

STATISTICAL LEARNING THEORY

Summer Term 2025

Book: Shai Shalev-Schwartz

Understanding Machine Learning

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

OTHER COURSE NAMES:

- 1) Advanced machine learning (for AI students)
- 2) Machine learning (for Data Science students)

MOTIVATION

Example:

	$\in \mathbb{R}$ Weight (in g)	$\in \mathbb{R}$ Color ($\in [0, 1]$)	Tasty?
Popeye 1	800	0.1	0
:	:	:	:
:	:	:	:
:	:	:	:
Popeye N	1200	0.7	1

0... non-tasty
1... tasty

lets call the information in this table our training data.
Based on that, we try to find

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

i.e., a function that will take a two-dimensional vector as input (weight & color) and output a prediction of whether a popeye is tasty (1) or not (0).

We call such a function a hypothesis.

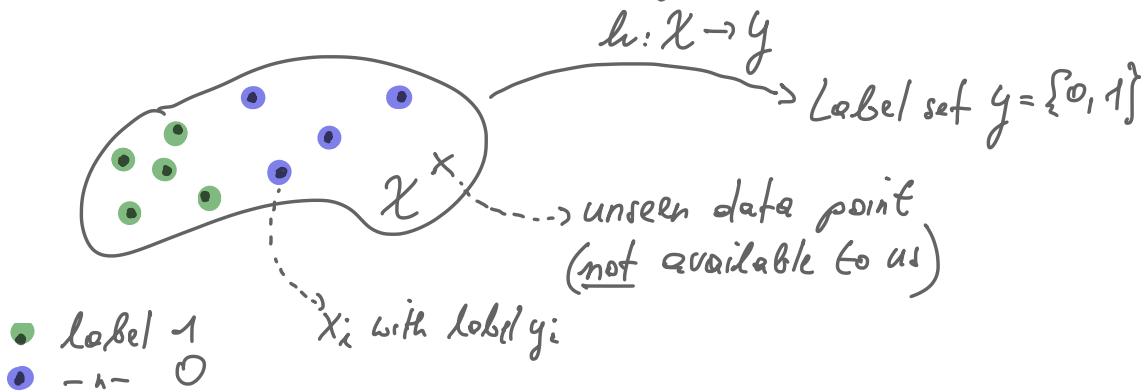
More formally, the data that is available to us comes in the form

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

with $x_i \in \mathbb{R}^2$ and $y_i \in \{0, 1\}$.

Domain $X = \mathbb{R}^2$ Label set $Y = \{0, 1\}$

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

independently and identically (iid)

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

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

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

called the true labeling function. This means

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

What do we care about?

We care about

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

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

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

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

We are interested in:

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

"Generalization error"

The empirical version of that is

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

"Empirical error"
(Empirical risk)

Notation:

$$[m] = \{1, \dots, m\}, S|_{\mathcal{X}} = (x_1, \dots, x_m)$$

$$\underline{\text{Claim: }} \mathbb{E} [L_S(h)] = L_{D,f}(h)$$

$\underbrace{S|x \sim D^m}_{(x_1, \dots, x_m)}$ $\frac{1}{m} \cdot |\{i \in [n] : h(x_i) \neq f(x_i)\}|$

$$\begin{aligned} \mathbb{E} [L_S(h)] &= \mathbb{E} \left[\underbrace{\frac{1}{m} \cdot \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq f(x_i)}}_{S|x \sim D^m} \right] \quad // \text{by def.} \\ &= \frac{1}{m} \cdot \sum_{i=1}^m \mathbb{E} [\mathbb{1}_{h(x_i) \neq f(x_i)}] \quad // \text{by linearity of } \mathbb{E}. \\ &= \frac{1}{m} \cdot \sum_{i=1}^m \mathbb{E} [\mathbb{1}_{h(x) \neq f(x)}] \quad // \text{as all } x_i's \text{ are drawn i.i.d from } D \\ &= \frac{1}{m} \cdot \sum_{i=1}^m \mathbb{P}[h(x) \neq f(x)] \\ &= \frac{1}{m} \cdot \mathbb{P}[h(x) \neq f(x)] = \underbrace{L_{D,f}(h)}_{\text{generalization error}} \end{aligned}$$

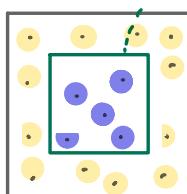
Our first learning paradigm

Empirical Risk Minimization (ERM). As a learner only has access to the training data (S), it is "natural" to try to select h (our hypothesis) such that the empirical risk (emp. error) is minimized. We call such a h an **empirical risk minimizer** (an ERM hypothesis).

