

DEFINITIONS

MEANING

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

DATASET, THE ONLY KNOWLEDGE THAT WE HAVE ABOUT THE SYSTEM (i.i.d.)

$$f(x) = \sum_{i=0}^p a_i x^i \in F$$

CLASS OF FUNCTIONS THAT WE USE FOR APPROXIMATING THE INPUT-OUTPUT RELATION

$$\ell(f(x), y) = (y - f(x))^2$$

MEASURE OF QUALITY FOR F . HOW MUCH THE THINGS THE f PREDICTS ARE CLOSE TO THE THINGS I MEASURE. **LOSS FUNCTION**

$$R(f) = \mathbb{E}_{x,y} (\ell(f(x), y))$$

RISK OF A FUNCTION. QUANTITY THAT WE WANT TO MINIMIZE. IT'S A CHOICE!

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

EMPIRICAL ERROR ON MY DATA, I CAN MEASURE.

IF D i.i.d. AND $\ell \in [a, b]$ (BOUNDED) THEN:

- IF $|F| = 1$ AND INDEPENDENT FROM $D \Rightarrow |R(f) - \hat{R}(f)|_{1-s} \leq \sqrt{\frac{1}{ms}}$

- IF $|F| = m_f$ AND INDEPENDENT FROM $D \Rightarrow |R(f) - \hat{R}(f)|_{1-s} \leq \sqrt{\frac{1}{ms}} + \sqrt{\frac{m_f - 1}{ms}}$

$$\min_{f \in F} R(f) \rightarrow \min_{f \in F} \hat{R}(f) + \lambda C(f)$$

TRADE OFF BETW. ACCURACY AND COMPLEXITY.

THE SQUARE LOSS IS:

- DIFFERENTIABLE
- CONVEX (good computational properties)

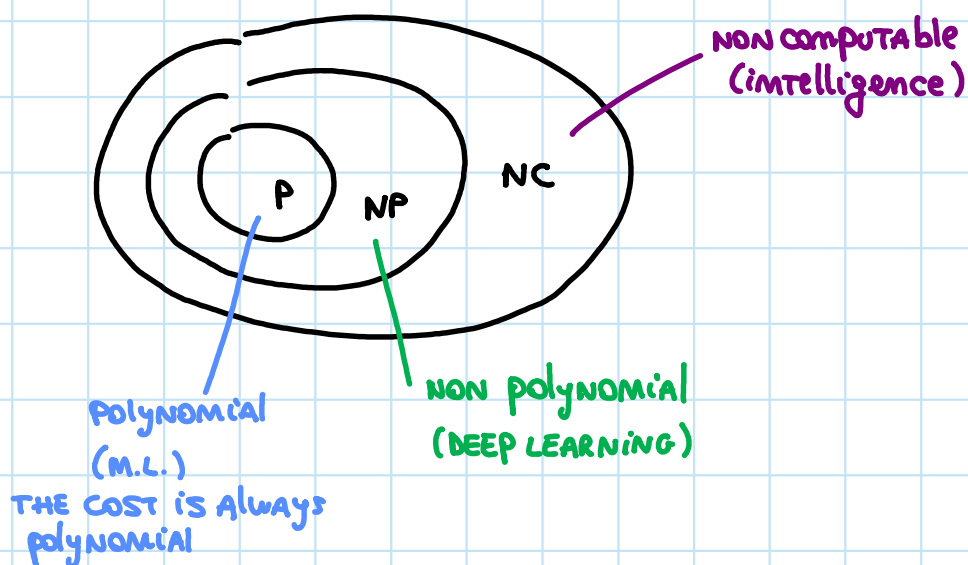
UNKNOWN, ERROR ON THE POPULATION GENERALIZATION ERROR

NO LEARNING, ONLY A STATISTICAL RISK (LIMITED)

