

Optimization for Machine Learning

September 20, 2023

Today
—

Syllabus

Motivation

Introduction to optimization

Syllabus

↳ The course is in English as most resources in the field are produced in English → Please play along!

↳ Questions / Discussions in French are always welcome at the end or via email

↳ Instructors:
Clément ROYER (First 4 lectures)
Florentin GOYENS (Last 4 lectures)

Contact: clement.royer@lamsade.dauphine.fr

↳ Schedule: 8 x 3 hours either 8.30am - 11.45am or 1.45pm - 5pm

↳ Tentative outline

09/20 (am): Intro to optimization

09/27 (am): A modern view on gradient descent

10/04 (am): Accelerated methods

10/05 (pm): Proximal methods and regularization

11/08 (am): Stochastic gradient (basics)

11/09 (pm): Stochastic gradient (advanced)

11/15 (am): Large-scale optimization

11/16 (pm): Overview

Assessment

- Exam (60%): December 13, 10am - 12pm
(will be given in both French and English)

- Course project (40%)
Individual, essentially a long Python notebook with coding tasks and questions
(loose timeline: between Oct 6 and Dec 23)

Course webpage

- Lecture notes to be updated
- Notebooks
- Exercises + solutions
- Virtual board

MOTIVATION

↳ Optimization: • Field of study that is concerned with making the best decision out of a set of alternatives

- As a subfield of theoretical computer science / computational mathematics, matured between 1980s - 2000s with efficient software (EX: Gurobi, CPLEX)

- Success stories:
 - Finance (portfolio allocation)
 - Supply chain
 - Facility location
 - Vaccine dispatch

- Shift in the past decades (2000s - 2020s)

Data-driven optimization

⇒ leads to specific optimization formulations

⇒ Changes the preferred classes of algorithms (compared to classical applications of optimization)

↳ Old algorithms became trendy (ex: stochastic gradient)

↳ Other algorithms are irrelevant (ex: Newton's method)

Typical ML optimization

L> Starts with data

$$D = \left\{ (x_i, y_i) \right\}_{i=1}^m$$

input output

m samples/
data points

L> Goal: Discover a mapping between the x_i s and the y_i s

$$\Rightarrow h \text{ such that } h(x_i) \approx y_i$$

Process: Training / Learning

L> Typically h is defined/parameterized with a vector $w \in \mathbb{R}^d$ and learning consists in finding the best value for w given the data

↓
Optimize!

An optimization problem associated with this task will have the form

Want w that gives the smallest objective value → minimize

variables/parameters → $w \in \mathbb{R}^d$

$$J_D(w) + \lambda r(w)$$

$f(w)$

↓
Objective function (to be optimized as a function of w)

$$w \in \mathbb{R}^d : w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \quad w_i \in \mathbb{R}$$

Specific objective structure: $J_D(w) + \lambda r(w)$

- $J_D(w)$ training / empirical loss that depends on the data

Average over samples
↓
n

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

Model
↑
Data points

Data points
↑
loss function
(measures agreement between $h(x_i; w)$ and y_i)

→ "data-fitting term"

⇒ The structure of J_D and its dependency on D are the reasons why some algorithms are better suited than others to solve this problem

- $r(w)$: "regularization term"

→ Does not depend on data

→ Is used to promote properties on w (and therefore on the mapping h)

- $\lambda > 0$: "regularization parameter"

⇒ Expresses a tradeoff between fitting the data and satisfying the properties represented by r

Examples

① Linear least squares

$$\text{Data} = \left\{ (x_i, y_i) \right\}_{i=1}^m \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R} \quad \forall i=1..m$$

$\mathbb{R}^{m \times d}$: matrices with m rows and d columns

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

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

$$x_i = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^d$$

$$x_i^T = [\dots \dots]$$

Goal: Find a linear relationship between $\{x_i\}$ and $\{y_i\}$
i.e. finding $w \in \mathbb{R}^d$ such that

$$Xw \approx y \quad (\Leftrightarrow) \quad \forall i=1..m$$

$$x_i^T w \approx y_i$$

\hookrightarrow There may not exist w such that $x_i^T w = y_i \quad \forall i$
but with optimization you can find the best possible w

Problem

$$\text{minimize}_{w \in \mathbb{R}^d} \frac{1}{2m} \|Xw - y\|^2 = \frac{1}{2m} \sum_{i=1}^m (x_i^T w - y_i)^2$$

