$$D = \{(\underline{x}_1, y_1), \dots (\underline{x}_N, y_N)\}$$

$$\underline{x} \in \mathbb{R}^d \qquad \hookrightarrow \text{i.i.d.}$$

$$y \in \begin{cases} \{0,1\} & \text{B.C.} \\ \{1,2,\dots,c\} & \text{M.C.C.} \\ \mathbb{R} & \text{REGR.} \end{cases}$$

← set of functions

L.M. : $f(\underline{x}) = \tilde{y}$ MAPS $\underline{x}$ INTO $\tilde{y}$, CHOOSEN IN $f \in F$

LOSS FUNCTION : $\ell(f(\underline{x}), y)$ HOW GOOD IS MY PREDICTION IN RESPECT TO MY TARGET

EMPIRICAL ERROR : $\hat{R}(f) = \frac{1}{m} \sum_{i=1}^{m} \ell(f(\underline{x}_i), y_i)$ ERROR THAT I COMPUTE ON MY DATA

(TRUE)
GENERALIZATION ERROR: $R(f) = \mathbb{E}_{\underline{x},y} \, \ell(f(\underline{x}), y)$ ERROR THAT I COMPUTE ON THE POPULATION, THE ONE TO MINIMIZE

CASE $|F| = 1$

$F = f \Rightarrow$ I HAVE ONLY 1 FUNCTION $\Rightarrow$ I HAVE NO LEARNING PROCESS

$$\hat{R}(f) = \frac{1}{m} \sum_{i=1}^{m} \ell_i$$

$$R(f) = \mathbb{E}_p \{\ell\}$$

VAR. OF RANDOM VAR $\ell$

$$P\{|R - \hat{R}| \geq \varepsilon\} \leq \frac{\sigma_p^2}{m \varepsilon^2} = \delta \quad \to \quad \varepsilon = \sqrt{\frac{\sigma_p^2}{m \delta}}$$

$$P\left\{|R - \hat{R}| \geq \sqrt{\frac{\sigma_p^2}{m \delta}}\right\} \leq \delta$$
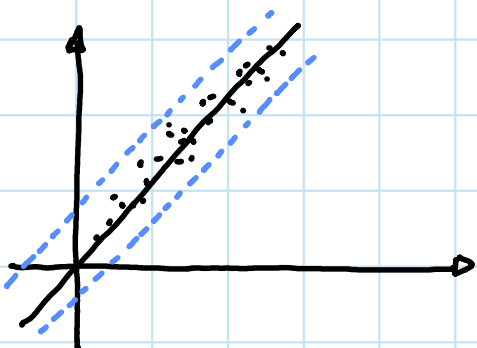
?: UNKNOWN

$$|R - \hat{R}|_{1-\delta} \leq \sqrt{\frac{\sigma_p^2}{m \delta}}$$

THE SMALLER IS $\delta \Rightarrow$ THE MORE PROBABLE THIS STATEMENT IS TRUE, THE LARGER IS THE INTERVAL BETWEEN R AND $\hat{R}$

$$\lim_{\delta \to 0} |R - \hat{R}|_{1-\delta} \leq \sqrt{\frac{\sigma_\rho^2}{m\delta}} = \infty$$

IF I HAVE A LOT OF SAMPLES $(m \to \infty)$ THE DISTANCE BECOMES SMALLER $(|R - \hat{R}| \to 0)$



IS $\ell \in [0, \infty)$ ? (UNBOUNDED)

<span style="color:blue">THE FAR I'M FROM THE PREDICTION, THE LARGER THE ERROR, BUT THEN AFTER A CERTAIN VALUE IS ONLY AN ERROR.</span>

$\hookrightarrow$ USUALLY THE LOSS IS BOUNDED, $\ell \in [0,1]$ IF I NORMALIZE. IN M.C.C. I JUST HAVE CORRECT OR NOT CLASSIFICATION.

