# EMPIRICAL RISK MINIMIZATION

Recall from the last lecture :

$$R(f) = \mathop{\mathbb{E}}_{(x,y) \sim P} \left[ \ell( f(x), y) \right]$$

for any $f: \mathbb{X} \to \mathbb{Y}$ (measurable). Then, if $P$ was known :

$$f^*(x') \in \mathop{\arg\min}_{y' \in \mathbb{Y}} \mathbb{E}[\ell( y', y) \mid x=x'], \quad \forall x' \in \mathbb{X},$$

is a Bayes optimal predictor, i.e., $R(f^*) = R^* = \inf_{\substack{f: \mathbb{X} \to \mathbb{Y} \\ f \text{ is meas.}}} R(f)$

Question : What if $P$ is unknown?

Only $(x_i, y_i) \overset{iid}{\sim} P$, $i=1,2,\ldots,n$ is provided. $\longrightarrow$ partial info.

$\longrightarrow$ Supervised learning

## Empirical Risk Minimization (ERM) :

Given $S = \{(x_i, y_i) \in \mathbb{X} \times \mathbb{Y} : i=1,2,\ldots,n\}$, let

$$\hat{R}_S(f) = \frac{1}{n} \sum_{i=1}^{n} \ell( f(x_i), y_i), \qquad f: \mathbb{X} \to \mathbb{Y}, \text{ meas.}$$

$\hat{R}_S(f)$ is the empirical risk of $f$.

Idea : Use $\hat{f}_{ERM} \in \mathop{\arg\min}_{\substack{f: \mathbb{X} \to \mathbb{Y} \\ f \text{ is meas.}}} \hat{R}_S(f)$.

Hope : The "excess risk" $R(\hat{f}_{ERM}) - R^*$ is small (in expectation, or with high probability, since $S$ and thus $\hat{R}_S$ and $\hat{f}_{ERM}$ are random.

In the following, we show that ERM (empirical risk minimization) is a reasonable idea.

__PROPOSITION__ ( $\hat{R}_S(f)$ is an unbiased and consistent estimate of $R(f)$ )

Let $f: \mathbb{X} \to \mathbb{Y}$ be a given measurable predictor, and $\ell: \mathbb{Y} \times \mathbb{Y} \to [-B, B]$ for some $B \in \mathbb{R}_+$.

Then,

(i) $\quad \mathbb{E}[\hat{R}_S(f)] = R(f),$

(ii) $\quad \lim_{n \to \infty} \hat{R}_S(f) = R(f)$ almost surely,

(iii) for any $\delta \in (0,1)$,

$$\mathbb{P}\left( |\hat{R}_S(f) - R(f)| \leq B\sqrt{\frac{2\log(2/\delta)}{n}} \right) \geq 1 - \delta.$$

**Proof**

(i) $\quad \mathbb{E}[\hat{R}(f)] = \mathbb{E}\left[ \frac{1}{n}\sum_{i=1}^{n} \ell(f(x_i), y_i) \right]$

$= \frac{1}{n}\sum_{i=1}^{n} \mathbb{E}[\ell(f(x_i), y_i)] = \frac{1}{n}\sum_{i=1}^{n} R(f) = R(f),$

since $(x_i, y_i) \overset{iid}{\sim} P$ for each $i = 1, 2, \dots, n$.

(ii) Let $z_i := \ell(f(x_i), y_i)$ for $i = 1, 2, \dots, n$. Since $(x_i, y_i) \overset{iid}{\sim} P$,

$(z_i)_i$ is an iid sequence, and $\mathbb{E}[z_i] = R(f), \forall i \in [n]$. Also,

since $\ell : \mathcal{X} \times \mathcal{Y} \to [-B, B]$, $|z_i| \leq B$ a.s. for all $i \in [n]$. Thus,

by the strong law of large numbers,

$$\frac{1}{n}\sum_{i=1}^{n} z_i \xrightarrow[n \to \infty]{} \mathbb{E} z_1 = R(f) \quad \text{almost surely}.$$

(iii) By (ii),

$$\mathbb{E}[z_i] = R(f), \forall i,$$
$$|z_i| \leq B, \quad \text{and} \quad (z_i)_i \text{ is an iid seq.}$$

Thus, by Hoeffding inequality,

$$\mathbb{P}\left( |\hat{R}_S(f) - R(f)| > B\sqrt{\frac{2\log(2/\delta)}{n}} \right)$$

$$= \mathbb{P}\left( \left| \frac{1}{n}\sum_{i=1}^{n} z_i - \mathbb{E}[z_1] \right| > B\sqrt{\frac{2\log(2/\delta)}{n}} \right)$$

$$\leq \delta.$$

**(Important) Remarks:**

① The performance criterion is $R(f)$, i.e., the population risk. We use $\hat{R}_S$ _only_ as a tool to construct a good predictor $\hat{f}$ to perform well in terms of $R$.

