# Empirical Risk Minimization

→ What is SUPERVISED LEARNING?

List of INPUTS:
$x_1, \ldots, x_n \in \mathcal{X}$

List of OUTPUTS:
$y_1, \ldots, y_n \in \mathcal{Y}$

GOAL: Given a new input $x$, predict the corresponding $y$.

Two large families:

CLASSIFICATION $\qquad \mathcal{Y} = \{a, b, c, \ldots\}$

REGRESSION $\qquad \mathcal{Y} = \mathbb{R}$

Examples:

| INPUTS | | OUTPUTS |
|---|---|---|
| Pictures | $\longrightarrow$ | Objects |
| Movie reviews | $\longrightarrow$ | Review rating |
| Patient | $\longrightarrow$ | Is the patient cured? |

# ① Risks and losses :

⚠️ There is no single good notion to quantify the quality of a prediction.

**Def** : A loss function is a function

$$\ell : \mathcal{Y} \times \mathcal{Y} \longrightarrow [0, \infty).$$

**Examples:** ① Binary classification $\mathcal{Y} = \{-1, +1\}$

- **Choice 1** : $\ell(y, y') = \begin{cases} 1 & \text{if } y \neq y' \\ 0 & \text{otherwise} \end{cases}$

- **Choice 2** : $\ell(y, y') = \begin{cases} a & \text{if } y = 1 \text{ and } y' = -1 \\ b & \text{if } y = -1 \text{ and } y' = 1 \\ 0 & \text{if } y = y' \end{cases}$

⇝ Choice 2 makes sense (for instance) in medical settings, where it is a more serious mistake to predict that a patient is not sick ($y' = -1$) whereas they are ($y = +1$) than the opposite.

(2) Regression $y = \mathbb{R}^D$     $y = (y_1 \cdots y_D)$    $y' = (y'_1 \cdots y'_D)$

- $L_\infty$ norm     $\|y' - y\|_\infty = \max\limits_{i=1 \cdots D} |y_i - y'_i|$

- $L_p$ norm     $\|y' - y\|_p = \left( \sum\limits_{i=1}^{D} |y_i - y'_i|^p \right)^{1/p}$

- Weighted $L_p$ norm     $\left( \sum\limits_{i=1}^{D} w_i |y_i - y'_i|^p \right)^{1/p}$

GOAL: Find a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$ such that,
$\ell(f(x), y)$ is small on <u>New</u> samples $(x'_1, y'_1) \cdots$
                                                       $(x'_m, y'_m)$.

TRAINING SAMPLES
$(x_1, y_1)$
$\cdots$ $(x_n, y_n)$

$\leadsto$ Find a predictor $f$.

NEW SAMPLES
$(x'_1, y'_1) \cdots (x'_m, y'_m)$

We want $\ell(f(x'_i), y'_i)$ to be small on average.
(i.e. good predictions on new samples)

**Example:** Electric consumption Forecasting

Training samples: $y_i$ = Electric consumption on Day $i$

$x_i$ = characteristics of Day $i$ (weather, day of the week, etc.)

TRAIN A MODEL

Testing sample: Predict tomorrow's consumption based on tomorrow's weather, etc.

ASSUMPTION: The observations $(x_1, y_1) \dots (x_n, y_n)$ are i.i.d. with law $P$.

independant identically distributed.

**Def** Given a function $f: \mathcal{X} \longrightarrow \mathcal{Y}$, the P-risk of $f$ is defined as:
$$R_P(f) := \mathbb{E}_P[\ell(f(X), Y)]$$

**Theorem:** The P-risk is minimized for the Bayes predictor $f_P^*$ defined by
$$f_P^*(x) \in \underset{z \in \mathcal{Y}}{\arg\min} \; \mathbb{E}_P[\ell(Y, z) \mid X = x].$$

proof: Let $\Psi(x,z) = E_P[\ell(Y,z)|X=x]$

$\leadsto \Psi(x,z) \geqslant \Psi(x, f_P^*(x))$    (by definition)

$$R_P(f) = E_P[\ell(Y, f(x))] \overset{(*)}{=} E_P[\Psi(X, f(x))]$$
$$\geqslant E_P[\Psi(X, f_P^*(x))]$$
$$\overset{(*)}{=} E_P[\ell(Y, f_P^*(x))] = R_P(f_P^*)$$

