

STATISTICAL LEARNING THEORY

Summer Term 2025

Book: Shai Shalev-Shwartz
Understanding Machine Learning

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

OTHER COURSE NAMES:

- 1) Advanced machine learning
- 2) Machine learning

MOTIVATION

Example:

| | $\in \mathbb{R}$ weight (in g) | $\in \mathbb{R}$ color ($\in [0,1]$) | Testy? |
|----------|-----------------------------------|---|----------|
| Peppya 1 | 800 | 0.1 | 0 |
| \vdots | \vdots | \vdots | \vdots |
| Peppya N | 1200 | 0.7 | 1 |

0... non-testy
1... testy

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 peppya is testy (1) or not (0).

we call such a function a hypothesis.

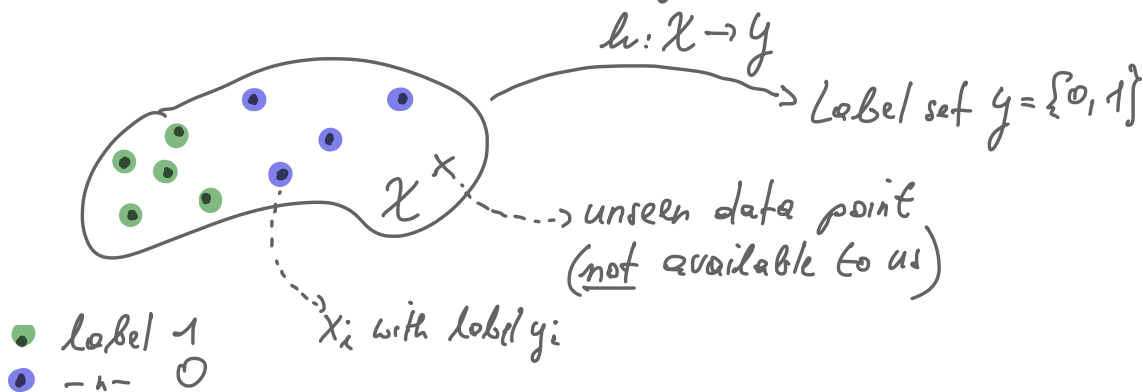
More formally, the data that is available to us comes in the form

$$((x_1, y_1), \dots, (x_n, y_n)) = S$$

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

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

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 X .

(2) The x_i 's are labeled by some (unknown) function
 $f: X \rightarrow Y$,

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

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

That is, all the points x in our domain 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 \mathcal{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:

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

"Generalization error"

The empirical version of that is

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

"Empirical error"
(Empirical risk)

Notation:

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

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

$S \sim \mathcal{D}^m$
(x_1, \dots, x_m)

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

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^m} [L_S(h)] &= \mathbb{E}_{S \sim \mathcal{D}^m} \left[\frac{1}{m} \cdot \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq f(x_i)} \right] \quad // \text{by def.} \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{1}_{h(x) \neq f(x)}] \quad // \text{by linearity of } \mathbb{E}[\cdot] \end{aligned}$$

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

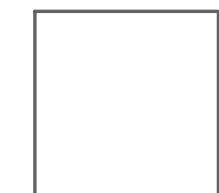
$$= \frac{1}{m} \cdot \sum_{i=1}^m \mathbb{P}[h(x) \neq f(x)]$$

$$= \frac{1}{m} \cdot \cancel{m} \cdot \mathbb{P}[h(x) \neq f(x)] = \underbrace{L_{D,f}(h)}_{\text{generalization error}}$$

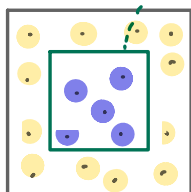
Our first learning paradigm

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 \mathcal{X}



true labeling function f

Say the distribution on \mathcal{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 \mathcal{D} (which is uniform on \mathcal{X}), h_S is only correct 50% of the time \rightarrow

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

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

Also, we define

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

Observation:

$$\{S/x: L_{D,f}(h_S) \geq \varepsilon\} \subseteq \{S/x: \exists h \in H_{\text{BAD}}, L_S(h) = 0\} = H$$

\uparrow
 ERM hypothesis
 (Empirical risk minimizer)

Since $\{S/x: \exists h \in H_{\text{BAD}}, L_S(h) = 0\} = \bigcup_{h \in H_{\text{BAD}}} \{S/x: L_S(h) = 0\}$

we have

$$\mathcal{D}^m\left(\{S/x: \exists h \in H_{\text{BAD}}, L_S(h) = 0\}\right) = \mathcal{D}^m\left(\bigcup_{h \in H_{\text{BAD}}} \{S/x: L_S(h) = 0\}\right)$$

(by σ -sub-additivity
"union" bound) $\leq \sum_{h \in H_{\text{BAD}}} \mathcal{D}^m(\{S/x: L_S(h) = 0\})$

Lets fix some $h \in H_{\text{BAD}}$:

$$\mathcal{D}^m(\{S/x: L_S(h) = 0\}) = \mathcal{D}^m(\{S/x: \forall i \in \{1, \dots, m\}: h(x_i) = \overset{\substack{\text{True labeling} \\ \text{function}}}{\downarrow} f(x_i)\})$$

by iid assumption \rightarrow

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

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

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

\hookrightarrow

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

Overall, we have

$$\begin{aligned}
\mathcal{D}^m(\{S|_x : L_{\mathcal{D}, P}(h_S) \geq \varepsilon\}) &\leq \sum_{h \in H_{\text{BAD}}} \mathcal{D}^m(\{S|_x : L_S(h) = 0\}) \\
&\leq \sum_{h \in H_{\text{BAD}}} e^{-\varepsilon m} \\
&= |H_{\text{BAD}}| \cdot e^{-\varepsilon m} \\
&\leq |H| \cdot e^{-\varepsilon m} \quad (\text{because } H_{\text{BAD}} \subseteq H)
\end{aligned}$$

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 ↓ confidence