② The results of the above proposition is pointwise (for every given $f$), _not_ uniform that hold simultaneously for all $f$.

# Level End Boss : Overfitting

$\hat{R}_s(f)$ is a consistent and unbiased estimator of $R(f)$, $\forall f$,

and $\hat{f}_{ERM} \in \underset{\substack{f: \mathcal{X} \to \mathcal{Y} \\ f \text{ is meas.}}}{\arg\min} \hat{R}_s(f)$.

Forget about the computational complexity of finding $\hat{f}_{ERM}$ (for now).

**Question:** Does $\hat{R}_s(f) = 0$ imply small (or vanishing with $n$) excess risk $R(f) - R^*$ ?
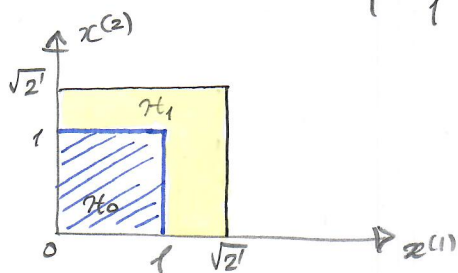
The answer is **no**.

**Example:** Let $\mathcal{X} = \mathbb{R}^2$, and $\mathcal{Y} = \{0, 1\}$.

The density function of the input is:

$$P_X(x) = \begin{cases} \frac{1}{2}, & x \in [0, \sqrt{2}] \times [0, \sqrt{2}], \\ 0, & \text{otherwise.} \end{cases}$$

Also, $y = f^*(x) = \begin{cases} 0, & \text{if } x \in [0,1] \times [0,1], \\ 1, & \text{otherwise.} \end{cases}$



Consider the following predictor:

$$\hat{f}(x) = \begin{cases} y_i, & \text{if } x = x_i \text{ for some } i = 1,2,\dots n \\ 1, & \text{otherwise.} \end{cases}$$

Then, if we consider $\ell(y, y') = \mathbb{1}\{y \neq y'\}$,

$$\hat{R}_s(\hat{f}) = \frac{1}{n} \sum_{i=1}^{n} \ell(\hat{f}(x_i), y_i) = 0.$$

On the other hand,

$$R(f) = \int_{\mathcal{H}_0} P_X(x)\,dx = \frac{1}{2}.$$

<u>Perfect fit</u> to the <u>training data</u> implied terrible population risk performance on a <u>test data</u>. $\longrightarrow$ overfitting

# Inductive Bias to Avoid Overfitting:

For a given problem, we previously had a humongous hypothesis class:

$$\mathcal{H}^* = \{f : \mathbb{X} \to \mathbb{Y} : f \text{ is measurable}\}.$$

This richness resulted in overfitting.

The traditional way to avoid overfitting is to use a **restricted** hypothesis class. For some well-chosen $\mathcal{H} \subset \mathcal{H}^*$, let

$$\hat{f}_{\mathcal{H}} \in \arg\min_{f \in \mathcal{H}} \hat{R}_S(f),$$

and also $\quad R_{\mathcal{H}}^* := \inf_{f \in \mathcal{H}} R(f).$ Then, we have:

$$\boxed{R(\hat{f}_{\mathcal{H}}) - R^* = \underbrace{R(\hat{f}_{\mathcal{H}}) - R_{\mathcal{H}}^*}_{\substack{\text{estimation} \\ \text{error}}} + \underbrace{R_{\mathcal{H}}^* - R^*}_{\substack{\text{approximation} \\ \text{error.}}}}$$

excess risk decomposition under inductive bias.

**Remarks**: There is a tradeoff between the estimation error and the approximation error in the decomposition above:

(i) Large $\mathcal{H}$ $\Rightarrow$ small approximation error
    but
    large estimation error (larger search)

(ii) Large $n = |S|$ $\Rightarrow$ smaller estimation error
    but
    no impact on the approximation error.