(*) LAW OF TOTAL EXPECTATION
$$\mathbb{E}[A] = \mathbb{E}[\mathbb{E}[A|x]]$$

Examples:

① Binary classification:

$$\ell(y, y') = \begin{cases} 1 & \text{if } y \neq y' \\ 0 & \text{otherwise} \end{cases}$$

$$E_P[\ell(Y,z)|X=x] = P(Y \neq z | X=x)$$

$z = -1 \text{ or } +1$
$$= 1 - P(Y=z|X=x)$$

$\longrightarrow$ if $\eta(x) = P(Y=1|X=x) \geqslant \frac{1}{2}$,

then $f_P^*(x) = 1$.

Otherwise, $f_P^*(x) = -1$.

② Quadratic loss $\mathcal{Y} = \mathbb{R}$ $\quad \ell(y, y') = (y - y')^2$

Question: A random variable.
What is the minimum of

$$a \longmapsto \mathbb{E}[(A-a)^2] = \mathbb{E}[A^2] - 2\mathbb{E}[A]a + a^2$$

$$\partial_a \quad = 2(a - \mathbb{E}[A]) \quad \leadsto \quad a = \mathbb{E}[A]$$

$$\longrightarrow \quad f_P^*(x) = \mathbb{E}_P[Y | X = x]$$

② Empirical risk

$P$ is unknown $\Longrightarrow$ $f_P^*$ is unknown

GOAL: Approximate $f_P^*$ using the observations

$$((x_1, y_1), \ldots, (x_n, y_n))$$

Def: The ==empirical risk== of a function $f$ is

$$\boxed{R_n(f) = \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i)}$$

# LAW OF LARGE NUMBERS

$$R_n(f) = \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i)$$

$$\downarrow n \to +\infty$$

$$R_p(f) = \mathbb{E}_p[\ell(f(x), y)]$$

$\leadsto$ Minimizing $R_n(f) \approx$ Minimizing $R_p(f)$

Def: Let $\mathcal{F}$ be a set of functions from $x \to y$.

The **empirical risk minimizer** of $\mathcal{F}$ is

$$\left[ \hat{f}_{\mathcal{F}} \in \underset{f \in \mathcal{F}}{\arg\min} \; R_n(f) \right]$$

Examples:

① Linear Regression: $x = \mathbb{R}^d \quad y = \mathbb{R}$

$$\ell(y, y') = (y - y')^2$$

$$\mathcal{F}_{lin} = \{ f_\theta : x \mapsto \langle x, \theta \rangle \; : \; \theta \in \mathbb{R}^d \}$$

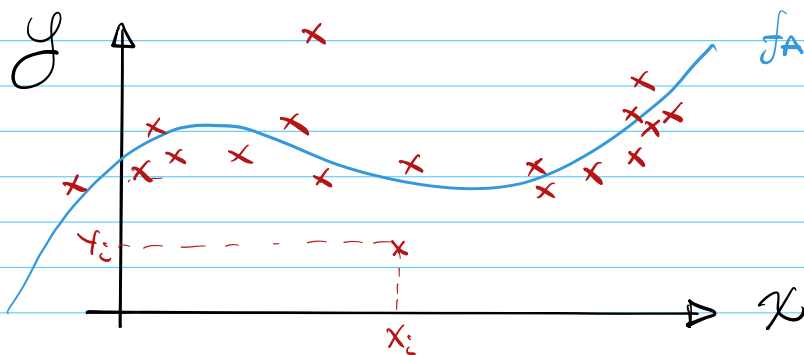$$R_n(f_\theta) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - f_\theta(x_i))^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (Y_i - \langle X_i, \theta \rangle)^2 \qquad \left.\begin{array}{c} \end{array}\right\} \text{ LINEAR REGRESSION}$$

$$= \frac{1}{n} \| \mathbb{Y} - \mathbb{X}\theta \|^2$$

$$\mathbb{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \qquad \mathbb{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} \updownarrow n$$