THE MORE FUNC. I USE TO CHECK IF A FUNC. IS GOOD ON THE DATA, THE MORE THE RISK OF FINDING

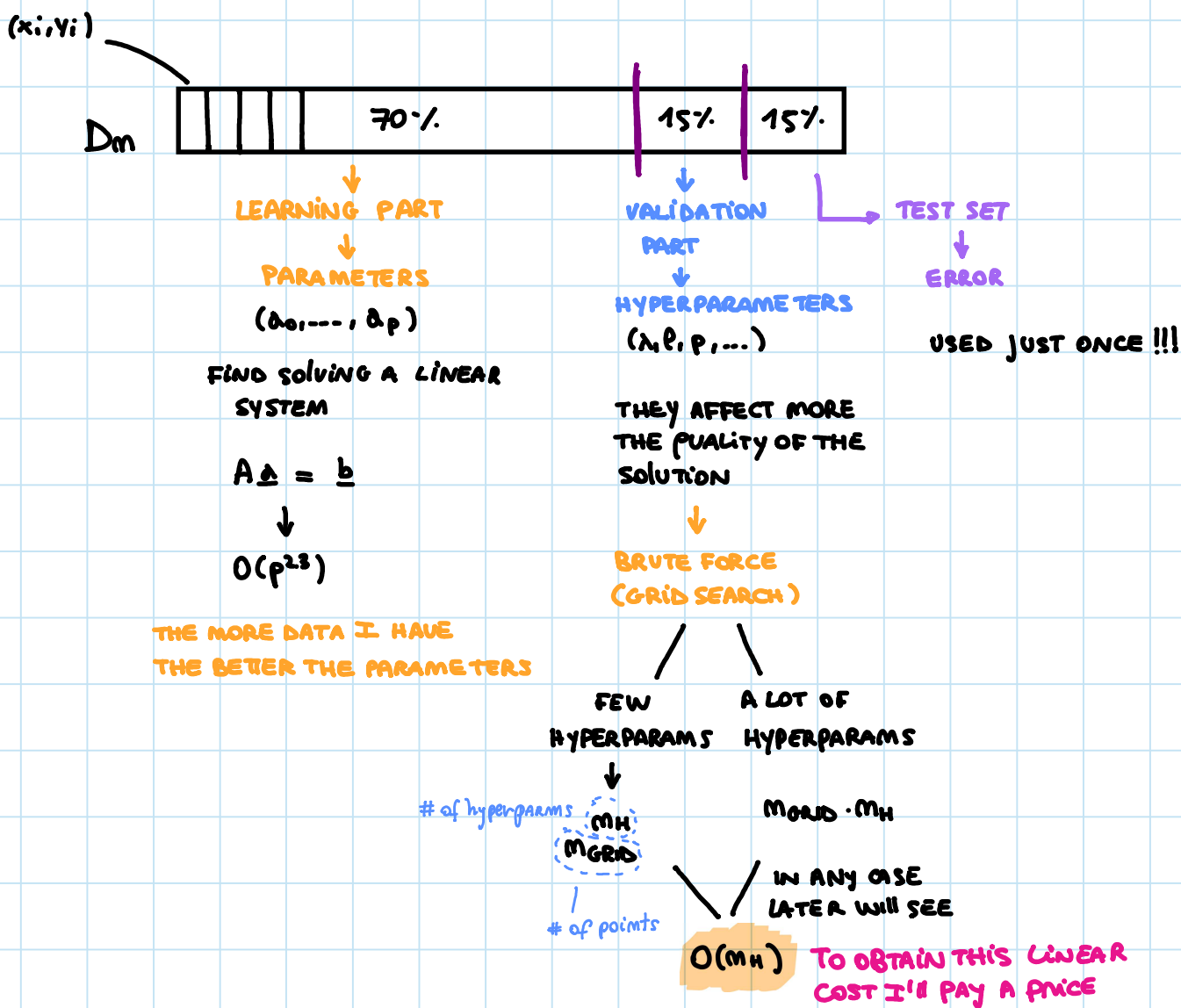
A FUNC. THAT RANDOMLY FITS MY DATA INCREASES \Rightarrow I DON'T HAVE ANY KNOWLEDGE OF THE DISTRIBUTION.

TYPES OF PROBLEMS:



WE ARE TRYING TO CONSTRUCT AN INTELLIGENT MACHINE WITH A POLYNOMIAL COST.

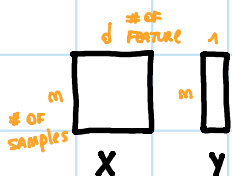
WE CANNOT GET THE CHOICES A PRIORI BUT WE NEED TO TEST THEM ON THE DATA.



IN THE LEARNING PART THE PROBLEM IS CONSTRAINED, BUT IN THE VALIDATION THE PROBLEM BURNS AS NON POLYNOMIAL. TO MAKE IT LINEAR I'LL HAVE TO PAY A PRICE.