## Some examples:

① Linear regressors: $\mathbb{X} = \mathbb{R}^d,\ \mathbb{Y} = \mathbb{R},$

$$\mathcal{H} = \{x \mapsto \theta^T x : \theta \in \mathbb{R}^d,\ \|\theta\|_2 \leq \rho\}$$

② Linear classifiers: $\mathbb{X} = \mathbb{R}^d,\ \mathbb{Y} = \{-1, 1\},$

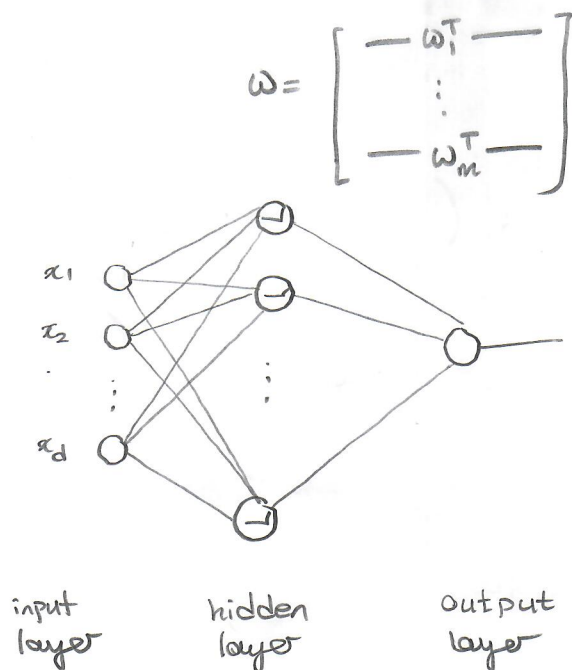$$\mathcal{H} = \{x \mapsto \text{sgn}(\theta^T x - b) : (\theta, b) \in \mathbb{R}^{d+1}\}$$

③ Quantized perceptrons : $\underline{\mathbb{X}} = \mathbb{R}^d$, $\underline{\mathbb{Y}} = \{-1, +1\}$.

$\mathcal{H} = \{ x \mapsto \text{sgn}(\omega^T x - b) : \omega_1, \omega_2, \ldots, \omega_d, b \in \mathbb{F}_k \text{ for some } k \in \mathbb{N} \}$

where $\mathbb{F}_k = \{0,1\}^k$, $k$

④ Shallow neural networks :

$\mathcal{H} = \{ x \mapsto \sum_{i=1}^{m} c_i \, \sigma(\omega_i^T x - b_i) : c, b \in \mathbb{R}^m, \omega \in \mathbb{R}^{m \times d} \}$

$$\omega = \begin{bmatrix} -\omega_1^T- \\ \vdots \\ -\omega_m^T- \end{bmatrix}$$



input    hidden    output
layer    layer    layer

## Capacity Control and Explicit Regularization :

As we mentioned, the traditional way to address overfitting is to consider a restricted hypothesis class $\mathcal{H} \subsetneq \mathcal{H}^*$.

**Parameterization :** Idea is to define the hypothesis class in a parametric way. Let $\Theta$ be a parameter set, and

$$\mathcal{H}_\Theta = \{ x \mapsto f_\theta(x) \in \mathbb{Y} : x \in \underline{\mathbb{X}}, \theta \in \Theta \}$$

is a parametric hypothesis class.

Then,

$$R_{\mathcal{H}_\Theta}^* = \inf_{f \in \mathcal{H}_\Theta} R(f) = \inf_{\theta \in \Theta} R(f_\theta),$$

$$\hat{\theta}_{ERM} \in \underset{\theta \in \Theta}{\text{argmin}} \; \hat{R}_S(f_\theta).$$

**Examples :** ① $\mathcal{H}_\Theta = \{ x \mapsto \theta^T x : \theta \in \Theta \}$ where $\Theta \subset \mathbb{R}^d$.

② $\Theta \subset \mathbb{R}^{m \times d} \times \mathbb{R}^m \times \mathbb{R}^m$, $m \in \mathbb{Z}+$,

$\mathcal{H}_\Theta = \{ x \mapsto \sum_{i=1}^{m} c_i \, \sigma(\omega_i^T x - b_i) : (\omega, c, b) \in \Theta \}$.

Recall that the rationale was to control the richness of the hypothesis class, $\mathcal{H}$ or $\mathcal{H}_\Theta$.

Question: Can we further control the richness of $\mathcal{H}_\Theta$?
"capacity control."
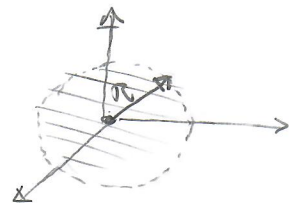
Answer: Yes. Via (explicit) regularization.