$$\xleftarrow{\quad d \quad}$$

② **Polynomial Regression**

$$\mathcal{X} = \mathbb{R} \qquad \mathcal{Y} = \mathbb{R} \qquad \ell(y, y') = |y - y'|^2$$

$$\mathcal{F}_d = \left\{ \underbrace{x \longmapsto \sum_{i=0}^{d} a_i x^i}_{= f_A} : \underbrace{a_0 \cdots a_d}_{= A \,\in\, \mathbb{R}^{d+1}} \in \mathbb{R} \right\}$$

d = 1
d = 2
d = 8
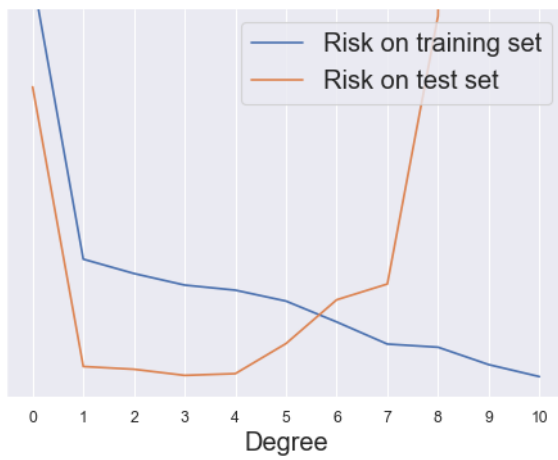


Let us explain this phenomenon.

⚠️ Training set $(X_1, Y_1) \ldots (X_n, Y_n)$

$\Rightarrow$ Find prediction $\hat{f}_{\mathcal{F}}$ **DEPENDING** on the Training set.

$\Rightarrow$ We can NOT apply the Law of Large numbers to say that $R_P(\hat{f}_{\mathcal{F}}) \approx R_n(\hat{f}_{\mathcal{F}})$

$$= \frac{1}{n} \sum_{i=1}^{n} \ell(\hat{f}_{\mathcal{F}}(x_i), Y_i)$$

not indep !!

$\Rightarrow$ On the testing sample $(X_1', Y_1') \ldots (X_m', Y_m')$ **INDEPENDENT** from the training sample,

LLN **CONDITIONALLY ON THE TRAINING SAMPLE**

$$\Downarrow$$

$$R_P(\hat{f}_{\mathcal{F}}) \approx \frac{1}{m} \sum_{i=1}^{m} \ell(\hat{f}_{\mathcal{F}}(x_i'), Y_i')$$

conditionally independent.

Good Approximation of $R_P(\hat{f}_{\mathcal{F}})$



Risk on training set
Risk on test set
Degree

# Decomposition of the Empirical Risk

Minimal Risk: $R_P^* = R_P(f_P^*) = \min\limits_f R_P(f)$

:

$$R_P(\hat{f}_{\mathcal{F}}) - R_P^* \qquad = R_P(f_{\mathcal{F}}^*)$$

$$= \left( R_P(\hat{f}_{\mathcal{F}}) - \inf\limits_{f \in \mathcal{F}} R_P(f) \right) + \left( \inf\limits_{f \in \mathcal{F}} R_P(f) - R_P^* \right)$$

$\underbrace{\hphantom{R_P(\hat{f}_{\mathcal{F}}) - \inf R_P(f)}}$  Estimation Error $\geqslant 0$ \qquad $\underbrace{\hphantom{\inf R_P(f) - R_P^*}}$  Approximation Error $\geqslant 0$



Approximation Error: "How far is the model $\mathcal{F}$ from the truth?"

Estimation Error: "How good can we estimate the best predictor in $\mathcal{F}$?"