ALL OF THIS WAS FOR JUST ONE OUTPUT!

$$x \in \mathbb{R} \rightarrow x \in \mathbb{R}^d$$



LET'S DO THE CHOICES NOW:

$$f(x) = a_0 \quad \text{FIRST element x vector} \quad \text{FLAT MODEL } O(1)$$

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1d}x_d \quad \text{LINEAR MODEL } O(d)$$

$$a_{11}x_1^2 + \dots + a_{1d}x_d^2 + a_{2d+1}x_1x_2 + \dots \quad \text{QUADRATIC } O(d + \binom{d}{2}) \approx O(d^2)$$

$$\vdots$$

$$\text{DEGREE } p : O\left(\binom{d}{p}\right) \approx O(d^p) \quad \text{MC. LAURIN SERIES IN DIMENSION } p$$

$$d = 10^{12}, p = 3 \Rightarrow (10^{12})^3 = (2^{36})^3 \approx 2^{108} \quad \text{JUST FOR STORING MY PARAMS}$$

FIRST APPROACH DOESN'T WORK

LET'S DO A LINEAR MODEL (# OF PARAMS SCALES LINEARLY WITH # FEATURES)

$$f(x) = \underline{w} x \quad \text{EASIEST THING THAT I CAN DO}$$

• CONSTRAINT: $f(x)$ PASSES IN 0 $\Rightarrow b = 0$

$$\ell(f(x), y) = (y - f(x))^2$$

$$C(f) = ? \quad \text{BEFORE WAS } \underline{w}' M \underline{w} \quad \Rightarrow C(f) = \|\underline{w}\|^2$$

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ \underline{w}' & I & \underline{w} \end{matrix} = \|\underline{w}\|^2$$

MY MODEL

$$\min_{f \in F} \hat{R}(f) + \lambda C(f) = \min_{f \in F} \frac{1}{m} \sum_{i=1}^m (\underline{w} x_i - y_i)^2 + \lambda \|\underline{w}\|^2$$

MATRICAL FORMAT:

$$X = \begin{bmatrix} \underline{x}_1' \\ \vdots \\ \underline{x}_m' \end{bmatrix} \in \mathbb{R}^{m \times d} \quad \underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

$$\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|^2$$

RIDGE REGRESSION (RLS)

PARABOLOID \rightarrow CONVEX
problem \rightarrow 1 solution

SOLUTION:

$$\nabla_{\underline{w}} (\underline{w}' \underline{X}' \underline{X} \underline{w} - 2 \underline{w}' \underline{X}' \underline{y} + \underline{y}' \underline{y} + \lambda \underline{w}' \underline{I} \underline{w}) = \underline{0}$$

$$2 \underline{X}' \underline{X} \underline{w} - 2 \underline{X}' \underline{y} + 2 \lambda \underline{I} \underline{w} = \underline{0}$$

$$(\underline{X}' \underline{X} + \lambda \underline{I}) \underline{w} = \underline{X}' \underline{y}$$

LINEAR SYSTEM $O(d^2)$

CAN BE GOOD OR BAD DEPENDING ON
THE SITUATION

$$\begin{matrix} A \\ \in \mathbb{R}^{d \times d} \end{matrix} \quad \begin{matrix} \underline{b} \\ \in \mathbb{R}^d \end{matrix}$$

$$\text{CAN BE ALSO REWRITTEN IN: } \underline{w} = (\underline{X}' \underline{X} + \lambda \underline{I})^+ \underline{X}' \underline{y} \quad O(d^2)$$

MERCER DISCOVERED KERNELS:

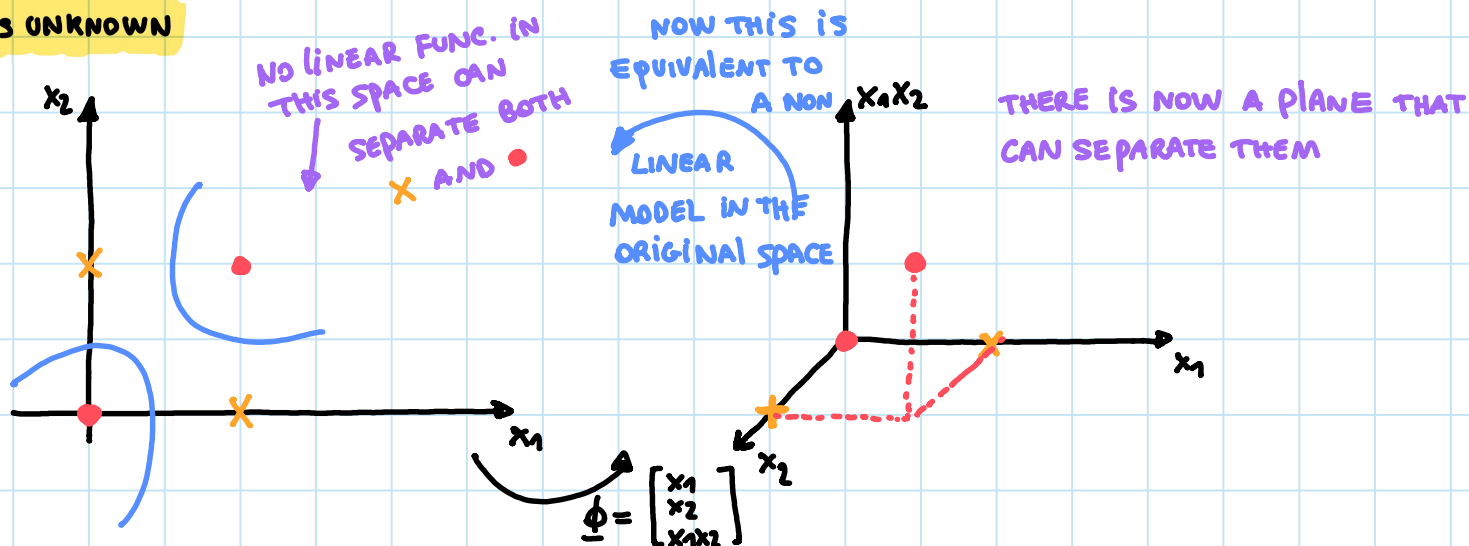
$$k(\underline{x}_i, \underline{x}_j) = e^{-\gamma \|\underline{x}_i - \underline{x}_j\|^2}$$

