finite hypothesis classes $(|H| < \infty)$

Assumption (realizability): $\exists h^* \in H$ with $L_{D,f}(h^*) = 0$.

Now, any ERM hypothesis h_S will attain 0 empirical error ($L_S(h_S) = 0$) as h_S competes with h^* (which obviously has 0 empirical error).

Hence, $L_{D,f}(h_S) \geq \varepsilon$ can only happen if we select a hypothesis with $L_S(h_S) = 0$ but $L_{D,f}(h_S) \geq \varepsilon$.

We define

$$H_{\text{BAD}} = \{h \in H : L_{D,f}(h) \geq \varepsilon\} \quad / \begin{array}{l} \text{set of BAD} \\ \text{hypothesis} \end{array}$$

Further, we define

$$M = \{S|x : \exists h \in \underbrace{H_{\text{BAD}}}_{L_S(h)=0}, L_S(h) = 0\}$$

these are the ones with generalization error of $\geq \varepsilon$

Observation:

$$\{S|x : L_{D,f}(h_S) \geq \varepsilon\} \subseteq \underbrace{\{S|x : \exists h \in H_{\text{BAD}}, L_S(h) = 0\}}_M \quad (\text{from before})$$

(A) empirical risk minimizer

$$M = \bigcup_{h \in H_{\text{BAD}}} \{S|x : L_S(h) = 0\}$$

We get (upon measuring with D):

$$D^m \left(\{S|_k : L_{D,f}(h_S) \geq \varepsilon\} \right) \leq D^m \left(\bigcup_{h \in H_{BAD}} \{S|_k : L_S(h) = 0\} \right)$$

"union bound"
(σ -sub-additivity)

$$\leq \sum_{h \in H_{BAD}} D^m \left(\{S|_k : L_S(h) = 0\} \right)$$

Let's fix some $h \in H_{BAD}$:

$$D^m \left(\{S|_k : L_S(h) = 0\} \right) = D^m \left(\{S|_k : \forall i \in [m] : h(x_i) = f(x_i)\} \right)$$

as all x_i 's are iid

$$= \prod_{i=1}^m D \left(\{x_i : h(x_i) = f(x_i)\} \right)$$

$$= \prod_{i=1}^m D \left(\{x : h(x) = f(x)\} \right)$$

$$= \prod_{i=1}^m (1 - L_{D,f}(h))$$

$$\leq \prod_{i=1}^m (1 - \varepsilon) \quad \begin{matrix} \text{as } h \text{ is a bad hypothesis} \\ h \in H_{BAD} \end{matrix}$$

$$= (1 - \varepsilon)^m$$

$$\leq e^{-\varepsilon m} \quad \text{(without proof)}$$

$$[m] = \{1, \dots, m\}$$

$$\Rightarrow \mathbb{D}^m \left(\{S|x : L_{D,f}(h_s) \geq \varepsilon\} \right) \leq \sum_{h \in H_{BAD}} \mathbb{D}^m \left(\{S|x : L_S(h) = 0\} \right) \leq e^{-\varepsilon m}$$

$$\leq |H_{BAD}| \cdot e^{-\varepsilon m}$$

$$\leq |H| \cdot e^{-\varepsilon m} \quad // \text{because } H_{BAD} \subseteq H$$

If we want $|H| \cdot e^{-\varepsilon m}$ to be less than some $\delta \in (0, 1)$, we can solve for m and get:

$$|H| \cdot e^{-\varepsilon m} < \delta$$

$$\Rightarrow m > \frac{1}{\varepsilon} \cdot \log \left(\frac{|H|}{\delta} \right)$$

error confidence

Corollary: Let $|H| < \infty$ and $\varepsilon, \delta \in (0, 1)$. Further, let m be an integer such that $m > \frac{1}{\varepsilon} \cdot \log \left(\frac{|H|}{\delta} \right)$. Then, for each $h_s: X \rightarrow Y$ and any distribution D over X (for which we have that with probability of at least $1 - \delta$ over the $S|x$ (of size m) every FERM_H hypothesis h_s satisfies

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

Interpretation: For sufficiently large m , FERM_H returns h_s that is PROBABLY APPROXIMATELY CORRECT (PAC).

This leads to:

Def. (PAC learnability): A hypothesis class H is **PAC learnable**, if there exists a function $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following properties: (I) for every $\varepsilon, \delta \in (0, 1)$ and (II) every distribution D over the domain X , and (III) every labeling function $f : X \rightarrow \{0, 1\}$, if (IV) **realizability** holds (with respect to H, D, f), then running A on $m \geq m_H(\varepsilon, \delta)$ iid instances from D (labeled by f) returns a hypothesis h such that with probability of at least $1 - \delta$ (over the choice of S)

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

Def. (Sample complexity): $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ is called the sample complexity function.

We will now move to a more general setting.

[1] We will first release the realizability assumption.

(in this setting, the best thing we can hope for are guarantees relative to the "best possible" hypothesis in the class H , i.e., relative to $\min_{h \in H} L_{D,f}(h)$).

