

# OPTIMIZATION FOR MACHINE LEARNING

September 16, 2024

Today:

- \* Logistics
- \* Introduction to optimization

# LOGISTICS

Course webpage: <https://www.lamsade.dauphine.fr/~croyer/teachOID.html>

My email: clement.royer@lamsade.dauphine.fr

Schedule:

Sep 16	8.30 - 11.45
Sep 23	8.30 - 11.45
Sep 24	8.30 - 11.45
Sep 30	8.30 - 11.45
Oct 1	17.15 - 20.30
Dec 16	13.45 - 17
Dec 17	17.15 - 20.30
Dec 20	13.45 - 17

## Grading

60% Exam (January 29)

40% Homework (done in pair in Oct 1st / Dec 17 sessions)

# INTRODUCTION TO OPTIMIZATION

Optimization: Field of study concerned with making the best decision out of a set of alternatives

→ Flattened in the 1980s as an area of computational mathematics

⇒ Success stories: Finance, Supply chain, Scheduling, Vaccine dispatch, ...

↳ Shift over the past 20 years: data-driven optimization

• changed the problems people are interested in  
⇒ Specific formulations

• changed the algorithms people develop and use  
⇒ Stochastic gradient methods are the most widely used algorithms in ML

⇒ Other methods that were highly popular in 1980s-2000s fell out of fashion (Newton's method)

## Typical "optimization for machine learning problem"

ML (Machine Learning): Extract information from data

↳ Data is typically available under the form of samples

$$D = \left\{ (x_i, y_i) \right\}_{i=1..m} \quad m \text{ samples}$$

input      output

↳ Goal: Learn a mapping between the  $x_i$ 's and the  $y_i$ 's

$\Rightarrow$  We assume that this mapping is defined (or parameterized) using a vector  $w \in \mathbb{R}^d$  of "weights" or parameters

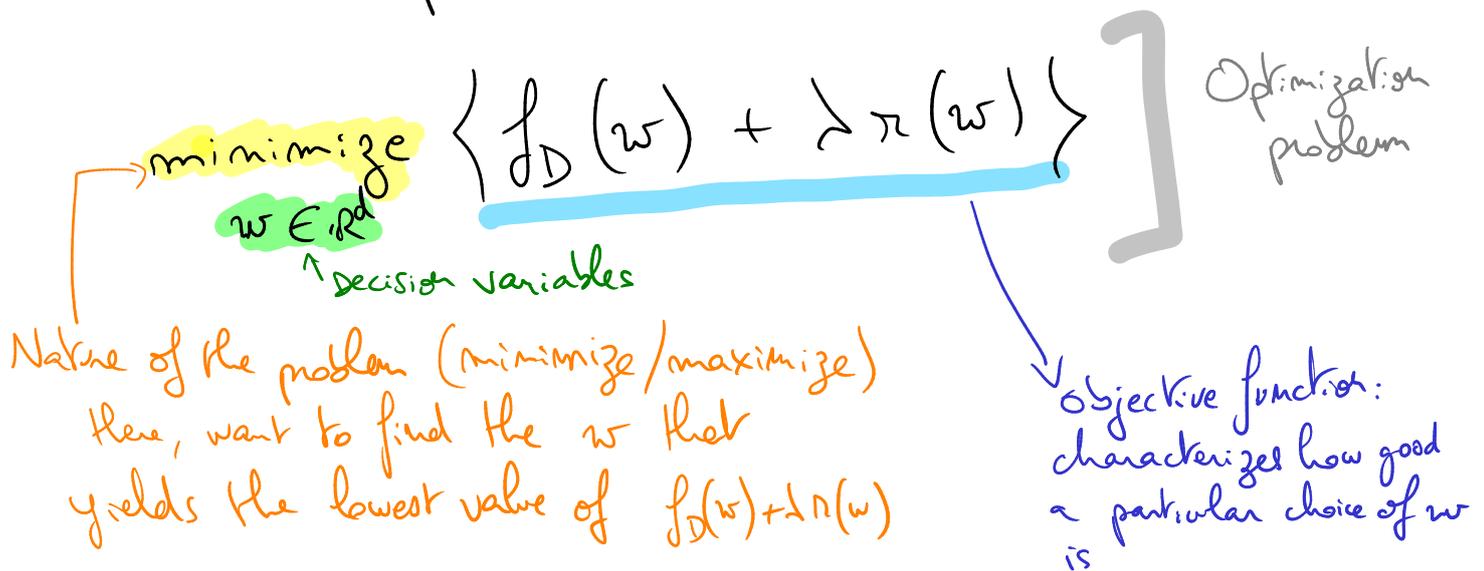
$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \begin{matrix} \uparrow \\ \downarrow \end{matrix} \quad \begin{matrix} d \\ d \end{matrix}$$