Harder if $\mathcal{F}$ is "large"

Bound on the Estimation Error:

$$R_P(\hat{f}_{\mathcal{F}}) - \inf_{f \in \mathcal{F}} R_P(f) = R_P(\hat{f}_{\mathcal{F}}) - R_P(f_{\mathcal{F}}^*)$$

$$= \left(R_P(\hat{f}_{\mathcal{F}}) - R_n(\hat{f}_{\mathcal{F}})\right)$$

$$\left[ + \left(R_n(\hat{f}_{\mathcal{F}}) - R_n(f_{\mathcal{F}}^*)\right)\right] \leq 0$$

$$+ \left(R_n(f_{\mathcal{F}}^*) - R_P(f_P^*)\right)$$

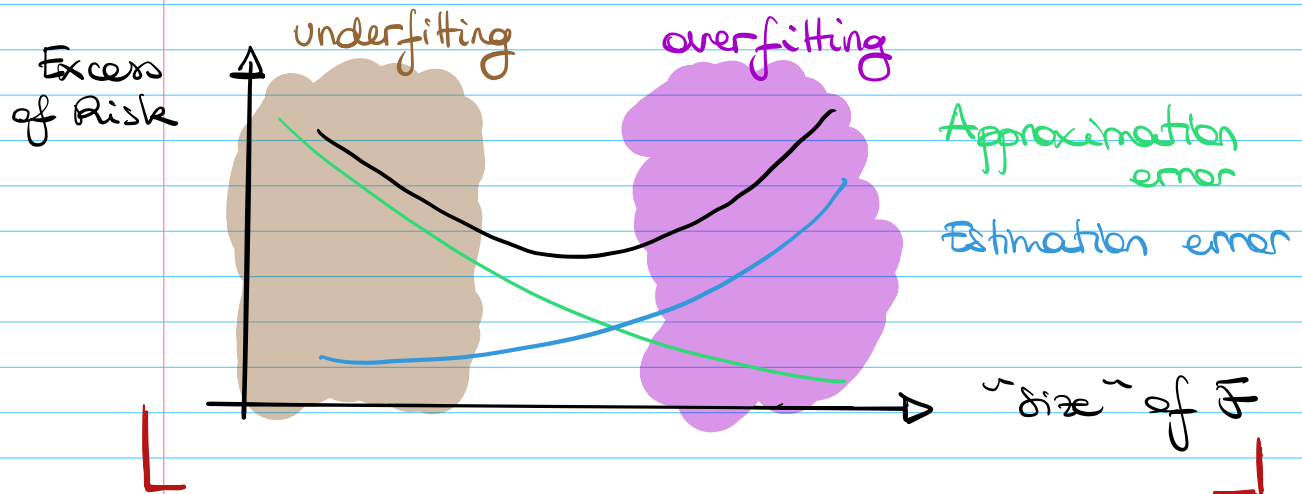$$\leq \sup_{f \in \mathcal{F}} \left(R_P(f) - R_n(f)\right) + \left(R_n(f_{\mathcal{F}}^*) - R_P(f_{\mathcal{F}}^*)\right)$$

uniform the deviation between empirical risk and P-risk.

↝ CENTRAL LIMIT THEOREM:

$$R_P(f) - R_n(f) = \frac{1}{n}\sum_{i=1}^{n}\left(E[\ell(f(x),y)] - \ell(f(x_i),y_i)\right)$$

$$\approx \frac{1}{\sqrt{n}} \times \text{Gaussian}$$

↳ But what about $\sup_{f \in \mathcal{F}}\left(R_P(f) - R_n(f)\right)$ ?

Excess
of Risk

underfitting

overfitting

Approximation
error

Estimation error

"size" of $\mathcal{F}$

③ Bound on the estimation error
in Binary Classification

$$y = \{-1, 1\} \qquad \ell(y, y') = \mathbb{1}\{y \neq y'\}$$

GOAL : Bound $\mathbb{E}\left[\sup_{f \in \mathcal{F}} (R_p(f) - R_n(f))\right]$

CASE 1 : Finite Number of Predictors

$$\mathcal{F} = \{f_1 \cdots f_k\}$$

## THEOREM (Maximal Inequality)

Let $z_1 \cdots z_k$ be random variables with

$$\forall d > 0, \quad \mathbb{E}[e^{d z_j}] \leq e^{d^2 \sigma^2 / 2} \qquad \text{(subgaussianity condition)}$$