Def. (Hoeffding inequality): Let X_1, \dots, X_m be iid random variables taking values in $[a_i, b_i]$ for $i \in [m]$. Then, given that

$$S_m = \sum_{i=1}^m X_i$$

it holds that

$$\mathbb{P}[S_m - \mathbb{E}[S_m] > \varepsilon] \leq e^{-\frac{-2\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}}$$

$$\mathbb{P}[S_m - \mathbb{E}[S_m] < -\varepsilon] \leq e^{-\frac{-2\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}}$$

$$\text{Also: } \mathbb{P}[|S_m - \mathbb{E}[S_m]| > \varepsilon] \leq 2 \cdot e^{-\frac{-2\varepsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}}$$

Another useful form of this inequality is

$$\mathbb{P}\left[\left|\frac{1}{m}S_m - \mu\right| > \varepsilon\right] \leq 2 \cdot e^{-\frac{2\varepsilon^2 m}{(b-a)^2}} \quad (*)$$

with $\mu = \mathbb{E}[X_i]$ and $\mathbb{P}[a \leq X_i \leq b] = 1$ for all $i \in [m]$.

As a consequence of $(*)$, we can say the following: fix $\varepsilon > 0$; then for a single $h: X \rightarrow Y$, we have

$$\mathbb{P}_{S \sim D^m} \left[\left| L_S(h) - L_{D,f}(h) \right| > \varepsilon \right] \leq 2 \cdot e^{-\frac{2\varepsilon^2 m}{(b-a)^2}} \quad (\text{as } b=1, a=0)$$

Corresponds
to S_m from above

If we would set $2e^{-2\varepsilon^2 m} = \delta$ and solve for ε , we would get

$$\varepsilon = \sqrt{\frac{\log(\frac{2}{\delta})}{2m}}$$

$\Rightarrow L_{D,f}(h) \leq L_S(h) + \sqrt{\frac{\log(\frac{2}{\delta})}{2m}}$; this holds with probability of at least $1-\delta$ over the choice of S !

Remark: This result only holds for a single h . However, we can easily get a bound that holds uniformly for all $h \in H$, with $|H| < \infty$:

$$\mathbb{P}_{S|x \sim D^m} \left[\exists h \in H : |L_S(h) - L_{D,f}(h)| \geq \varepsilon \right] \leq \underbrace{\sum_{h \in H} 2e^{-2\varepsilon^2 m}}_{\substack{\text{union} \\ \text{bound}}} \leq 2 \cdot |H| \cdot e^{-2\varepsilon^2 m}$$

(no realizability needed for flat result)

[2] Next, we release our requirement of a "true labeling function" f . We do this by letting D be a distribution over $X \times Y = \mathcal{Z}$. Hence, we need to adjust our definitions of the empirical error and the generalization error.

$$(1) \quad L_D(h) = \mathbb{P}_{(x,y) \sim D} [h(x) \neq y] = D(\{(x,y) \in X \times Y : h(x) \neq y\})$$

(generalization error)

$$(2) \quad L_S(h) = \frac{1}{m} \cdot \left| \{i \in [m] : h(x_i) \neq y_i\} \right|$$

(empirical error)

Def. (Agnostic PAC learnability): A hypothesis class H is **agnostic PAC learnable** if there exists $m_H: (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following properties: (I) for every $\epsilon, \delta \in (0,1)$ and (II) every distribution D over $X \times Y = \mathcal{Z}$, when running A on $m \geq m_H(\epsilon, \delta)$ iid instances from D , A returns a hypothesis h such that with probability of at least $1-\delta$ over the choice of S

$$L_D(h) \leq \min_{h' \in H} L_D(h') + \epsilon.$$

General "loss" functions

$$l: H \times (X \times Y) \rightarrow \mathbb{R}_+$$

Example: 0-1 loss

$$l^{0-1}(h, (x, y)) = \begin{cases} 1, & \text{if } h(x) \neq y \\ 0, & \text{else} \end{cases}$$

Square loss

$$l^{sp}(h, (x, y)) = (h(x) - y)^2 \quad \begin{matrix} \text{used in regression} \\ \text{problems} \end{matrix}$$

Now, even with more generality:

$$L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)], \quad l_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, t_i) \quad z = (x, y)$$

<u>1st exam date:</u>	July, 8, 2024	(10^{am})
<u>2nd exam date</u>	July, 16, 2024	(10^{am})

Remark: $\mathbb{E}_{x \sim D} [l^{0-1}(h, x)] = 0 \cdot \cancel{\mathbb{P}_x [h(x) = y]} + 1 \cdot \cancel{\mathbb{P}_x [h(x) \neq y]}$

$$= \mathbb{P}_x [h(x) \neq y]$$

Uniform convergence

Def. (ϵ -representative sample): A sample S is called ϵ -representative with respect to $\mathcal{Z} = X \times Y$, hypothesis class H , loss function l and distribution D , if

$$\forall h \in H: |L_S(h) - L_D(h)| \leq \epsilon.$$

Lemma: Assume that S is $\frac{\epsilon}{2}$ -representative wrt. \mathcal{Z}, H, l and D . Then, any hypothesis h_S returned by $\text{ERM}_H(S) \in \underset{h' \in H}{\operatorname{arg\,min}} L_S(h')$ satisfies

$$L_D(h_S) \leq \min_{h \in H} L_D(h) + \epsilon.$$

Proof: for any $h \in H$:

$$L_D(h_s) \leq L_S(h_s) + \frac{\epsilon}{2} \quad // by \text{ def. of } \epsilon\text{-rep.}$$