$w_j \in \mathbb{R} \forall j$

$\Rightarrow$  Notation:  $h(\cdot; w): x \mapsto h(x; w)$

Learning problem / Training problem

Find the best parameter values given the data  
We represent this problem as follows:



Decision variables: Parameters we can choose to affect the objective function

NB: The optimization problem reads  
"minimize  $f_D(w) + \lambda r(w)$  over  $w \in \mathbb{R}^d$ "

Specific objective function structure in ML:  $f_D(w) + \lambda r(w)$

$\bullet$   $f_D(w)$ : training loss / empirical loss / data-fitting term  
 $\Rightarrow$  Depends on the data  $D$ , measures how good the choice of  $w$  is at learning / fitting relationships between

inputs and outputs

⇒ Generic form:

$$D = \{(x_i, y_i)\}_{i=1..m}$$

Input data (blue arrow to  $x_i$ )  
Output data (green arrow to  $y_i$ )

$$J_D(w) = \frac{1}{m} \sum_{i=1}^m \ell(h(x_i; w), y_i)$$

Loss function (pink arrow to  $\ell$ )

Average/mean over all data points

Prediction of  $y_i$  given  $x_i$  made using the mapping  $h(\cdot; w)$  parameterized by  $w$

$\ell$ : loss function quantifies how close  $h(x_i; w)$  and  $y_i$  are

$\lambda(w)$ : regularization term, penalizes choices of  $w$  that do not satisfy desired properties ⇒ Not data-dependent!

$\lambda > 0$ : tradeoff parameter

$\lambda \gg 1$ : more weight on the regularization term, less on data

$\lambda \ll 1$ : more weight on the data-fitting term, less on regularization

## Examples

### ① Linear least squares / Linear regression

Data:  $\{(x_i, y_i)\}_{i=1..m}$      $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$

⇒ Representation:  $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}$

matrices with  $m$  rows and  $d$  columns (orange arrow pointing to  $X$ )

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

$$x_i = \begin{bmatrix} \uparrow 1 \\ \vdots \\ \uparrow d \end{bmatrix}$$

$$x_i^T = \begin{bmatrix} \leftarrow d \end{bmatrix} \uparrow 1$$

Goal: Find a linear relationship between  $x_i$ 's and  $y_i$ 's

Optimization problem:

minimize  $\frac{1}{2m} \|Xw - y\|^2$   
 $w \in \mathbb{R}^d$

Norm  $\|v\|^2 = \sum_{i=1}^m v_i^2 \quad \forall v \in \mathbb{R}^m$

where  $\|Xw - y\|^2 = \sum_{i=1}^m (x_i^T w - y_i)^2$

loss function

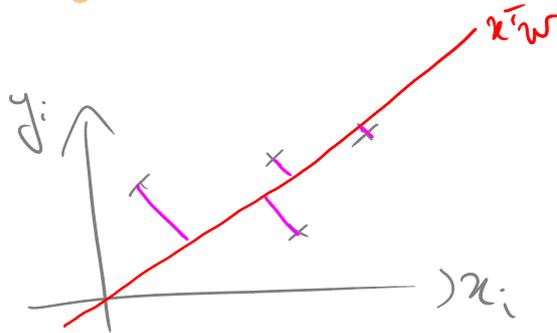
↳ Fits the general framework with  $r(w) = 0 \quad \forall w$

and

$J_0(w) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (x_i^T w - y_i)^2$

flapping  $x \mapsto x^T w = \sum_{j=1}^d x_j w_j$   
 "scalar product"  
 / "inner product"

NB:  $x^T =$  "x transpose"



↳ Variants with regularization

\* Ridge regression

minimize  $w \in \mathbb{R}^d$

$\frac{1}{2m} \|Xw - y\|^2 + \lambda \underbrace{\sum_{j=1}^d w_j^2}_{\|w\|^2}$

$\lambda > 0$   
 $l_2$  regularization  
 (reduces overfitting)

