

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \text{i.i.d. Hypothesis}$$

MORE DATA I HAVE, MORE INFORMATION ABOUT THE SYSTEM

SYSTEM DOESN'T CHANGE IN TIME.

$$f(x) = \tilde{y} \in F$$

PARAMETERS w, x
HYPERPARAMETERS $w, \phi(x)$

HERE WE SEARCH FOR A GOOD QUALITY OF THE SOLUTION.

LOSS FUNCTION: $\ell(f(x), y)$ (HYPERPARAMETER)

COMPLEXITY MEASURE: $C(f)$ (HYPERPARAMETER)

HERE WE SEARCH FOR THE SIMPLICITY OF THE SOLUTION

WE WANT THE MODEL TO BE GOOD ON PREVIOUSLY UNSEEN DATA, SO WE WANT TO MINIMIZE THE RISK OF OUR FUNCTION.

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

I JUST CARE TO BE MORE GOOD ON THE DATA THAT IS MORE FREQUENT (HYPER PARAMETER)

WE CANNOT MEASURE THIS RISK, I ONLY HAVE THE DATA, SO I CAN MEASURE THE RISK ON MY DATA.

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

USING STATISTICALLY LEARNING THEORY: WHEN F IS INDP. FROM D_m

TRADE OFF: COMPLEX FUNC \rightarrow REDUCE RISK ON MY DATA, SIMPLER FUNC \rightarrow MORE RISK ON MY DATA

$$R(f) \leq \hat{R}(f) + \lambda C(f) + \Phi(m, \delta)$$

CHOICE RISK (DUE TO THE LEARNING PHASE)

DATA RISK

I NEED TO BALANCE THIS TWO NEEDS.

I SELECT AS BEST FUNCTION:

$$f^* = \underset{f \in F}{\text{ARG min}} \hat{R}(f) + \lambda C(f)$$

HOW TO SOLVE THIS MINIMA?

IF: $\begin{cases} f(x) = \underline{w}x \\ e = (y - f)^2 \\ C = \|\underline{w}\|^2 \end{cases} \rightarrow \begin{array}{l} \text{THE SOLUTION IS } O(d^2) \text{ OF } A\underline{w} = b \\ \text{THE SOLUTION IS } O(m^2) \text{ OF } A\underline{x} = b \end{array} \quad \left. \vphantom{\begin{matrix} f(x) = \underline{w}x \\ e = (y - f)^2 \\ C = \|\underline{w}\|^2 \end{matrix}} \right\} \text{OPTIMAL SOLUTION}$

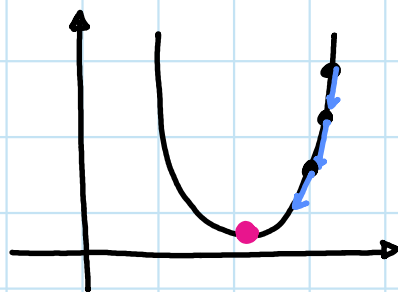
$\min_{\underline{w}} \|\underline{x}\underline{w} - \underline{y}\|^2 + \lambda \|\underline{w}\|^2$ THIS FORMULATION IS GOOD BECAUSE I DON'T HAVE CONSTRAINTS. IT'S THE MINIMA OF A PARABOLOID.

$$\begin{cases} \min M(\underline{x}) \\ \text{s.t.} \\ I(\underline{x}) \leq 0, E(\underline{x}) = 0 \end{cases}$$

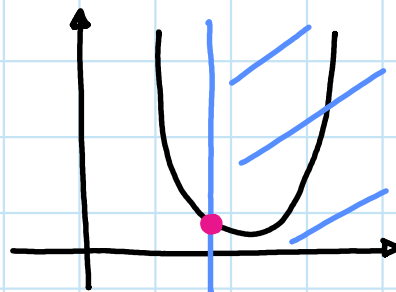
INEQUALITY AND EQUALITY CONSTRAINTS

I WOULD LIKE TO DEAL ALWAYS WITH NON CONSTRAINED OPTIMIZATION PROBLEMS.

$$\min_{\underline{x}} M(\underline{x})$$



UNCONSTRAINED
PROBLEM



CONSTRAINED
PROBLEM

GRADIENT DESCENT:

$$\begin{cases} \underline{x}_0 = \underline{0} \\ \underline{x}_{i+1} = \underline{x}_i - \eta \nabla_{\underline{x}} M(\underline{x})|_{\underline{x}_i} \end{cases}$$

NEWTON METHOD

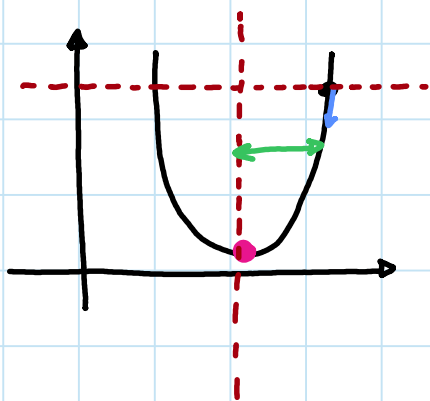
LEARNING RATE
(HYPERPARAMETER)

WHY USE GRADIENT DESCENT AND NOT THE METHOD ABOVE, SOLVING THE LINEAR SYSTEM?

IF MY FUNCTION $M(\underline{x})$ IS L SMOOTH, SO IF $\forall \underline{x}$ IT'S BOUNDED BY A PARABOLOID

$$M(\underline{x}') \leq M(\underline{x}) + \nabla_{\underline{x}} M(\underline{x}) \cdot (\underline{x}' - \underline{x}) + \frac{L}{2} \|\underline{x}' - \underline{x}\|^2$$

THERE'S A PARABOLA SMALL ENOUGH THAT CAN BE INSERTED IN EACH POINT.