ERT hyp. \leftarrow

$$\begin{aligned} &\leq L_S(h) + \frac{\epsilon}{2} \quad // \text{as } h_s \text{ is an ERT hypothesis} \\ &\leq L_D(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} \quad // \text{by def. of } \epsilon\text{-rep.} \\ &= L_D(h) + \epsilon \end{aligned}$$

$$\Rightarrow L_D(h_s) \leq L_D(h) + \epsilon$$

As this inequality chain holds for any $h \in H$, we can conclude that

$$L_D(h_s) \leq \min_{h \in H} L_D(h) + \epsilon$$

◻

Def. (uniform convergence): A hyp. class H has the uniform convergence property (UC) with respect to ℓ and loss function l , if there exists $m_H^{uc}: (0, 1)^2 \rightarrow \mathbb{N}$, such that for any $\epsilon, \delta \in (0, 1)$ and every distribution D over $Z = X \times Y$, if S is an iid sample of size $m \geq m_H^{uc}(\epsilon, \delta)$ from D , then with probability of at least $1 - \delta$ (over the choice of S), S is ϵ -representative.

Corollary: If H has the UC property with m_H^{uc} , then H is agnostic PAC learnable with

$$m_H(\epsilon, \delta) \leq m_H^{uc}\left(\frac{\epsilon}{2}, \delta\right)$$

We also know that ERM is a successful agnostic PAC learner.

Claim Finite hypothesis classes (H/∞) are agnostic PAC learnable.

Fix $\varepsilon, \delta \in (0, 1)$. We want to show that

$$\mathcal{D}^m \left(\{S : \forall h \in H : |L_S(h) - L_D(h)| \leq \varepsilon\} \right) \geq 1 - \delta \quad (\times)$$

Rerphrasing (\times) gives

$$\mathcal{D}^m \left(\underbrace{\{S : \exists h \in H : |L_S(h) - L_D(h)| > \varepsilon\}}_{\bigcup_{h \in H} \{S : |L_S(h) - L_D(h)| > \varepsilon\}} \right) < \delta$$

By the union bound

$$\mathcal{D}^m \left(\bigcup_{h \in H} \{S : |L_S(h) - L_D(h)| > \varepsilon\} \right) \stackrel{(\text{blue})}{\leq} \sum_{h \in H} \mathcal{D}^m \left(\{S : |L_S(h) - L_D(h)| > \varepsilon\} \right)$$

We now fix $\ell := \ell^{0-1}$ (0/1 loss). We know that $L_D(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$ and $L_S(h) = \frac{1}{m} \cdot \sum_{i=1}^m \ell(h, z_i)$. Hence, we can

apply Höffding's inequality:

$$\mathcal{D}^m \left(\{S : |L_S(h) - L_D(h)| > \varepsilon\} \right) \leq 2 \cdot e^{-2\varepsilon^2 m}$$

$$\Rightarrow \mathcal{D}^m \left(\{S : \exists h \in H : |L_S(h) - L_D(h)| > \varepsilon\} \right) \leq 2 \cdot |H| \cdot e^{-2\varepsilon^2 m} \quad (\dagger\dagger)$$

(fix 2023 notes)!

What remains is to let $2|H|e^{-2\varepsilon_m^2}$ be smaller than $\delta \in (0, 1)$ and solve for m :

$$2|H|e^{-2\varepsilon_m^2} < \delta$$
$$\Rightarrow m > \frac{1}{2\varepsilon^2} \cdot \log\left(\frac{2|H|}{\delta}\right)$$

This yields the following sample complexity function m_H^{UC} for uniform convergence:

$$m_H^{UC}(\varepsilon, \delta) \leq \left\lceil \log\left(\frac{2|H|}{\delta}\right) \cdot \frac{1}{2\varepsilon^2} \right\rceil$$

→ By our previous corollary on the relation between UC and diagnostic PAC learnability, we get

$$m_H(\varepsilon, \delta) \leq m_H^{UC}\left(\frac{\varepsilon}{2}, \delta\right) \leq \left\lceil \frac{2}{\varepsilon^2} \cdot \log\left(\frac{2|H|}{\delta}\right) \right\rceil$$

This establishes our claim.

Question: Is there a "universal" learner?

By **universal** we mean a learner A without any prior knowledge of a learning task, but can be challenged by any task and still achieve low $L_D(A(s))$.

defined by D

Theorem (No-Free-Lunch (NFL)): Let A be a learning algorithm for the task of binary (0/1) classification with respect to the 0-1 loss over a domain X . Also, let m be any number smaller than $|X|/2$, representing the training set size. Then, there exists a distribution D over $X \times \{0,1\}$ such that

1. $\exists f : X \rightarrow \{0,1\}$ with $L_D(f) = 0$ / 0 generalization error

2. with probability of at least $\frac{1}{7}$ over the choice of $S \sim D^m$, we have

$$L_D(A(s)) = \frac{1}{8}$$

e.g. FRk_H with $H = \{f\}$

Here, 1. means that the task can be successfully learned by another learner; 2. means that A fails on that task.