IF $\ell \in [0,1] \Rightarrow \sigma_\rho^2 \leq 1$

$\Downarrow$

$$P\{|R - \hat{R}| \geq \varepsilon\} \leq \frac{\sigma_\rho^2}{m\varepsilon^2} \leq \frac{1}{m\varepsilon^2}$$

<span style="color:blue">$\delta$ CONFIDENCE</span>

<span style="color:orange">PROB. OF TRUE ERROR TO BE FAR AWAY FROM THE EMPIRICAL $\leq \frac{1}{m\varepsilon^2}$ ERROR MORE THAN $\varepsilon$</span>

$$|R - \hat{R}|_{1-\delta} \leq \sqrt{\frac{\sigma_\rho^2}{m\delta}} \leq \sqrt{\frac{1}{m\delta}}$$

<span style="color:red">EVERY ML ALGORITHM CAN BE RECONDUCTED TO THIS FORMULA. THERE'S RELATION BETWEEN THE MAX. DISTANCE FROM R AND $\hat{R}$ AND THE # OF SAMPLES COLLECTED.</span>

$\longrightarrow$ I EITHER HAVE MORE DATA $(m \to \infty)$ OR LESS CONFIDENCE $(\delta \to 0)$

$$P\{|R - \hat{R}| \geq \varepsilon\} \leq \frac{1}{m\varepsilon^2}$$   <span style="color:orange">CHERNOFF BOUND</span> $\longrightarrow$   HYPOTHESIS: i.i.d. $\ell \in [0,1]$

OTHER BOUNDS ARE ALSO USED:   $P\{|R - \hat{R}| \geq \varepsilon\} \leq e^{-2m\varepsilon^2}$   $\uparrow$   <span style="color:orange">HOEFFDIG BOUND</span>

$\boxed{\text{CASE } |F| = m_f}$

INDEPENDENT FROM
THE DATA $D_m$ !!

ASSUMPTIONS: $D_m$ is I.I.D, $\ell \in [0,1]$, $|F| = m_f$

$\Rightarrow \hat{f}$ IS LEARNED, FROM $D_m$, $F \longrightarrow \hat{f} \in F$

THE CHOSEN FUNCTION CAN
DEPEND FROM THE DATA $D_m$

THE FUNCTION SPACE $F$ IS CHOSEN BEFORE OBSERVING THE DATASET $\Rightarrow$ $F$ IS INDEPENDENT
FROM THE DATA $f \in F$ IS CHOSEN BY LOOKING AT THE DATA.

$\forall f \in F \quad P\{|R(f) - \hat{R}(f)| \geq \varepsilon\} \leq \dfrac{1}{m\varepsilon^2}$   BECAUSE EACH FUNCTION IS INDEPENDENTLY FROM THE DATASET.

WORST CASE SCENARIO: $\forall f \in F$ THE TRUE ERROR IS FAR AWAY FROM THE OBSERVED.

$\hat{f}: \quad P\{|R(\hat{f}) - \hat{R}(\hat{f})| \geq \varepsilon\} \leq \displaystyle\sum_{i=1}^{m_f} \dfrac{1}{m\varepsilon^2} = \dfrac{m_f}{m\varepsilon^2}$

$\rightarrow$ UNDERSTAND THE FORMULA

$\dfrac{m_f}{m\varepsilon^2} = \delta \Rightarrow \varepsilon = \sqrt{\dfrac{m_f}{m\delta}}$

BY PITAGORA

$|R(\hat{f}) - \hat{R}(\hat{f})|_{1-\delta} \leq \sqrt{\dfrac{m_f}{m\delta}} \leq \sqrt{\dfrac{m_f + 1 - 1}{m\delta}} \leq \sqrt{\dfrac{m_f - 1}{m\delta}} + \sqrt{\dfrac{1}{m\delta}}$

STATISTICAL PRICE:
THE MORE DATA I HAVE THE
LESS RISK I HAVE.

DEPENDS FROM THE LEARNING ($m_f$)

• IF $m_f = 1 \Rightarrow$ CLASSICAL STATISTICS $\rightarrow |R(\hat{f}) - \hat{R}(\hat{f})|_{1-\delta} \leq \sqrt{\dfrac{1}{m\delta}}$   (NO LEARNING)
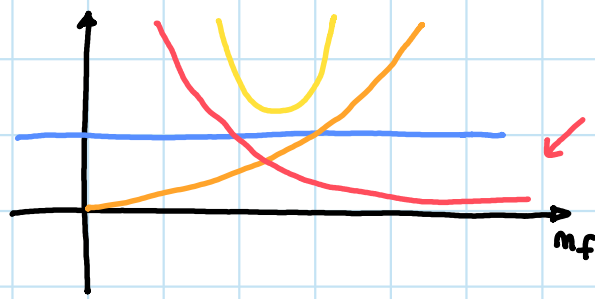
• IF $m_f > 1 \Rightarrow$ LEARNING PROCESS

$\boxed{D_m} \qquad \boxed{F} \qquad \hat{f} = \underset{f \in F}{\text{ARG MIN }} \hat{R}(f)$ — THE MINIMUM ERROR ON MY DATA

EMPIRICAL RISK MINIMIZATION (E.R.M.)