I KNOW WHICH IS THE MAX VALUE OF η SO THAT I DON'T START TO GO IN THE OPPOSITE DESIDERATE DIRECTION.

$$\eta_{\max} = \frac{1}{L}$$

AND WE CAN PROVE THAT:

$$M(\underline{x}_t) - M(\underline{x}^*) \leq \frac{L}{t} \|\underline{x}^*\|^2$$

REACHED
POINT t

OPTIMAL
POINT

THE DISTANCE BETWEEN THE POINT THAT I AM AND THE POINT THAT I WANT TO REACH, GOES DOWN LINEARLY IN THE NUMBER OF ITERATIONS.

↳ IF I WANT TO BE 2 TIMES CLOSER, I HAVE TO DOUBLE THE NUMBER OF ITERATIONS.

LINEAR SYSTEM

$O(m^2)$

$$Ax = b$$

↳ $x'x$ or xx'

$$X = \begin{bmatrix} x_1' \\ \vdots \\ x_n' \end{bmatrix} \in \mathbb{R}^{n \times d}$$

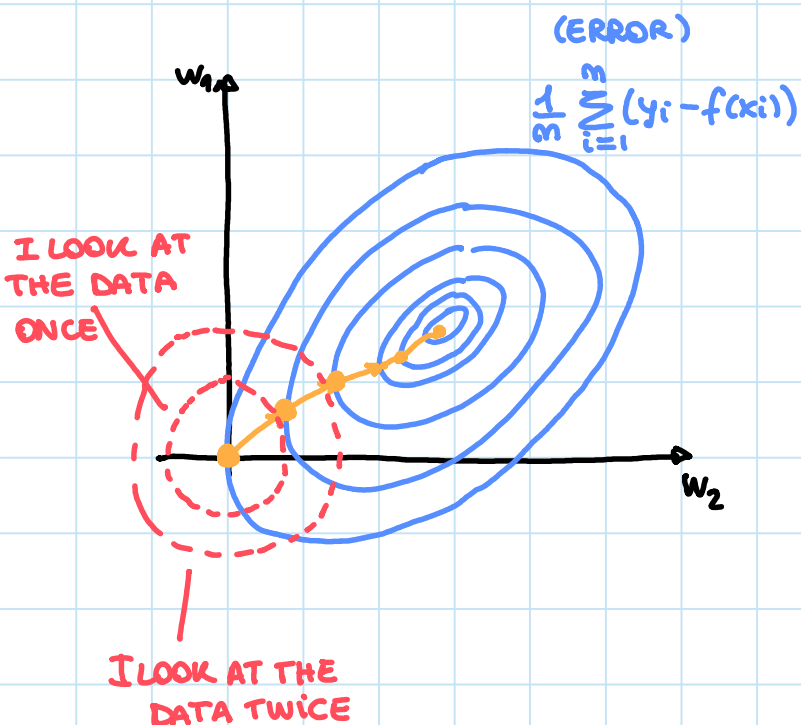
I NEED TO HAVE ALL THE DATA ON MY P.C.

A ALSO NEED TO BE SEMIDEFINITE POSITIVE



NOT VERY PARALLELIZABLE

VERY TIME CONSUMING!



GRADIENT DESCENT

$1/t$

$$x_{i+1} = x_i - \eta \nabla_x M(x) \Big|_{x_i}$$



$$\min_{\underline{w}} \frac{1}{m} \sum_{i=1}^m (y_i - \underline{w} \phi(\underline{x}_i))^2 + \lambda \|\underline{w}\|^2$$

I APPLY HERE THE GRADIENT:

$$\frac{1}{m} \sum_{i=1}^m \nabla_{\underline{w}} (y_i - \underline{w} \phi(\underline{x}_i))^2 + 2\lambda \underline{w}$$

THIS IS PARALLELIZABLE!!

OSS: EVERY STEP THAT I DO I'M RE-

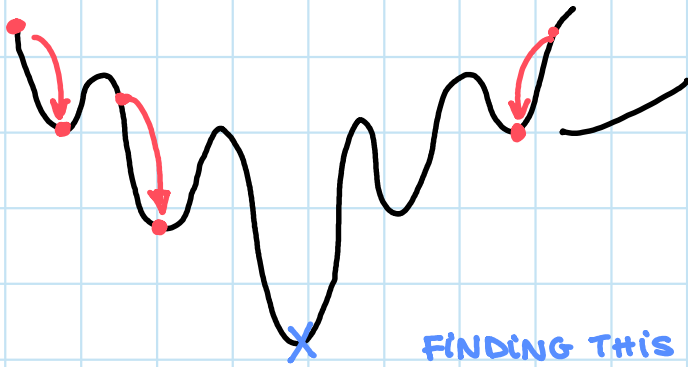
LOOKING AT MY DATA → I'M OVERFITTING MY DATA

THE NUMBER OF STEPS ARE EQUIVALENT

TO $\|\underline{w}\|^2 \rightarrow$ IMPLICIT REGULARIZER

⇒ RIDGE REGRESSION IF I DON'T REACH THE MINIMA.

IF MY PROBLEM IS NOT CONVEX:



FROM THE OPTIMIZATION POINT OF VIEW THEY ARE WORSE, BUT FROM THE PRACTICAL THEY COULD BE BETTER (BEST TRADE OFF BETWEEN SIMPLICITY AND QUALITY)

FINDING THIS POINT IS
NON POLYNOMIAL

MINIMA THAT ARE EASIER TO FIND ARE LESS PROBABLE TO OVERFIT.

NUMBER OF ITERATION IS A COMPLEXITY MEASURE (IN GRADIENT DESCENT)