\* LASSO

minimize  $w \in \mathbb{R}^d$   $\frac{1}{2m} \|Xw - y\|^2 + \lambda \|w\|_1$

$l_1$  regularization

→ Produces solutions ( $w$ ) with many zero coefficients

$\|w\|_1 = \sum_{j=1}^d |w_j|$

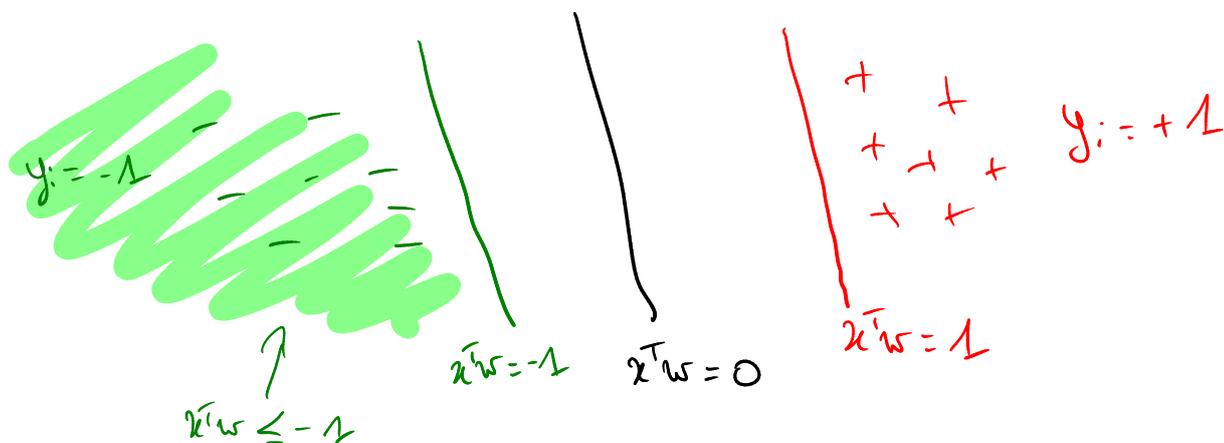
⇒ Sparse solutions / Feature selection

## ② Support vector machines (SVM)

Data:  $\{(x_i, y_i)\}_{i=1, \dots, m}$      $x_i \in \mathbb{R}^d$      $y_i \in \{-1, +1\}$   
↑  
Binary label

Goal (classification): Compute  $y_i$  given  $x_i$  through a linear model  $x \mapsto x^T w$

Idea: For  $(x_i, y_i)$ , we should have  
 $x_i^T w \geq 1$  if  $y_i = 1$   
and  $x_i^T w \leq -1$  if  $y_i = -1$



Optimization problem

$$\text{minimize}_{w \in \mathbb{R}^d} \frac{1}{m} \left\{ \sum_{i=1}^m \max(1 - y_i x_i^T w, 0) \right\} + \lambda \|w\|_1$$

↑ hinge loss

↓  $\ell_1$  regularization for robustness to outliers

Extensions

- \* Use  $\ell_2$  regularization instead of  $\ell_1$
- \* Use a nonlinear function of  $x$  (kernel SVM)

### ③ Deep neural networks

Illustration: Feed-forward NN, multi class classification

Data:  $\{ (x_i, y_i) \}_{i=1..m}$   $x_i \in \mathbb{R}^{d_0}$   
 $y_i \in \mathbb{R}^K$   $y_i = [y_{ik}]_{k=1..K}$

$K$  classes,  $y_i$  1-hot encoding of the class of  $x_i$   
 $y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \rightarrow$  class of  $x_i$

Goal: Classification (between  $K$  classes)

Mapping: Neural network with  $L$  layers

$$x^0 \in \mathbb{R}^{d_0} \mapsto x^1 \in \mathbb{R}^{d_1} \mapsto \dots \mapsto x^L \in \mathbb{R}^{d_L}$$

$$\forall l=1..L, \quad x^l = \sigma \left( \underbrace{W^l x^{l-1} + b^l}_{\substack{\text{linear transformation} \\ \text{of } x^{l-1}}} \right)$$

$\uparrow$  weight matrix in  $\mathbb{R}^{d_l \times d_{l-1}}$        $\uparrow$  bias vector in  $\mathbb{R}^{d_l}$

$\sigma$ : nonlinear activation applied component-wise

Ex) ReLU,  $\sigma(x^{l-1}) = \begin{bmatrix} \max(x_1^{l-1}, 0) \\ \vdots \\ \max(x_{d_{l-1}}^{l-1}, 0) \end{bmatrix} \in \mathbb{R}^{d_{l-1}}$

Given one input  $x_i$  as  $x^0$ , the NN outputs  $x^L \in \mathbb{R}^{d_L}$  with  $d_L = K$  (number of classes)