Proof: Let $C \subset X$ of size 2^m , i.e., $|C| = 2^m$. The number of possible labelings of C is $2^{2^m} = T$. We denote the functions that realize these T labelings as

$$f_1, \dots, f_T$$

For each f_i , we define

$$\mathcal{D}_i(\{(x_i, y)\}) = \begin{cases} \frac{1}{|C|}, & \text{if } y = f_i(x) \\ 0, & \text{else} \end{cases}$$

defines the learning task

Consequence of this construction: we have $\forall i : L_{\mathcal{D}_i}(f_i) = 0$!

We will show now that for every learning algorithm A that receives samples from $X \times \{0, 1\}$, there exists a function $f: X \rightarrow \{0, 1\}$ and a distribution \mathcal{D} over $X \times \{0, 1\}$, such that

$$L_{\mathcal{D}}(f) = 0 \quad \text{and} \quad \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] \geq \frac{1}{4} \quad (\times)$$

In particular, we will show that for every learning algorithm A receiving m samples from $C \times \{0, 1\}$ and returning a function $A(S): X \rightarrow \{0, 1\}$, we have

$$\max_{i \in [T]} \mathbb{E}_{S \sim \mathcal{D}_i^m} [L_{\mathcal{D}_i}(A(S))] \geq \frac{1}{4} \quad (\times)$$

$$[T] = \{1, \dots, T\}$$

Remark: (x) suffices for $L_D(A(s)) \geq \frac{1}{8}$ with prob. of at least $\frac{1}{7}$ over the choice of $s \sim D^m$ (see PS).

lets continue with (xx) : we have $(2m)^m = k$ possible training sequences of size m . we are going to call them

$$S_1, \dots, S_k$$

lets also call S_j^i the sequence S_j labeled by f_i , i.e.,

$$S_j^i = ((x_1, f_i(x_1)), \dots, (x_m, f_i(x_m)))$$

so, in case the distribution is D_i , then S_1^i, \dots, S_k^i are all possible training sequences (of size m) that A can receive. Also, by construction (of the D_i 's), all these training sequences have equal probability of being drawn.

$$\Rightarrow \mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(s))] = \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(S_j^i)) \quad (\text{xxx})$$

$$\begin{aligned} \text{so, } \max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(s))] &= \max_{i \in [T]} \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i)) // \text{by (xxx)} \\ \text{as } \max &\geq \text{avg.} \end{aligned}$$

$$\geq \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k L_{D_i}(A(S_j^i))$$

$$= \frac{1}{k} \cdot \sum_{j=1}^k \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j^i))$$

$$\geq \min_{j \in [k]} \frac{1}{T} \cdot \sum_{i=1}^T L_{D_i}(A(S_j^i))$$

Let's fix some $j \in [k]$: As S_j is of size m , but C is of size $2m$, there are instances from C which we have not seen. Let's call them

$$v_1, \dots, v_p$$

We know that $p \geq m$.

Now, for every $h: C \rightarrow \{0,1\}$ and every i , it holds that

$$L_{D_i}(h) = \frac{1}{2m} \cdot \sum_{x \in C} \mathbb{1}_{h(x) \neq f_i(x)}$$

$$\geq \frac{1}{2m} \sum_{h=1}^p \mathbb{1}_{h(v_h) \neq f_i(v_h)} \quad \begin{matrix} \text{as the } v_h \text{'s are} \\ \text{less than } 2m \end{matrix}$$

$$\geq \frac{1}{2p} \cdot \sum_{h=1}^p \mathbb{1}_{h(v_h) \neq f_i(v_h)} \quad \text{as } p \geq m$$

Combining the results so far:

$$\frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j^i)) \geq \frac{1}{T} \sum_{i=1}^T \frac{1}{2p} \sum_{h=1}^p \underbrace{\mathbb{1}_{A(S_j^i)(v_h) \neq f_i(v_h)}}_h$$

$$\Rightarrow \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j^i)) \geq \frac{1}{2} \cdot \frac{1}{\rho} \sum_{r=1}^{\rho} \frac{1}{T} \cdot \sum_{i=1}^T \frac{1}{A(S_j^i)(v_r)} \neq f_i(v_r)$$

$$\geq \frac{1}{2} \cdot \min_{r \in [\rho]} \frac{1}{T} \cdot \sum_{i=1}^T \frac{1}{A(S_j^i)(v_r)} = f_i(v_r)$$

(opinion, as $\min \geq \min$)

Let's try fixing $r \in [\rho]$: we have T labeling functions f_1, \dots, f_T (we have 2^{2n} of them)

Let's say, we would have $m=1$, so $|C|=2$ (as an example)

		unseen
		v_1
		0 1
f_1		0
f_2		0
f_3		1
f_4		1

partition into 2 pairs
 $(f_1, f_2), (f_3, f_4)$
 pair 1 pair 2

where (f_1, f_2) only differ on v_1
and (f_3, f_4) only differ on v_1 .