Example (of a problematic case):



Domain X



→ true labeling function f

Say the distribution on X is uniform.

- label 1
- label 0

Also, assume that the area of the domain (square) is 2 and the area of \square is 1.

Now, say we have an ERM algorithm that returns h_S such that

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

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

$$L_S(h_S) = 0!$$

But, on unseen instances from D (which is uniform on X), h_S is only correct 50% of the time \rightarrow

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

That's what we call overfitting!

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

$$\text{ERM}_H(S) \in \arg \min_{h \in H} L_S(h)$$

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 against h^* (which has $L_{D,f}(h^*) = 0$ and, obviously, $L_S(h^*) = 0$).

We know that $L_{D,f}(h_S) > \varepsilon$, $\varepsilon \in (0, 1)$, can only happen if our learner selects a hypothesis h_S with $L_S(h_S) = 0$, BUT $L_{D,f}(h_S) > \varepsilon$.

We define $H_{\text{BAD}} = \{h \in H : L_{D,f}(h) > \varepsilon\}$ set of bad hypothesis!

Also, we define

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

Observation:

$$\left\{ S|x : L_{D,f}(h_s) > \varepsilon \right\} \subseteq \left\{ S|x : \exists h \in H_{BAD}, L_s(h) = 0 \right\} = M$$

↑
ERM hypothesis
(Empirical risk minimizer)

Since $\left\{ S|x : \exists h \in H_{BAD}, L_s(h) = 0 \right\} = \bigcup_{h \in H_{BAD}} \left\{ S|x : L_s(h) = 0 \right\}$

we have

$$D^m \left(\left\{ S|x : \exists h \in H_{BAD}, L_s(h) = 0 \right\} \right) = D^m \left(\bigcup_{h \in H_{BAD}} \left\{ S|x : L_s(h) = 0 \right\} \right)$$

(by σ -sub-additivity
"union" bound) $\leq \sum_{h \in H_{BAD}} D^m \left(\left\{ S|x : L_s(h) = 0 \right\} \right)$

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

True labeling
function
↓

$$D^m \left(\left\{ S|x : L_s(h) = 0 \right\} \right) = D^m \left(\left\{ S|x : \forall i \in \{1, \dots, m\} : h(x_i) = f(x_i) \right\} \right)$$

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

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

(By definition
of $L_{D,f}$) $= \prod_{i=1}^m \left(1 - \underbrace{\frac{L_{D,f}(h)}{\varepsilon}}_{> \varepsilon} \right)$ (remember that
 $h \in H_{BAD}$)

$$\begin{aligned}
 &\leq \prod_{i=1}^m (1-\varepsilon) \\
 &= (1-\varepsilon)^m \\
 &\leq e^{-\varepsilon m} \quad (\text{without proof})
 \end{aligned}$$

Overall, we have

$$\begin{aligned}
 D^m \left(\{S|_k : L_{D, \delta}(h_s) > \varepsilon \} \right) &= \sum_{h \in H_{BAD}} D^m \left(\{S|_k : L_S(h) = 0\} \right) \\
 &\leq \sum_{h \in H_{BAD}} e^{-\varepsilon m} \\
 &= |H_{BAD}| \cdot e^{-\varepsilon m} \\
 &\leq |H| \cdot e^{-\varepsilon m} \quad \left(\text{because } H_{BAD} \subseteq H \right)
 \end{aligned}$$

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

$$\begin{aligned}
 |H| \cdot e^{-\varepsilon m} &< \delta \\
 \Rightarrow m &> \frac{1}{\varepsilon} \cdot \log \left(\frac{|H|}{\delta} \right)
 \end{aligned}$$

error

we call $(1-\delta)$ confidence

In words, the probability of the generalization error of an ERk hypothesis being larger or equal to $\varepsilon \in (0, 1)$ is upper bounded by $\|f\| \cdot e^{-mc}$ (to be understood as the probability over a random draw of a dataset S of size m , iid from D).