$\hookrightarrow$  The parameters of the network are the weight matrices  $(W^l)$  and the bias vectors  $(b^l)$



# BASICS OF OPTIMIZATION

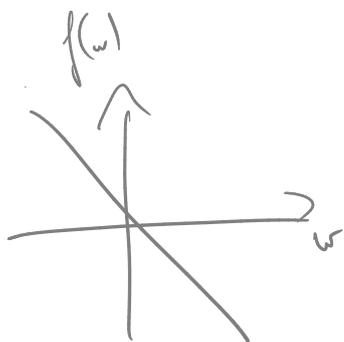
We consider a general optimization problem

$$(P) \quad \underset{w \in \mathbb{R}^d}{\text{minimize}} \quad f(w) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

Definitions:

- $\min_{w \in \mathbb{R}^d} f(w)$ : optimal value of (P)  
best value of  $f$  (smallest) obtained  
among all  $w$  in  $\mathbb{R}^d$

In general,  $\min_{w \in \mathbb{R}^d} f(w) \in \mathbb{R} \cup \{-\infty, +\infty\}$



$\min_{w \in \mathbb{R}^d} f(w) = -\infty$  : (P) is unbounded

$\min_{w \in \mathbb{R}^d} f(w) = +\infty$  : (P) is infeasible  
(only occurs for special  $f$ 's)

- $\underset{w \in \mathbb{R}^d}{\text{argmin}} f(w) \subseteq \mathbb{R}^d$  : set of solutions of (P)

$$\underset{w \in \mathbb{R}^d}{\text{argmin}} f(w) = \left\{ \bar{w} \in \mathbb{R}^d \mid f(\bar{w}) = \min_{w \in \mathbb{R}^d} f(w) \right\}$$

↳ This set can be empty (No solutions)  
finite (1 or more solutions, but finitely many)  
or infinite

Ex)  $f: w \mapsto 0$

$$\underset{w \in \mathbb{R}^d}{\text{argmin}} f(w) = \mathbb{R}^d$$

Properties:  $\forall \alpha > 0,$

$$\min_{w \in \mathbb{R}^d} \{ \alpha \times f(w) \} = \alpha \times \min_{w \in \mathbb{R}^d} f(w)$$

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \{ \alpha f(w) \} = \operatorname{argmin}_{w \in \mathbb{R}^d} f(w)$$

$\Rightarrow$  We say that (P) is equivalent to minimize  $\{ \alpha f(w) \}$  for any  $\alpha > 0$

Ex) minimize  $\frac{1}{n} \sum_{i=1}^n \dots$  equivalent to minimize  $\sum_{i=1}^n \dots$

! In general, it is intractable to compute the set of solutions directly  $\Rightarrow$  we need to use algorithms to compute approximate solutions in finite time

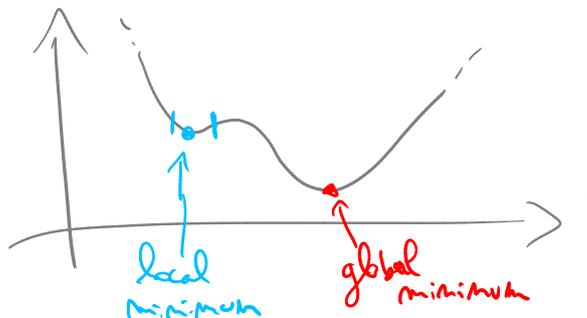
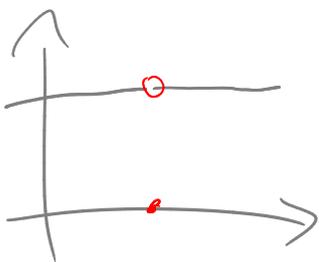
Def: • A solution  $\bar{w} \in \operatorname{argmin}_{w \in \mathbb{R}^d} f(w)$  is called a **global minimum**

$$\forall w \in \mathbb{R}^d, \underline{f(\bar{w}) \leq f(w)}$$

• A vector  $\bar{w} \in \mathbb{R}^d$  is called a **local minimum** if

$\exists w \in \mathbb{R}^d$  such that  $\|w - \bar{w}\|$  is small enough,

$$f(\bar{w}) \leq f(w)$$



NB  
Global  
 $\Downarrow$   
Local

↳ Most optimization algorithms in this course are designed to find local minima

↳ But in general, finding local minima is not tractable  $\Rightarrow$  Need more assumptions on  $f$