In general, we can always partition into $\frac{T}{2}$ disjoint pairs of labeling functions. For every pair $(f_i, f_{i'})$, we have that for every $c \in C$, $f_i(c) \neq f_{i'}(c)$ if and only if $c = v_r$. Further, for such a pair, we also have $S_j^i = S_j^{i'}$.

$$\Rightarrow \sum_{i=1}^T \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)} + \sum_{i=1}^T \mathbb{1}_{A(S_j^{i'})}(v_r) \neq f_{i'}(v_r) = 1$$

Overall, we now have:

$\frac{1}{T} \cdot \sum_{i=1}^T \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)}$ and we re-write this as
 a sum over all $T/2$ disjoint pairs $(f_i, f_{i'})$, knowing that
 $(\times \times \times)$ holds; we get

$$\frac{1}{T} \cdot \frac{T}{2} = \frac{1}{2}$$

$$\Rightarrow \frac{1}{T} \sum_{i=1}^T L_{D_i}(A(S_j^i)) \geq \frac{1}{4}$$

Remark: Our earlier steps of fixing j and κ are irrelevant as we have a constant as a lower bound now.

Combining all partial results, yields:

$$\max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(S))] \geq \frac{1}{4}$$

Corollary: let X be an infinitely large domain and let H be the set of all functions from $X \rightarrow \{0, 1\}$. Then, H is not PAC learnable.

Vapnik - Chervonenkis Dimension (VC Dimension)

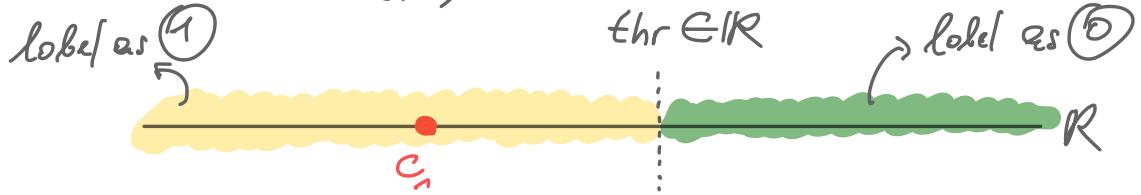
(Knowing that in the NFL theorem what mattered most was the behavior of labeling functions on a dataset, it seems intuitive to tie learnability to this behavior.)

Def. Let H be a class of functions from $X \rightarrow \{0, 1\}$ and let $C \subset X$, $C = \{c_1, \dots, c_m\}$. We define

$$H_C = \{(h(c_1), \dots, h(c_m)) : h \in H\}$$

as the restriction of H to C .

Example : $C = \{c_1\}$, $c_1 \in \mathbb{R}$ and lets look at the class of thresholds on the real line.



lets call this hypo. class H^{thr} . Then

$$H_C^{thr} = \{(1), (0)\} \Rightarrow |H_C^{thr}| = 2^1 = 2$$

lets take $C = \{c_1, c_2\}$



$$H_C^{thr} = \{(1, 1), (0, 0), (1, 0)\} \Rightarrow |H_C^{thr}| = 3$$

(all possible labelings would be $2^2 = 4$)

Def. H **shatters** a finite set C of size m ($|C|=m$), if

$$|H_C| = 2^m$$

Example (continued): The class of thresholds on \mathbb{R} ($=X$) is
PAC-learnable.

$$H^{thr} = \{h_\alpha : \alpha \in \mathbb{R}\},$$

$$h_\alpha : \mathbb{R} \rightarrow \{0, 1\}$$

$$x \mapsto h_\alpha(x) = \mathbb{1}_{x < \alpha} = \begin{cases} 1, & \text{if } x < \alpha \\ 0, & \text{else} \end{cases}$$

Claim: H^{thr} is PAC learnable with

$$m_{H^{thr}}(\varepsilon, \delta) \leq \lceil \log\left(\frac{2}{\delta}\right) \cdot \frac{1}{\varepsilon} \rceil$$

Proof: We assume realizability, hence $\exists h^* \in H^{thr}$ such that

$$L_{D,f}(h^*) = 0.$$

we let α^* be the corresponding threshold.

$$\begin{array}{c} \alpha^* \\ | \\ \mathbb{R} \end{array}$$



let q_0 be such that $D(\{x \in \mathbb{R} : x \in (q_0, q^*)\}) = \varepsilon$
 let q_1 be - - - $D(\{x \in \mathbb{R} : x \in (q^*, q_1)\}) = \varepsilon$

In case $D(\{x \in \mathbb{R} : x \in (q_0, q^*)\}) < \varepsilon$, set $q_0 = -\infty$ } special cases
 In case $D(\{x \in \mathbb{R} : x \in (q^*, q_1)\}) < \varepsilon$, set $q_1 = +\infty$ }

We have $S = ((x_1, y_1), \dots, (x_m, y_m))$. An ERh algorithm would be, for instance, to first pick b_0 and b_1 as follows:

$$b_0 = \max \{x : (x, 1) \in S\}$$

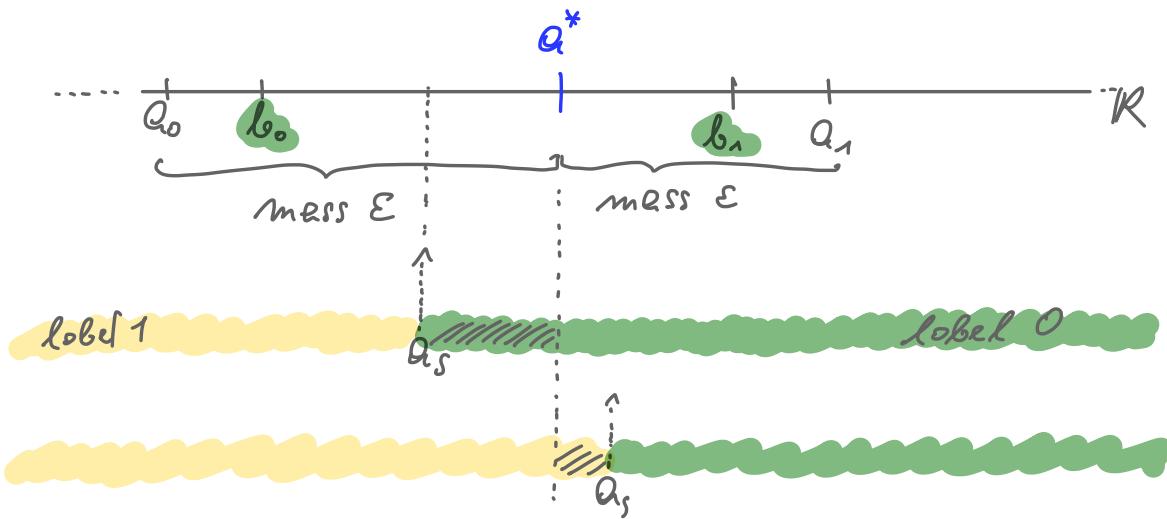
$$b_1 = \min \{x : (x, 0) \in S\}$$

Then, ERh can choose any threshold within (b_0, b_1) . We will call the chosen threshold q_s and the corresponding hypothesis h_s .

By these choices, we get 0 empirical error.

For h_s (with threshold q_s) to have $L_{D,f}(h_s) \leq \varepsilon$, it suffices that

- 1) $b_1 \leq q_s$ AND
- 2) $b_0 \geq q_s$



In other words,

$$\begin{aligned} \mathbb{P}[L_{D,f}(h_s) > \varepsilon] &\leq \mathbb{P}\left[\left(b_0 < q_0\right) \stackrel{(CV)}{\text{OR}} \left(b_1 > q_1\right)\right] \\ &\stackrel{\text{union bound}}{\leq} \mathbb{P}[b_0 < q_0] + \mathbb{P}[b_1 > q_1] \end{aligned}$$

Hence, to upper bound $\mathbb{P}[L_{D,f}(h_s) > \varepsilon]$, we need to check when $b_0 < q_0$ (or $b_1 > q_1$) actually happens. Answer: when there is no data point x in S that is labeled 1 s.t. $x \in (q_0, q^*)$.

We know that $D(\{x \in \mathbb{R}: x \in (q_0, q^*)\}) = \varepsilon$, hence not seeing a data point in (q_0, q^*) has probability of $1 - \varepsilon$. Consequently, not seeing a data point in m i.i.d samples from D has probability of $(1 - \varepsilon)^m$.

Remark: $(1-\varepsilon)^m \leq e^{-\varepsilon m}$

Overall, $\underset{S \sim D^m}{\mathbb{P}} [h_0 < q_0] = (1-\varepsilon)^m \leq e^{-\varepsilon m}$

$$\underset{S \sim D^m}{\mathbb{P}} [h_0 > q_1] = (1-\varepsilon)^m \leq e^{-\varepsilon m}$$

Combined, we get

$$\mathbb{P}[L_{D, f}(h_S) > \varepsilon] \leq \underbrace{\lambda \cdot e^{-\varepsilon m}}$$

FRh hyps.

let this be $< \delta$ and solve for
 $m \Rightarrow$

$$m > \log\left(\frac{2}{\delta}\right) \cdot \frac{1}{\varepsilon}$$

We have seen that "finiteness" of $H(H/\infty)$ is sufficient, but not necessary (for PAC learnability).

Def. (VC-Dimension): The VC dimension of H (i.e., a class of functions from $X \rightarrow \{0, 1\}$), written as $VC(H)$,

is the maximal size of a set $C \subset X$
that is shattered by H .

Theorem : Let H be a class of functions from $X \rightarrow \{0,1\}$.
 if H has infinite VC dimension, then H is not PAC learnable.

(follows immediately from NFL)

Def. (growth function):

Let H be a class of functions from $X \rightarrow \{0,1\}$. The growth function of H , $\gamma_H: \mathbb{N}_0 \rightarrow \mathbb{N}_0$, is defined as

$$\gamma_H(m) = \max_{C \subset X, |C|=m} |H_C| \quad \left(\begin{array}{l} \text{with} \\ \gamma_H(0) \stackrel{\text{def.}}{=} 1 \end{array} \right)$$

Lemma (Sauer, Shelah, Peretz; "Sauer's lemma") : Let H be a hyp. class of functions from $X \rightarrow \{0,1\}$ with $VC(H)=d$.

Then

$$\gamma_H(m) = 2^m \text{ if } m \leq d, \text{ but}$$

$$\gamma_H(m) = \left(\frac{e^m}{d}\right)^d \text{ if } m > d$$

(without proof)

VC-Dimension of finite H :

It's obvious that, if we take any set C (of size m),

$$|H_C| \leq |H|$$

So, if $|H| < 2^m$, then H cannot shatter C (of size m).