Then,

$$\mathbb{E}\left[ \max_{j=1\cdots k} z_j \right] \leq \sigma \sqrt{2 \log k}$$

**proof:**

$$\max_{j=1\cdots k} e^{d a_j} \leq \sum_{j=1}^{k} e^{d a_j} \qquad (*)$$

$$\Rightarrow \mathbb{E}\left[ \max_{j=1\cdots k} z_j \right] = \frac{1}{d} \mathbb{E}\left[ \log\left( \max_{j=1\cdots k} e^{d z_j} \right) \right]$$
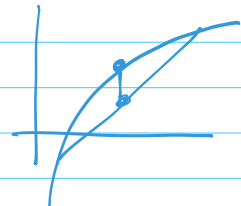
$$\frac{1}{d} \log(e^{d x}) = x$$

JENSEN
$$\leq \frac{1}{d} \log\left( \mathbb{E}\left[ \max_{j=1\cdots k} e^{d z_j} \right] \right)$$

$$\overset{(*)}{\leq} \frac{1}{d} \log\left( \mathbb{E}\left[ \sum_{j=1}^{k} e^{d z_j} \right] \right)$$

$$\underbrace{\sum_{j=1}^{k} \mathbb{E}[e^{d z_j}] \leq k e^{d^2 \sigma^2 / 2}}$$

$$\leq \frac{\log k}{d} + \frac{1}{d} \frac{d^2 \sigma^2}{2}$$

$\leadsto$ Choose $d = \frac{\sqrt{2 \log k}}{\sigma}$. ◼

Recall

$$\mathbb{E}\left[R_P(\hat{f}_{\mathcal{F}}) - \inf_{f \in \mathcal{F}} R_P(f)\right]$$

$$\leq \mathbb{E}\left[\sup_{f \in \mathcal{F}} \left(R_P(f) - R_n(f)\right)\right] + \mathbb{E}\left[\left(R_n(f_{\mathcal{F}}^*) - R_P(f_{\mathcal{F}}^*)\right)\right]$$

$$= \frac{1}{n}\sum_{i=1}^{n} \mathbb{1}\{f_{\mathcal{F}}^*(X_i) \neq Y_i\} - \mathbb{E}_P\left[\mathbb{1}\{f_{\mathcal{F}}^*(X) \neq Y\}\right]$$

$$\Rightarrow \mathbb{E}\left[\quad\right] = 0$$

See notes: $R_P(f) - R_n(f)$ is $\frac{1}{\sqrt{n}}$-subgaussian.
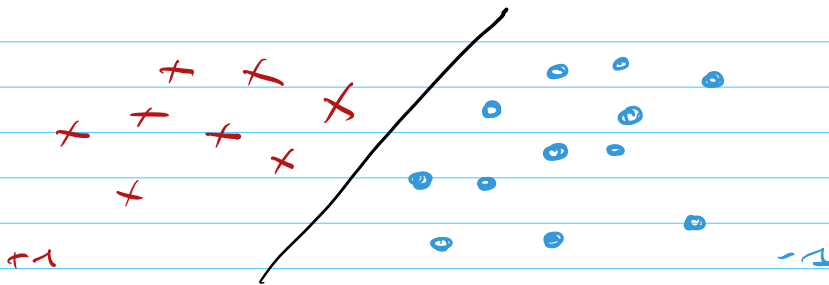
BOUND ON THE ESTIMATION ERROR

$$\Rightarrow \left[\mathbb{E}\left[R_P(\hat{f}_{\mathcal{F}}) - \inf_{f \in \mathcal{F}} R_P(f)\right] \leq \sqrt{\frac{2\log k}{n}}\right]$$

size of $\mathcal{F}$.

CASE 2: With VC-dimension

What if $\mathcal{F}$ is infinite?

ex: $\mathcal{F}_{lin} = \{$ hyperplane classifiers $\}$



Even if $\mathcal{F}$ is infinite, there is only a
finite number of classifications.

$$\underbrace{\mathcal{C}_{\mathcal{F}}(x_1 \cdots x_n)}_{} = \{(f(x_1), \ldots, f(x_n)): f \in \mathcal{F}\}$$
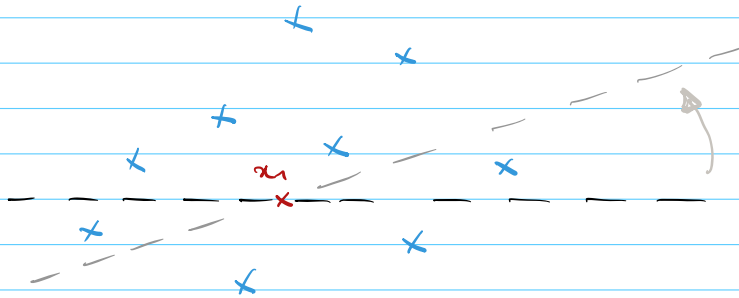$$\subseteq \{-1, +1\}^n$$