EXAMPLE
Consider $\mathcal{H}_\Theta = \{ x \mapsto \theta^T x : \theta \in \textcircled{H} \subset \mathbb{R}^d \}$,

$$R(f) = \mathop{\mathbb{E}}_{(x,y) \sim P} \left[ (\theta^T x - y)^2 \right].$$

How to control the richness?

①  $\mathcal{H}_{\Theta,R} = \{ x \mapsto \theta^T x : \theta \in \textcircled{H} \subset \mathbb{R}^d, \; \underline{\|\theta\|_2 \leq R} \}$,

by a design parameter $R > 0$.

Hypotesis class is now a ball of radius $R > 0$ around the origin in $d$-dimensional Euclidean space.

increasing $R$ : (i) decreases the approximation error



$$\inf_{f \in \mathcal{H}_{\Theta,R}} R(f) - R^*$$

since  $\mathcal{H}_{\Theta,R} \subset \mathcal{H}_{\Theta,R'}$ for any $R \leq R'$.

(ii) increases the optimization and generalization errors, since we have a larger set of candidates.

② (Tikhonov regularization)

For a design parameter $\lambda > 0$, consider

$$R^\lambda(f) = R(f) + \lambda \cdot \|\theta\|_2^2, \quad f \in \mathcal{H}_\Theta.$$

In this case, the (regularized) objective is

$$\min_{\theta \in \textcircled{H} = \mathbb{R}^d} R^\lambda(f_\theta).$$

As you notice, $\lambda > 0$ implies larger penalty for large $\|\theta\|_2$.
⇒ incentivizes using small $\theta \in \mathbb{R}^d$ ⟶ (Norm-based control)

# PAC (Probably Approximately Correct) Learnability :

In a given restricted hypothesis class $\mathcal{H} \subset \mathcal{H}^*$, the goal is to come up with $\hat{f} \in \mathcal{H}$ s.t.

$$\mathbb{P}\left( R(\hat{f}) \leq R_{\mathcal{H}}^* + \varepsilon \right) \geq 1 - \delta, \qquad \varepsilon > 0, \ \delta \in (0,1),$$

by using a sufficiently large but <u>finite</u> training set $S$.

<u>Question</u> : Is this possible for $\mathcal{H} \subset \mathcal{H}^*$?

<u>DEF)</u> (Learning algorithm) For a given learning problem specified by $(\mathbb{X}, \mathbb{Y}, P)$, and a class of hypotheses $\mathcal{H} \subset \mathcal{H}^*$, a learning algorithm is a sequence $A = (A_n)_{n=1}^{\infty}$ of mappings

$$A_n : (\mathbb{X} \times \mathbb{Y})^n \longrightarrow \mathcal{H}.$$

<u>DEF)</u> (PAC - Learnability) Let $\mathcal{H} \subset \mathcal{H}^*$. $\mathcal{H}$ is PAC-learnable if: there is a learning algorithm $A$ such that :
for any $\varepsilon > 0$, $\delta \in (0,1)$, there is an integer $n_0(\varepsilon, \delta)$ s.t. for any $P$ on $\mathbb{X} \times \mathbb{Y}$, if for any $n \geq n_0$, $S_n \sim P^n$,
then $\qquad \mathbb{P}\left( R(A_n(S_n)) < R_{\mathcal{H}}^* + \varepsilon \right) \geq 1 - \delta.$

<u>Important remark</u> : PAC-learnability does <u>not</u> always hold. For some important hypothesis classes, we will prove their PAC-learnability.

As we will see,

- Any finite hypothesis classes,
- Linear regressors with bounded parameter norm :
$$\mathcal{H} = \{ x \mapsto \theta^T x \ : \ \|\theta\|_2 \leq R \},$$
- Neural network with large width $m \geq 0$ $m$ a certain operating regime,
- ReLU networks with bounded parameter norm,

will be PAC-learnable.

# Sufficiency of Uniform Convergence for PAC-Learnability:

Now, we are given a training set $S$, and for a hypothesis class $\mathcal{H} \subset \mathcal{H}^*$, an algorithm $L_0$ returns $\hat{f}$.