This implies

$$VC(H) \leq \log_2(|H|)$$

Theorem: Let H be a hyp. class of functions from $X \rightarrow \{0,1\}$ and $\ell: H \times X \times Y \rightarrow [0,c]$, $c > 0$ be a loss function. For any distribution D over $X \times Y$ and $\delta \in (0,1)$, we have with probability of at least $1-\delta$ over the choice of $S \sim D^m$

$$\forall h \in H: \left| L_D(h) - L_S(h) \right| \leq c \cdot \sqrt{\frac{8 \cdot \log(\gamma_H(2m) \cdot \frac{4}{\delta})}{m}}$$

growth function, evaluated at $2m$.

(without proof)

Fundamental theorem of statistical learning

Let H be a hyp. class of functions from $X \rightarrow \{0, 1\}$, and let ℓ be the 0-1 loss. Then, the following statements are equivalent:

1. H has the uniform convergence property.
2. Any ERM algorithm is a successful agnostic PAC learner for H .
3. H is agnostic PAC learnable
4. Any ERM algorithm is a successful PAC learner for H .
5. H is PAC learnable
6. H has finite VC dimension ($VC(H) < \infty$)

Example: Hyp. class of LINEAR PREDICTORS

We define

$$\{h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R}\} =: L_d, \text{ with}$$

$$h_{w,b}(x) = \langle w, x \rangle + b, \quad x \in \mathbb{R}^d$$

as the class of **affine functions**. From L_d , we get different predictors via composition with some

$$\phi: \mathbb{R} \rightarrow Y$$

We can also use

$$w' = (b, w_1, \dots, w_d)^T \in \mathbb{R}^{d+1}$$

$$x' = (1, x_1, \dots, x_d)^T \in \mathbb{R}^{d+1}$$

to write

$$h_{w,b}(x) = \langle w', x' \rangle$$

If $\phi: \mathbb{R} \rightarrow Y$ is $\phi = \text{sign}$, then we get halfspace predictors

$$HS_d = \phi \circ L_d - \{x \mapsto \text{sign}(h_{w,b}(x)) : h_{w,b} \in L_d\}$$

Question: What is the VC dimension of

$$\mathcal{F} = \left\{ x \mapsto \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d \right\},$$

that is $\text{VC}(\mathcal{F}) = ?$

Claim: $\text{VC}(\mathcal{F}) = d$

Part 1 : Show that $\text{VC}(\mathcal{F}) \geq d$

$$X = \{c_1, \dots, c_d\} \text{ with } c_i \in \mathbb{R}^d \text{ for all } i \in [d]$$

with

$$c_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad c_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad \dots, \quad c_d = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}$$

and

$$w = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{pmatrix} \rightarrow \begin{array}{l} \text{label of } c_1 \quad (\in \{\pm 1\}) \\ \text{label of } c_d \end{array}$$

Remark: c_{ij} refers to the j-th coordinate of c_i

$$\Rightarrow \langle w, c_i \rangle = \sum_{j=1}^d w_j \cdot c_{ij} = 1 \cdot w_i = y_i \Rightarrow \text{VC}(\mathcal{F}) \geq d$$

Part 2

We need to show $VC(\mathcal{F}) < d+1$. We will show this via contradiction. \hookrightarrow no set of size $d+1$ is shattered by \mathcal{F} .

! Assume that $d+1$ points

$$X = \{c_1, \dots, c_{d+1}\}$$

! are shattered by \mathcal{F} .

This means that $\mathcal{F}w_1, \dots, w_{2^{d+1}}$ weight vectors that will yield all (2^{d+1}) different labelings of our $d+1$ points.

Let's write down all possible inner products:

$$M = \begin{pmatrix} w_1^T c_1 & w_2^T c_1 & \cdots & \cdots & \cdots & w_{2^{d+1}}^T c_1 \\ w_1^T c_2 & \vdots & & & & \vdots \\ w_1^T c_3 & \vdots & & & & \vdots \\ \vdots & & & & & \vdots \\ w_1^T c_{d+1} & w_2^T c_{d+1} & \cdots & \cdots & \cdots & w_{2^{d+1}}^T c_{d+1} \end{pmatrix}$$

$\xrightarrow{2^{d+1} \text{ columns}}$

d+1 rows

We can also write this matrix as

$$Xw = H$$
$$X = \begin{pmatrix} c_1 \\ \vdots \\ c_{d+1} \end{pmatrix} \quad \text{d+1 rows of columns}$$
$$w = \begin{pmatrix} | & | & & | \\ w_1 & w_2 & \dots & w_{2^{d+1}} \\ | & | & & | \end{pmatrix} \quad \begin{matrix} d \text{ rows} \\ 2^{d+1} \text{ columns} \end{matrix}$$

If we would take $\text{sign}(Xw)$ we would get all possible labelings. We know that

$$\text{rank}(H) \leq \min \{ \text{rank}(X), \text{rank}(w) \} = d$$

06/18/29

LD

$$H = \begin{pmatrix} \omega_1^T c_1 & & & \omega_2^T c_{d+1} \\ \omega_1^T c_2 & \cdots & & \vdots \\ \vdots & & & \vdots \\ \omega_1^T c_{d+1} & & & \omega_2^T c_{d+1} \end{pmatrix}$$