$\hat{R}(f)$ AND $m_f$ ARE RELATED: THE MORE FUNCTION I HAVE, BIGGER IS THE RISK, BUT IT'S

PROBABLE THAT MY $\hat{R}(f)$ WILL BE SMALLER.

$$|R(\hat{f}) - \hat{R}(\hat{f})| \leq \frac{1}{1-\delta}\sqrt{\frac{m_f-1}{m\delta}} + \sqrt{\frac{1}{m\delta}}$$
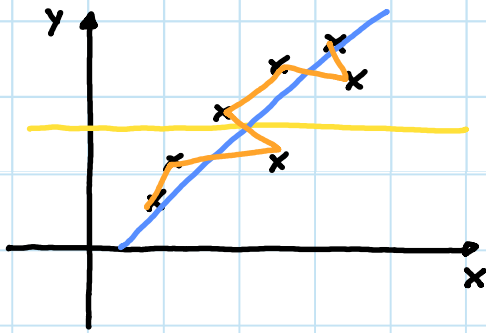
$$R(\hat{f}) \leq \hat{R}(\hat{f}) + \sqrt{\frac{m_f-1}{m\delta}} + \sqrt{\frac{1}{m\delta}}$$
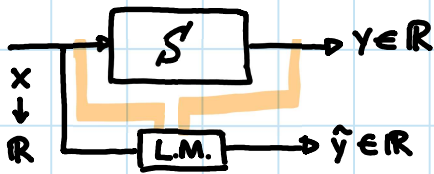


THE MORE FUNCTION I LEARN THE MORE PROB. A FUNC. WILL WORK WELL ON MY DATA

OCCAM'S RAZOR PRINCIPLE: GIVEN A PROBLEM IF YOU HAVE A LIST OF SOLUTIONS THAT MORE OR LESS BEHAVES THE SAME, THE ONE TO SELECT IS THE SIMPLEST ONE

IN M.L.: STRUCTURAL RISK MINIMIZATION



THE BLU IS THE CORRECT ONE



$$D_m = \{(x_1, y_1), \dots, (x_N, y_N)\} \quad i.i.d$$

$$f(x) = \sum_{i=0}^{P} a_i x^i \quad \text{HYPERPARAMETER}$$

PARAMETERS

ALSO THE FUNC. FORM IS AN HYPER PARAMETER (MCLAURIN, FOURIER...)

$$\ell(f(x), y) = (y - f(\underline{x}))^2 \quad \text{M.S.E.}$$

WHAT WE WANT TO FIND: $f^* = \text{ARG min } R(f)$
$$f \in F$$
$$a \in R^P$$

ORACLE

WITH $R(f) \leq \hat{R}(f) + C(m_f) + \sqrt{\frac{1}{m\delta}}$

?

COMPLEXITY OF MY MODEL

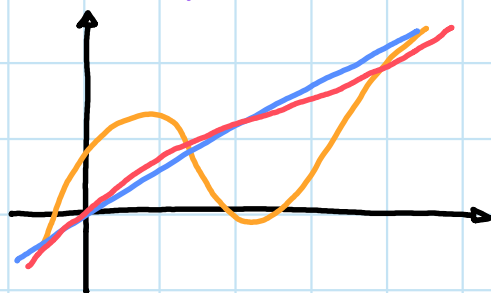WITH $R(f) \leq \hat{R}(f) + \lambda C(m_f) + \sqrt{\frac{1}{m\delta}}$

$\lambda$: REGULARIZATION PARAMETER (HYPERPARAMETER)

TRADES OFF

ACCURACY AND COMPLEXITY: HOW MUCH I WANT TO LEARN FROM MY DATA AND HOW MUCH I WANT MY MODEL COMPLEX

# HOW TO MEASURE COMPLEXITY

$C(f)$

$f(x) = \sum\limits_{i=0}^{p} a_i x^i$



• MORE COMPLEX THAN • (DEGREE OF THE POLY)

FOR • I USE THE FLACTUATION ($f''$)