Interlude (MARKOV inequality)

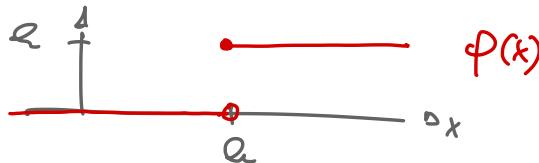
We have a probability space (S, \mathcal{F}, P) and a non-negative random variable X ($X \geq 0$).

Claim: for all $a > 0$, we have

$$\boxed{P[X \geq a] \leq \frac{\mathbb{E}[X]}{a}}$$

Proof: Define $\varphi: S \rightarrow \mathbb{R}$

$$\varphi(x) = \begin{cases} a, & \text{if } x \geq a \\ 0, & \text{else } (x < a) \end{cases}$$



We have $0 \leq \varphi(X(x)) \leq X(x)$ (for all $x \geq 0$)

$$\Rightarrow \int X dD \geq \underbrace{\int \varphi(X(\omega)) dD}_{\alpha \cdot D(\{\omega \in S : X(\omega) \geq \alpha\})} \quad (\text{by monotonicity})$$

$$P[X \geq \alpha]$$

Since $\alpha > 0$:

$$\frac{1}{\alpha} \cdot \underbrace{\int X dD}_{E[X]} \geq P[X \geq \alpha]$$

$E[X]$ by def.

$$\Rightarrow P[X \geq \alpha] \leq \frac{E[X]}{\alpha} \quad \square$$

Corollary: Let $|H| < \infty$ and $\epsilon, \delta \in (0, 1)$. Further, let m be an integer such that $m > \frac{1}{\epsilon} \cdot \log\left(\frac{|H|}{\delta}\right)$. Then, for any labeling function $f: X \rightarrow \{0, 1\}$ and any distribution D (for which realizability holds), we have that with probability of at least $1 - \delta$ over the choice of S/X of size m , every ERM hypothesis h_S satisfies

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

Interpretation: For sufficiently large m , ERM returns a hypothesis h_S that is $\underbrace{1 - \delta}_{\epsilon}$ PROBABLY APPROXIMATELY CORRECT (PAC)

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 $\epsilon, \delta \in (0, 1)$ and (II) every distribution D over the domain \mathcal{X} , and (III) every labeling function $f: \mathcal{X} \rightarrow \{0, 1\}$, if (IV) realizability holds (with respect to H, D, f), then running A on

$$m \geq m_H(\epsilon, \delta)$$

iid instances from D (labeled by f), returns a hypothesis h such that with probability of at least $1 - \delta$, it holds that

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

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

We will now move to a more general setting.

- ① we will first remove the realizability assumption; this means, there is no longer a $h^* \in H$ with $L_{D,f}(h^*) = 0$. Now, the best possible thing that we can hope for is a guarantee relative to $\min_{h \in H} L_{D,f}(h)$.

↳ before, with realizability, this was 0!!!

Def. (Hoeffding inequality): Let X_1, \dots, X_m iid random variables (RVs) taking values in $[a_i, b_i]$ for $i \in \{1, \dots, m\}$. Then, given that

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

it holds that

$$(a) \quad \mathbb{P}[S_m - \mathbb{E}[S_m] > \varepsilon] \leq e^{-\frac{-2\varepsilon^2}{\sum_i (b_i - a_i)^2}}$$

$$(b) \quad \mathbb{P}[S_m - \mathbb{E}[S_m] < -\varepsilon] \leq \text{---}$$

and, upon combination (using union bound)

$$\mathbb{P}\left[\left|S_m - \mathbb{E}[S_m]\right| > \varepsilon\right] \leq 2 \cdot e^{\frac{-2\varepsilon^2}{\sum_i (b_i - a_i)^2}}$$

An alternative (quite useful) form is:

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

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

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

$$S|x \sim D^m \quad \mathbb{P}\left[\left|L_S(h) - L_{D,f}(h)\right| > \varepsilon\right] \leq 2 \cdot e^{-2\varepsilon^2 m}$$

since $a=0$ and $b=1$.