\downarrow

$x_{w_1} \in \mathbb{R}^{d+1}$

Remark: Any v_1, \dots, v_n (set of vectors) are lin. independent iff

$$\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n = \textcircled{0}$$

only when $\theta_i : \alpha_i = 0$

In our case (looking at the rows of H):

$$\begin{aligned} Q_1 \cdot [\omega_1^T c_1 \quad \cdots \quad \omega_2^T c_{d+1}] &+ \\ Q_2 \cdot [\omega_1^T c_2 \quad \cdots \quad \omega_2^T c_{d+1}] &+ \\ &\vdots \\ Q_{d+1} \cdot [\omega_1^T c_{d+1} \quad \cdots \quad \omega_2^T c_{d+1}] &- \\ [\alpha^T X_{w_1} \quad \alpha^T X_{w_2} \quad \dots \quad \alpha^T X_{w_{d+1}}] & \end{aligned}$$

by our assumption

As c_1, \dots, c_{d+1} are shattered by F , there $\exists k \in \{1, \dots, 2^{d+1}\}$ such that

$$\text{sign}(\alpha) = \text{sign}(x_{w_k})$$

This means that

$$a^T X w_k$$

is a sum of positive numbers $\rightarrow +\infty$!

$\Rightarrow d+1$ rows of H are lin. independent $\Rightarrow \text{rank}(H) = d+1$

\Rightarrow Contradiction (with $\text{rank}(H) \leq d$).

We conclude that $\text{VC}(T) = d$. (as $\text{VC}(T) \geq d$ and $\text{VC}(T) < d+1$)

RADEMACHER complexity

Unif. con. (ε -rep.)

$$\forall h: \underset{z \in H}{\sup} |L_D(h) - L_S(h)| \leq \varepsilon$$

We can also say

$$\underset{h \in H}{\sup} |L_D(h) - L_S(h)| \leq \varepsilon$$

$$\left. \begin{array}{l} z = x \vee y, \text{ class } H \\ \text{loss } \ell: H \times z \rightarrow \mathbb{R}^+$$

$$\text{Def: } \mathcal{F} = \ell \circ H = \left\{ z \mapsto \ell(h, z) : h \in H \right\}_{(x, y)}$$

$$\text{and for } f \in \mathcal{F}: \boxed{L_D(f) = \mathbb{E}_z [f(z)]} \quad L_S(f) = \frac{1}{m} \cdot \sum_{i=1}^m f(z_i)$$

$$\text{Def: } \text{Rep}_D(\mathcal{F}, S) = \sup_{f \in \mathcal{F}} (L_D(f) - L_S(f)) \quad (\times)$$

Take $S = S_1 \cup S_2$ with $|S_1| = |S_2|$ and estimate (\times) by

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left(\underbrace{L_{S_1}(f)}_{\substack{z \in S_1}} - \underbrace{L_{S_2}(f)}_{\substack{z \in S_2}} \right) \\ & + 1 \cdot \frac{2}{m} \cdot \underbrace{\sum_{z \in S_1} f(z)}_{\substack{z \in S_1}} + (-1) \cdot \frac{2}{m} \cdot \underbrace{\sum_{z \in S_2} f(z)}_{\substack{z \in S_2}} \\ & = \frac{2}{m} \cdot \underbrace{\sum_{z \in S_1} (\textcolor{red}{+1}) \cdot f(z)}_{\substack{z \in S_1}} + \frac{2}{m} \cdot \underbrace{\sum_{z \in S_2} (-1) \cdot f(z)}_{\substack{z \in S_2}} \\ & \text{collect into a vector } \sigma \in \{-1\}^m \end{aligned}$$

Set $S_1 = \{z_i : \sigma_i = +1\}$, $S_2 = \{z_i : \sigma_i = -1\}$. We rewrite

$$\sup_{f \in \mathcal{F}} \frac{2}{m} \cdot \sum_{i=1}^m \sigma_i f(z_i)$$

Def.

$$\mathcal{F} \circ S = \{(f(z_1), \dots, f(z_m)) : f \in \mathcal{F}\}$$

$$RAD(\mathcal{F} \circ S) = \frac{1}{m} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(z_i) \right]$$

$$\sigma \sim \text{Uniform}(\{\pm 1\}^m)$$

$\langle \sigma_1, (f(z_1), \dots, f(z_m)) \rangle$

RADEMACHER
COMPLEXITY

$$\mathbb{P}[\sigma_i = +1] = \mathbb{P}[\sigma_i = -1] = \frac{1}{2}$$

Lemma

$$\mathbb{E}_{S \sim D^m} [Rep_D(\mathcal{F}, S)] \leq 2 \cdot \mathbb{E}_{S \sim D^m} [RAD(\mathcal{F} \circ S)]$$

$$\sup_{f \in \mathcal{F}} (L_D(f) - L_S(f))$$

Theorem (informal)

Assume for all $t \in \mathcal{T}$ and $h \in H$: $|l(h, t)| \leq c$.

Then

(1) with prob. of $1-\delta$ over the choice of S , for all $h \in H$:

$$L_D(h) - L_S(h) \leq 2 \cdot \mathbb{E}_{S \sim D^m} [\text{RAD}(l \circ H \circ S)] + c \cdot \sqrt{\frac{2 \log(\frac{2}{\delta})}{m}}$$

(2) $-n - \leq 2 \cdot \text{RAD}(l \circ H \circ S) + 4c \cdot \sqrt{\frac{2 \log(\frac{4}{\delta})}{m}}$