$\rightarrow$ of size $\mathcal{N}_{\mathcal{F}}(x_1, \ldots, x_n) \leq 2^n$

$\Rightarrow$ Actually the size of $\mathcal{F}$ can be
replaced by $\mathcal{N}_{\mathcal{F}}(x_1, \cdots x_n)$!

THEOREM:

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} (R_p(f) - R_n(f))\right] \leq 2 \mathbb{E}\left[\sqrt{\frac{2 \log \mathcal{N}_{\mathcal{F}}(x_1, \cdots x_n)}{n}}\right]$$

**Example:** $\mathcal{F} = \mathcal{F}_{ein}$ in $\mathbb{R}^2$



Rotation around $x_1$ : $\leq n$ different classifications
Rotation around $x_2$ : $\leq n$ different classifications
$\qquad \vdots \qquad\qquad\qquad\qquad\qquad \vdots$
Rotation around $x_n$ : $\leq n$ different classifications

$\Longrightarrow$ At most $n^2$ different classifications!

$$N_{\mathcal{F}}(x_1 \cdots x_n) \leq n^2$$

$\Longrightarrow$ $\mathbb{E}[\text{estimation error}] \leq 4\sqrt{\dfrac{\log n}{n}}$

And for a general $F$ ?

- We say that $F$ shatters $(x_1, \ldots, x_n)$
  if $\quad N_F(x_1, \ldots x_n) = 2^n$.

Vapnik - Chervonenkis

- The VC dimension of $F$ is the longest
$n$ such that there exists a configuration
$(x_1 \ldots x_n)$ of $n$ points being shattered by $F$.

$\leadsto$ $VC(F)$

For $n \leq VC(F)$, we cannot bound
$\qquad N_F(x_1, \ldots x_n)$ meaningfully ☹

What if $n > VC(F)$ ?

$\Downarrow$ MAGICAL RESULT

SAUER LEMMA: if $n > VC(F)$,

$$\log N_F(x_1, \ldots, x_n) \leq VC(F) \log\left(\frac{en}{VC(F)}\right)$$

$\leadsto$ if $n > 2 VC(F)$,
$\qquad N_F(\ldots) \lesssim n^{VC(F)} <<< 2^n \; !!$

**THEOREM:** if $n > 2 VC(\mathcal{F})$

$$\left[ E\left[ R_P(\hat{f}_{\mathcal{F}}) - \inf_{f \in \mathcal{F}} R_P(f) \right] \leq 2 \sqrt{\frac{2 VC(\mathcal{F})}{n} \log\left(\frac{en}{VC(\mathcal{F})}\right)} \right]$$

$\underbrace{\phantom{E\left[ R_P(\hat{f}_{\mathcal{F}}) - \inf_{f \in \mathcal{F}} R_P(f) \right]}}_{}$

Estimation error in binary classification $\lesssim \sqrt{\frac{VC(\mathcal{F})}{n}}$

# TAKE - HOME MESSAGES :

- The quality of a prediction is measured by a loss $\ell$.

$$P-Risk = \mathbb{E}_P[loss]$$

- The best theoretical predictor is the ~~Bayes~~ predictor.
    ⤳ cannot be computed

- A strategy to find a good predictor is to minimize the empirical risk on a model $\mathcal{F}$.

- The "size" of $\mathcal{F}$ has to be properly chosen to avoid

| UNDERFITTING | and | OVERFITTING |
|---|---|---|
| Large approximation error | | Large estimation error |
| ↳ The model $\mathcal{F}$ is too simple to capture the complexity of the dataset | | ↳ in binary classification, bounded by $\approx \sqrt{\dfrac{VC(\mathcal{F})}{P}}$ |
| | | Large if $\mathcal{F}$ is too complex. |