Importantly, this only holds for a single h !!

What we want is a bound that holds uniformly over all $h \in H$.

So, for $|H| < \infty$ (finite classes), we can get (by the union bound):

$$\underset{S|x \sim D^m}{\mathbb{P}} \left[\exists h \in H : |L_S(h) - L_{D,f}(h)| > \varepsilon \right] \leq \sum_{h \in H} 2e^{-2\varepsilon^2 m}$$
$$\leq |H| \cdot 2e^{-2\varepsilon^2 m}$$

Remark: we did not need realizability for that!!! But, we pay the price of ε^2 vs. ε .

- ② Next, we will get rid of the true labeling function $f: X \rightarrow Y$ (which labeled our training instances x_i) by letting D be a distribution over $X \times Y$, ie., over domain and label space (Y).

Let $Z := X \times Y$. We need to adjust our definitions of empirical error and generalization error.

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

$$L_S(h) = \frac{1}{m} \cdot \left| \left\{ i \in \{1, \dots, m\} : h(x_i) \neq y_i \right\} \right|$$

Def. (Agnostic PAC learnability): A hypothesis class H of functions $h: X \rightarrow Y$ is **agnostic PAC learnable** if $\exists m_H: (0, 1)^2 \rightarrow \mathbb{N}$ and an algorithm A with the following properties:

- (I) for every $\varepsilon, \delta \in (0, 1)$ and every distribution D over $X \times Y = Z$, when running A on m iid instances from D with $m \geq m_H(\varepsilon, \delta)$, A returns a hypothesis h such that with probability of at least $1 - \delta$ (over the choice of S)

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

One more generalization:

Loss function:

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

Example: "0-1 loss"

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

With this, we can write

$$L_D(h) = \mathbb{E}_{z \sim D} [\ell(h, z)] \quad \text{with} \quad z = (x, y)$$

$$L_S(h) = \frac{1}{m} \cdot \sum_{i=1}^m \ell(h, z_i)$$

If we set in $\ell \equiv \ell^{0-1}$, then

$$L_D(h) = \mathbb{E}_{z \sim D} [\ell(h, z)] = \mathbb{E}_{z \sim D} [\underbrace{\ell^{0-1}(h, z)}_{\text{only takes values 0 or 1}}]$$

$$= 0 \cdot \mathbb{P}_{(x, y) \sim D} [h(x) = y] + 1 \cdot \mathbb{P}[h(x) \neq y]$$

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

Uniform Convergence

Def. (ϵ -representative sample): A sample S is called ϵ -representative with respect to $Z = X \times Y$, 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 $\epsilon/2$ -representative wrt. Z, H, l and D . Then, any hypothesis h_S returned by $\text{ERM}_H(S) \in \arg\min_{h' \in H} L_S(h')$ satisfies

$$L_D(h_S) \leq \min_{h \in H} L_D(h) + \frac{\epsilon}{2}$$

Proof: for any $h \in H$:

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

$$\text{ERM hyp.} \quad \leftarrow \quad \leq L_S(h) + \frac{\epsilon}{2} \quad // \text{since } h_S \text{ is an ERM hypothesis}$$

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

$$= L_D(h) + \epsilon$$

$$\text{Overall} \quad L_D(h_S) \leq L_D(h) + \epsilon \quad (\times)$$

Finally, as (x) holds for any $h \in H$, we conclude that

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

□

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

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

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

by ERH.

Question: Is there a universal learner? By universal we mean a learner without any prior knowledge of the learning task (given by D), but can be challenged by any task and still returns $A(S)$ with low $L_D(A(S))$.

This question is answered by the following theorem:

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

1. $\exists f: X \rightarrow \{0, 1\}^m$ with $L_D(f) = 0$!

and

2. with probability of at least $\frac{1}{4}$ over the choice of $S \sim D^m$, we have

$$L_D(A(S)) \geq \frac{1}{8}.$$

Interpretation: 1 means that the task can be learned successfully by another learner (e.g. using $H - \{f\}$) and 2 means that A fails on that task.