WORKS IN THE ORIGINAL SPACE BUT EQUIVALENT TO THIS:

$$k(\underline{x}_i, \underline{x}_j) = \underline{\phi}'(\underline{x}_i) \underline{\phi}(\underline{x}_j)$$

PROJECT MY DATA IN A NEW SPACE INDUCED BY ϕ AND COMPUTING THE
SCALAR PRODUCT IN THIS SPACE

ϕ IS UNKNOWN



WITH KERNELS I CAN DO SPECIFIC COMPUTATION WITHOUT COMPUTING ϕ .

$$e^{-\gamma \|x_i - x_j\|^2}$$

WHAT IS $|\phi| = ?$ IN HOW MANY DIMENSIONS THIS KERNEL IS PROJECTING MY DATA

CASE $d=1$

$$e^{-\gamma \|u-v\|^2} = e^{-\gamma u^2} e^{-\gamma v^2} e^{2\gamma uv} = e^{-\gamma u^2} e^{-\gamma v^2} \sum_{i=0}^{\infty} \frac{(2\gamma uv)^i}{i!} =$$

$$e^x = \sum_{i=0}^{\infty} \frac{1}{i!} \left. \frac{d^i e^x}{dx^i} \right|_{x=0} x^i = \sum_{i=0}^{\infty} \frac{x^i}{i!}$$

MACLAURIN SERIES

SO $|\phi| = \infty \Rightarrow$ I CAN ALWAYS LINEARIZE

$$\begin{aligned} e^{-\gamma u^2} \begin{bmatrix} 1 \\ \sqrt{2\gamma} u \\ \frac{2\gamma}{\sqrt{2}} u^2 \\ \frac{\sqrt{(2\gamma)^3}}{\sqrt{3}} u^3 \\ \vdots \end{bmatrix}^T & e^{-\gamma v^2} \begin{bmatrix} 1 \\ \sqrt{2\gamma} v \\ \frac{2\gamma}{\sqrt{2}} v^2 \\ \frac{\sqrt{(2\gamma)^3}}{\sqrt{3}} v^3 \\ \vdots \end{bmatrix} \\ \phi^T(u) & \phi(v) \end{aligned}$$

IN RLS:

$$\min_{\underline{w}} \|X\underline{w} - y\|^2 + \lambda \|\underline{w}\|^2$$

$$\underline{w} = (X'X + \lambda I)^+ X'y \quad (d^2)$$

WHAT IS THE SHAPE OF THIS SOLUTION?

I SUPPOSE THAT THE SOLUTION IS LINEAR IN THE DATA THAT I HAVE (THE SOLUTION DEPENDS ON THE DATA I COLLECTED)

$$\underbrace{\underline{w}}_{d \times 1} = \sum_{i=1}^n \underbrace{a_i}_{d \times 1} \underbrace{x_i}_{m \times 1} = \underbrace{X'}_{d \times m} \underbrace{\underline{\alpha}}_{m \times 1} \Rightarrow \text{THE SOLUTION LIVES IN THE SPACE INDUCED BY THE MATRIX } X$$

FOR THE POINTS THAT ARE PERPENDICULAR TO THE SPACE GENERATED BY X

$$X\alpha = 0$$



my hypothesis is that my solution is parallel to the space induced by the vectors inside X