↳ With differentiability and convexity, we can

- characterize local minima (or approximations thereof)
- in finite time / in an easy way

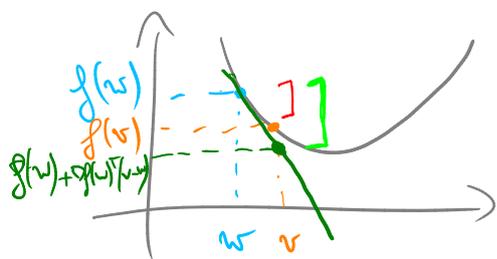
## a) Differentiability

Def:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C^1$  if its derivative exists for any  $w \in \mathbb{R}^d$  and the derivative is a continuous function



$\Rightarrow \forall w \in \mathbb{R}^d$ , there exists a vector called the gradient of  $f$  at  $w$ , denoted by  $\nabla f(w) \in \mathbb{R}^d$  that represents the derivative of  $f$  at  $w$

Fundamental property: If  $f$  is  $C^1$  and  $w \in \mathbb{R}^d$ , then for any  $v \in \mathbb{R}^d$  such that  $\|v-w\|$  is small enough, then



$$f(v) - f(w) \approx \nabla f(w)^T (v-w)$$

"approximates"  $\uparrow$   $\mathbb{R}^d$   $\mathbb{R}^d$

↳ Computing  $f(v) - f(w)$  involves computing  $f$

↳ Computing  $\nabla f(w)^T (v-w)$  for  $v$  involves only 1 calculation of  $\nabla f(w)$

→ The fact that  $f(v) \approx f(w) + \nabla f(w)^T (v-w)$  for small  $\|v-w\|$  is a local property that has connections to local optimality

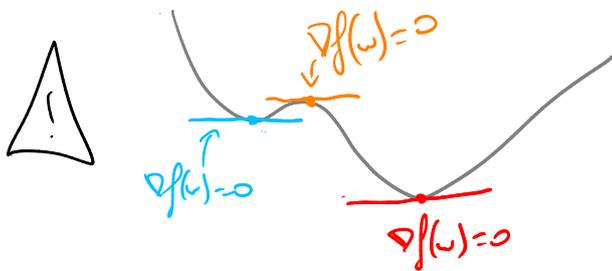
Theorem: [First-order necessary conditions]

Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be  $C^1$ .

If  $\bar{w} \in \mathbb{R}^d$  is a local minimum of  $f$ ,

then  $\nabla f(\bar{w}) = 0$ .

$$\nabla f(\bar{w}) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$



\* There can be  $\bar{v} \in \mathbb{R}^d$  that are not local minima but for which  $\nabla f(\bar{v}) = 0$   
( $\bar{x}$ : local maxima)

\* This is why the condition is called necessary (IF... THEN...)

↳ Even if the condition is only necessary, it gives a way to check whether a vector is not a local minimum in finite time (time required to check  $\nabla f(w) = 0$ )

IF local minimum THEN  $\nabla f(w) = 0$

IF  $\nabla f(w) \neq 0$  THEN NOT local minimum

## b) Convexity

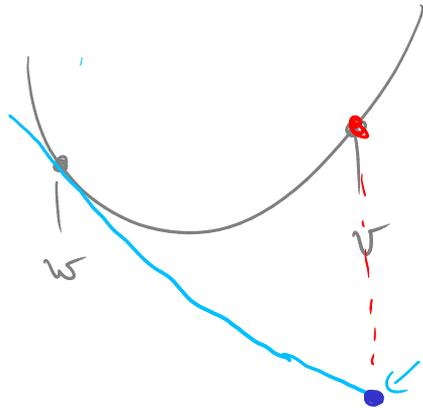
Def:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $C^1$  is convex if

$$\forall (w, v) \in (\mathbb{R}^d)^2, \quad f(v) \geq f(w) + \nabla f(w)^T (v-w)$$

Convex



not convex



Theorem: Let  $f$  be  $C^1$  convex. Then

$$(i) \quad [\bar{w} \in \mathbb{R}^d \text{ local minimum}] \Leftrightarrow [\nabla f(\bar{w}) = 0]$$

$$(ii) \quad [\bar{w} \in \mathbb{R}^d \text{ local minimum}] \Leftrightarrow [\bar{w} \in \mathbb{R}^d \text{ global minimum}]$$

Takeaway:

- Finding global minima is intractable in practice...
  - ... but not for  $C^1$  and convex functions  
(like the objective function of linear least squares)
- For those functions, checking the gradient is enough to determine whether a point is a global minimum