Corollary: Let X be infinitely large and H be the class of all functions from $X \rightarrow \{0, 1\}$, then H is not PAC learnable.

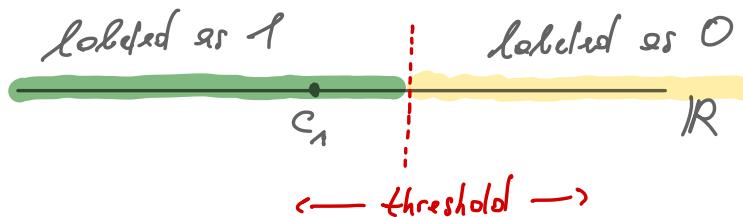
VAPNIK - CHERVONENKIS Dimension (VC-Dimension)

Daf.: 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

$$H_C = \{(h(c_1), h(c_2), \dots, h(c_m)) : h \in H\}$$

as the restriction of H to the set C .

Example: $C = \{c_1\}$, $c_1 \in \mathbb{R}$ and let H^{thr} be the class of thresholds on the real line.



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

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



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

But, we would have $2^2 = 4$ possible labelings of $\{c_1, c_2\} = C$.

Def. (Shattering) H shatters a finite set C of size m
if $|H_C| = 2^m$.

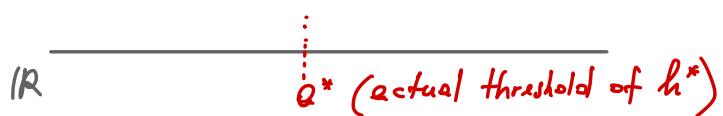
Claim: The class H^{thr} on $X = \mathbb{R}$ is PAC learnable
with $m_{H^{thr}}(\epsilon, \delta) \leq \left\lceil \log\left(\frac{2}{\delta}\right) \cdot \frac{1}{\epsilon} \right\rceil$.

$$H^{thr} = \{h_a : a \in \mathbb{R}\}$$

$$h_a : \mathbb{R} \rightarrow \{0, 1\}, \quad h_a(x) = \begin{cases} 1, & \text{if } x < a \\ 0, & \text{else} \end{cases}$$

$$\left(\text{or } h_a(x) = \underline{\underline{1}}_{x < a} \right)$$

We assume realizability, i.e., $\exists h^* \in H^{\text{thr}}$ such that $L_{\text{Df}}(h^*) = 0$.



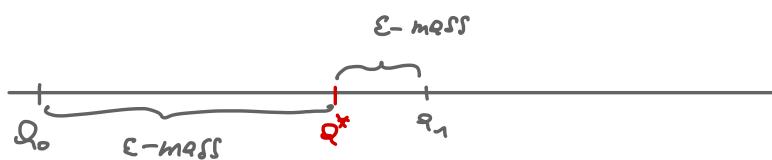
We let $q_0 \in \mathbb{R}$ be such that

$$\mathcal{D}\left(\{x \in \mathbb{R} : x \in (q_0, q^*)\}\right) = \varepsilon$$

and let $q_1 \in \mathbb{R}$ be such that

$$\mathcal{D}\left(\{x \in \mathbb{R} : x \in (q^*, q_1)\}\right) = \varepsilon$$

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



Next, we are going to define an Err algorithm (needed for PAC)



Pick $b_0 \in \mathbb{R}$ and $b_1 \in \mathbb{R}$ from $S = ((x_1, y_1), \dots, (x_m, y_m))$ as

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

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

Then, pick any threshold within (b_0, b_1) and call the corresponding hypothesis h_S . We see, h_S always has 0 empirical error on S .