$\dfrac{d^2 f}{dx^2}$ : NOT SO GOOD SINCE IS A FUNCTION $\Rightarrow$ I'll CALCULATE THE INTEGRAL IN THE DOMAIN (0 AND 1 SINCE I CAN ALWAYS NORMALIZE)

$\int\limits_{0}^{1} \left( \dfrac{d^2 f}{dx^2} \right)^2 dx$ : FIRST MEASURE OF COMPLEXITY

$\int\limits_{0}^{1} \left( \dfrac{d^2 f}{dx^2} \right)^2 dx$ WITH $f(x) = \sum\limits_{i=0}^{p} a_i x^i$

$\dfrac{df}{dx} = \sum\limits_{i=1}^{p} i \, a_i x^{i-1}$

↑ DERIVATIVE of $1 = 0$

$\dfrac{df^2}{dx^2} = \sum\limits_{i=2}^{p} i(i-1) a_i \cdot x^{i-2}$

$\left( \dfrac{df^2}{dx^2} \right)^2 = \sum\limits_{i=2}^{p} \sum\limits_{j=2}^{p} ij \, (i-1)(j-1) \, a_i a_j \, x^{i-2} x^{j-2}$

$\int\limits_{0}^{1} \left( \dfrac{df^2}{dx^2} \right)^2 = \sum\limits_{i=2}^{p} \sum\limits_{j=2}^{p} ij (i-1)(j-1) a_i a_j \int\limits_{0}^{1} x^{i+j-4} dx = \sum\limits_{i=2}^{p} \sum\limits_{j=2}^{p} ij(i-1)(j-1) a_i a_j \left. \dfrac{x^{i+j-3}}{i+j-3} \right|_{0}^{1} =$

$= \sum\limits_{i=2}^{p} \sum\limits_{j=2}^{p} \dfrac{(i-1)(j-1)ij}{i+j-3} a_i a_j$

WHAT WE WANT TO MINIMIZE?

$\overbrace{\qquad\qquad}^{a^T M a}$

$\hat{f} = \min\limits_{\substack{f \in F \\ a \in \mathbb{R}^d}} \hat{R}(f) + \lambda C(f) \Rightarrow \min\limits_{\underline{a}} \dfrac{1}{m} \sum\limits_{i=1}^{m} \Big( \underbrace{\sum\limits_{j=0}^{p} a_j x_i^j}_{f(x_i)} - y_i \Big)^2 + \lambda \sum\limits_{i=2}^{p} \sum\limits_{j=2}^{p} \dfrac{ij(i-1)(j-1)}{i+j-3} a_i a_j$

REWRITE THIS FORMULA INTO MATRICIAL FORM

$$\underline{a} = \begin{bmatrix} a_0 \\ \vdots \\ a_p \end{bmatrix} \in \mathbb{R}^{p+1} \qquad \underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m \qquad X = \begin{bmatrix} x_1^0 & \cdots & x_1^p \\ & \vdots & \\ x_N^0 & \cdots & x_N^p \end{bmatrix} \in \mathbb{R}^{m \times (p+1)}$$

$$M = \begin{cases} M_{ij} = 0 & i,j < 2 \\ M_{ij} = \dfrac{ij\,(i-1)(j-1)}{i+j-3} & i,j \geqslant 2 \end{cases}$$

 $\in \mathbb{R}^{(p+1) \times (p+1)}$

$$\min_{\underline{a}} \; \frac{1}{m} \| X\underline{a} - \underline{y} \|_2^2 + \lambda \underline{a}' M \underline{a}$$

$$\min_{\underline{a}} \; \underbrace{\| X\underline{a} - \underline{y} \|_2^2}_{R(f)} + \overbrace{\lambda \underline{a}' M \underline{a}}^{C(f)}$$

PARABOLOID

$\hookrightarrow \; I :$ RIDGE REGRESSION (REGULARIZED RIDGE SQUARE)

TO FIND THE MINIMUM I COMPUTE THE GRADIENT:

$$\nabla_{\underline{a}} \left( \underline{a}' X' X \underline{a} - 2\underline{a}' X' \underline{y} + \underline{y}' \underline{y} + \lambda \underline{a}' M \underline{a} \right) = \underline{0}$$

$$2 X' X \underline{a} - 2 X' \underline{y} + 2 M \underline{a} = 0$$

$$\boxed{(X'X + \lambda M)\,\underline{a} = X'\underline{y}} \qquad \text{LINEAR SYSTEM, GAUSS JORDAN } O(m^2)$$