$\underline{\text{Recall}}$:

$$\hat{f}_{\mathcal{H}} \in \underset{f \in \mathcal{H}}{\arg\min} \; \hat{R}_S(f),$$

$$R_{\mathcal{H}}^* = \inf_{f \in \mathcal{H}} R(f), \qquad f_{\mathcal{H}}^* \in \underset{f \in \mathcal{H}}{\arg\min} \; R(f)$$

Then, we have the following error decomposition:

$$R(\hat{f}) - R_{\mathcal{H}}^* = \underbrace{R(\hat{f}) - \hat{R}(\hat{f})}_{(1)} + \underbrace{\hat{R}_S(\hat{f}) - \hat{R}_S(\hat{f}_{\mathcal{H}})}_{(2)} + \underbrace{\hat{R}_S(\hat{f}_{\mathcal{H}}) - \hat{R}_S(f_{\mathcal{H}}^*)}_{(3)} + \underbrace{\hat{R}_S(f_{\mathcal{H}}^*) - R_{\mathcal{H}}^*}_{(4)}$$

Let's examine these error terms:

$(3)$: Note that $\hat{f}_{\mathcal{H}} \in \underset{f \in \mathcal{H}}{\arg\min} \; \hat{R}_S(f)$, and $f_{\mathcal{H}}^* \in \mathcal{H}$. Thus, this term is $\underline{\text{non-positive.}}$

$(4)$: Recall from Lecture :

$$Z_i = \ell\left(f_{\mathcal{H}}^*(x_i), y_i\right), \quad i = 1, 2, \dots, n.$$

since $(x_i, y_i)_{i=1}^n$ are iid, $(Z_i)_i$ are iid, also $\mathbb{E}[Z_i] = R_{\mathcal{H}}^*$. Thus,

$$\left| \hat{R}_S(f_{\mathcal{H}}^*) - R_{\mathcal{H}}^* \right| \leq B \sqrt{\frac{2 \log(2/\delta)}{n}} \qquad \text{w.p.} \geq 1 - \delta.$$

$(2)$: Finding $\hat{f}_{\mathcal{H}}$ may be an NP-hard problem. This term accounts for the optimization error in finding $\hat{f}_{\mathcal{H}}$.

$(1)$: Same as $(4)$? $\underline{\text{Not quite.}}$

$$\tilde{Z}_i = \ell\left(\hat{f}(x_i), y_i\right), \quad i = 1, 2, \dots, n.$$

Note that $\hat{f}$ is $\sigma(S)$-measurable, thus $\tilde{Z}_i$ and $\tilde{Z}_j$ are correlated for $i \neq j$. Furthermore,

$$\mathbb{E}\left[\ell\left(\hat{f}(x_i), y_i\right)\right] \neq R(\hat{f}) \quad \text{since } \hat{f} \text{ is } \sigma(S)\text{-measurable.}$$

Uniform convergence over $\mathcal{H}$:

$$\underbrace{R(\hat{f}) - R_{\mathcal{H}}^*}_{\text{excess risk}} \leq R(\hat{f}) - \hat{R}_s(\hat{f}) + \hat{R}_s(\hat{f}) - \hat{R}_s(\hat{f}_{\mathcal{H}}) + \hat{R}_s(f_{\mathcal{H}}^*) - R(f_{\mathcal{H}}^*)$$

$$\leq |R(\hat{f}) - \hat{R}_s(\hat{f})| + |\hat{R}_s(f_{\mathcal{H}}^*) - R(f_{\mathcal{H}}^*)| + \hat{R}_s(\hat{f}) - \hat{R}_s(\hat{f}_{\mathcal{H}})$$

$$\leq 2 \cdot \sup_{f \in \mathcal{H}} |R(f) - \hat{R}_s(f)| + \underbrace{\hat{R}_s(\hat{f}) - \hat{R}_s(\hat{f}_{\mathcal{H}})}_{\text{optimization error}}.$$

## THEOREM 1 (Finite $\mathcal{H}$)

Suppose that $|\mathcal{H}| < \infty$. Then, if $\ell : \mathcal{Y} \times \mathcal{Y} \to [-B, B]$,

$$\mathbb{P}\left( \sup_{f \in \mathcal{H}} |R(f) - \hat{R}_s(f)| \leq B\sqrt{\frac{2 \log(2|\mathcal{H}|/\delta)}{n}} \right) \geq 1-\delta.$$

Thus,

$$\hat{R}(\hat{f}_{\mathcal{H}}) - R_{\mathcal{H}}^* \leq \frac{2B\sqrt{2 \log(2 \cdot |\mathcal{H}|/\delta)}}{\sqrt{n}} \qquad \text{w.p. } \geq 1-\delta,$$

hence any finite $\mathcal{H}$ is PAC-learnable since

$$n \geq \frac{8B^2 \log(2 \cdot |\mathcal{H}|/\delta)}{\varepsilon^2}$$

implies

$$\hat{R}(\hat{f}_{\mathcal{H}}) \leq R_{\mathcal{H}}^* + \varepsilon \qquad \text{w.p. } \geq 1-\delta.$$