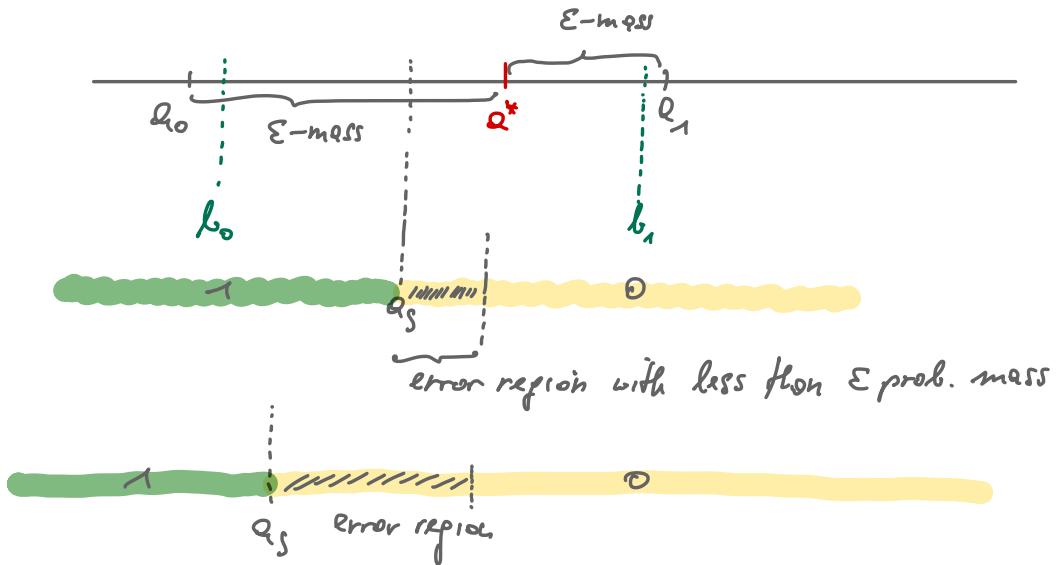
IRL hyp.

By our construction of q_0 and q_1 , the claim is that for $\underline{h_S}$ to have $L_{D,f}(h_S) \leq \varepsilon$, it suffices that

\downarrow
threshold
 q_S

- (a) $b_1 \leq q_1$ AND
- (b) $b_0 \geq q_0$

(x)



1st exam date: June 29 (10am) - last lecture

If we write down (x) formally, we have

$$\begin{aligned} \mathbb{P}[L_{D,f}(h) > \varepsilon] &\leq \mathbb{P}[(b_0 < q_0) \text{ or } (b_1 > q_1)] \\ &\leq \mathbb{P}[b_0 < q_0] + \mathbb{P}[b_1 > q_1] \quad // \text{union bound} \end{aligned}$$

To upper-bound $\mathbb{P}[L_{D,f}(h) > \varepsilon]$, we need to check when $b_0 < q_0$ (or $b_1 > q_1$). Answer: when there is no datapoint x in S that is labeled '1' such that $x \in (q_0, q^*)$.

We know that $D(\{x \in \mathbb{R} : x \in (q_0, q^*)\}) = \varepsilon$ (by construction of q_0); hence, not seeing a data point in (q_0, q^*) has probability q_0 ; consequently, not seeing a data point in m iid samples (inside (q_0, q^*)) has probability $(1-\varepsilon)^m$.

\Rightarrow Overall, we have

$$\begin{aligned} \mathbb{P}[L_{D,f}(h) > \varepsilon] &\leq 2 \cdot \underbrace{(1-\varepsilon)^m}_{\leq e^{-\varepsilon m}} \leq \underbrace{2e^{-\varepsilon m}}_{\delta} \\ &\Rightarrow m > \frac{1}{\varepsilon} \log\left(\frac{2}{\delta}\right) \end{aligned}$$

Sample complexity

We have seen that "finiteness" of H 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 $\text{VC}(H)$,

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 theorem).

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

$$\gamma_H(m) = \max_{\substack{C \subset X \\ |C|=m}} |H_C|$$

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

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

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

VC-dimension

Let's take any set C (of size m) and $|H| < \infty$ (finite).

We know

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

So, if $|H| < 2^m$, then H can obviously not shatter C of size m . This implies

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

Theorem: Let H be a class of functions from $X \rightarrow \{0,1\}$ and $\ell : H \times X \times Y \rightarrow [0, c]$, $c > 0$, 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 $S \sim D^m$)

$$\forall h \in H : |L_D(h) - L_S(h)| \leq c \cdot \sqrt{\frac{\delta \cdot \log(\tilde{V}_H(2m))^4}{m}}$$

(without proof).