$\underbrace{\hspace{10em}}_{\mathbb{R}^m} \quad \quad \quad \underbrace{\hspace{10em}}_{\mathbb{R}}$

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

$$\|v\| = \sqrt{v^T v}$$

(Euclidean norm)

$$\forall (v, u) \in (\mathbb{R}^m)^2, v^T u = \sum_{i=1}^m v_i u_i$$

scalar product

⇒ Basic formulation, also called / equivalent to linear regression
 ⇒ Variants

• Ridge regression

minimize $w \in \mathbb{R}^d$ $\frac{1}{2n} \|Xw - y\|^2 + \underbrace{\frac{\lambda}{2} \|w\|^2}_{\text{regularization}}$ with $\lambda > 0$

$\xrightarrow{\quad} \sum_{i=1}^d w_i^2$

(1996) • LASSO / l_1 regression

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

$\|w\|_1 := \sum_{i=1}^d |w_i|$

② Support vector machines (SVM)

Goal: Given $\{(x_i, y_i)\}_{i=1}^m$, $x_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$,
 find a mapping that classifies the x_i s
 correctly

Approach: Consider linear mappings that define hyperplanes

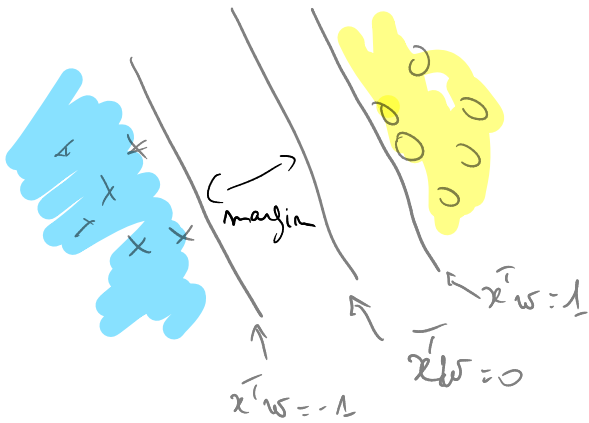


Seek $h(x; w) = x^T w$ (also possible: $x^T w + v$)
 $w \in \mathbb{R}^d$ $v \in \mathbb{R}$

Such that

$$x_i^T w \geq 1 \quad \text{if } y_i = 1$$

$$x_i^T w \leq -1 \quad \text{if } y_i = -1$$



Problem:

"Linear SVM"

↓
connected
with optimization
in the late 1990s

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

$$\frac{1}{n} \sum_{i=1}^n \max(1 - y_i x_i^T w, 0)$$

hinge loss $(x_i^T w, y_i)$

if $x_i^T w > 1$ and $y_i = 1 \Rightarrow 0$

if $x_i^T w < -1$ and $y_i = -1 \Rightarrow 0$

if $x_i^T w > 1$ and $y_i = -1 \Rightarrow \text{loss} > 0$
 $x_i^T w < -1$ and $y_i = 1 \Rightarrow$

if $|x_i^T w| \leq 1 \Rightarrow \text{loss "small" (between 0 and 2)}$

↳ Can add regularization terms (as in least squares)

↳ Generalized formulations \Rightarrow kernel SVM

Remark:

$$\forall f: \mathbb{R}^d \rightarrow \mathbb{R}$$

minimize
 $w \in \mathbb{R}^d$ $f(w)$

and

minimize
 $w \in \mathbb{R}^d$ $\lambda \times f(w)$ with $\lambda > 0$

have the same set of solutions

3

Deep neural networks

↳ Feed-forward architecture, multi-class classification

↳ Data: $\{(x_i, [y_{ij}]_j)\}_{i=1}^n$

$$x_i \in \mathbb{R}^{d_0}$$

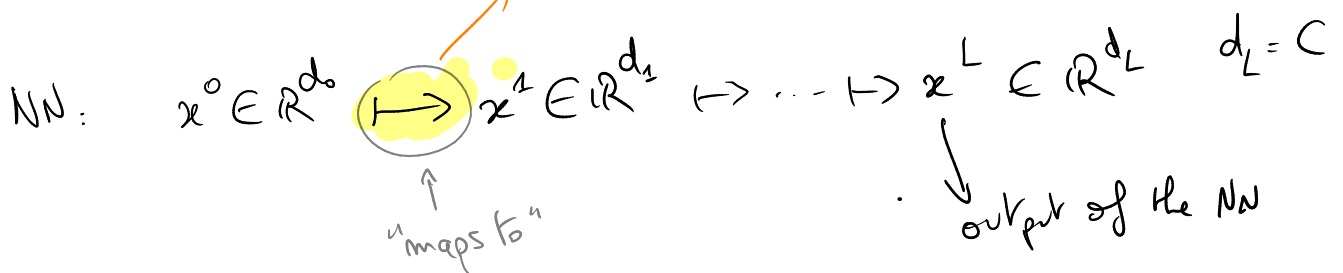
$$y_i = [y_{ij}]_j \in \mathbb{R}^C$$

C number of classes

1-hot encoding $y_{ij} = 1$ if x_i is in class j
 $y_{ij} = 0$ otherwise

↳ Neural architecture
L layers

First layer of the NN computes x^1 from x^0



$\forall l = 1..L,$

$$x^l = \sigma(W^l x^{l-1} + b^l)$$

activation function applied componentwise

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

bias vector $b^l \in \mathbb{R}^{d_l}$

$$\forall u \in \mathbb{R}^{d_l}, \sigma(u) := [\sigma(u_i)]_{i=1..d_l}$$

Ex): $\sigma(t) = t$ (identity)

$\sigma(t) = \max(t, 0)$ (ReLU)

$\sigma(t) = \tanh(t)$

Model: Takes $x_i \in \mathbb{R}^{d_0}$ as input ($x^0 = x_i$)
 and outputs $x_i^L \in \mathbb{R}^{d_L}$ $d_L = C$

Parameters: $(w^l, b^l)_{l=1..L} \Rightarrow$ stacked into a vector

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

$$d = \underbrace{d_2}_{b^2} + \underbrace{d_1}_{w^1} d_0 + \dots + \underbrace{d_L}_{b^L} + \underbrace{d_L}_{w^L} d_{L-1}$$

$\rightarrow L$ can be very large (hundreds at least)

$\rightarrow d$ usually is very large

Optimization problem

minimize $w \in \mathbb{R}^d$ $\frac{1}{n} \sum_{i=1}^n \left\{ \log \left(\sum_{j=1}^C \exp(x_{ij}^L) \right) - \sum_{j=1}^C y_{ij} x_{ij}^L \right\}$

$x_i^L = \left[x_{ij}^L \right]_{j=1}^C$

$\left\{ \right\} \Rightarrow$ negative log-likelihood

Want x_i^L to represent probabilities that x_i belongs to each of the classes

$$y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{correct class } j$$

ideal model

$$x_i^L = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{what the network learns}$$

$$x_i^L = \begin{bmatrix} p_1 \\ \vdots \\ p_c \end{bmatrix} \leftarrow \text{probabilities of belonging to a class}$$

Goal: Maximize P_j
(want $P_j > P_i \forall i \neq j$)

① ② ③

↳ All these examples involve:

- functions of real variables that are continuous
- structure (finite sum, regularization, layer calculation, ...)

Goal of the work: Understand how to choose the best algorithms for problems of this form

BASICS OF CONTINUOUS OPTIMIZATION

An optimization problem (for us):

$$(P) \left\{ \begin{array}{l} \text{minimize} \\ w \in \mathbb{R}^d \end{array} f(w) \right.$$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$
 f : objective function
 w : (vector of) variables

(More general formulation: minimize $f(w)$ s.t. $w \in F$)

↑ "Subject to" ↑ Constraints on w

⇒ We focus on the unconstrained setting ($F = \mathbb{R}^d$)

Solution set of (P)

$$\cdot \operatorname{argmin}_{w \in \mathbb{R}^d} f(w) = \left\{ \bar{w} \in \mathbb{R}^d \mid f(\bar{w}) \leq f(w) \forall w \in \mathbb{R}^d \right\}$$

set $\subseteq \mathbb{R}^d$

↳ can be empty ($f(w) = w$ and $d = 1$)

↳ can be infinite ($f(w) = 0$)

$$\cdot \min_{w \in \mathbb{R}^d} f(w) : \text{minimal value of } f$$

$\mathbb{R}^d \cup \{+\infty, -\infty\}$

$\lambda > 0$,

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

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

Example of
equivalent
formulations

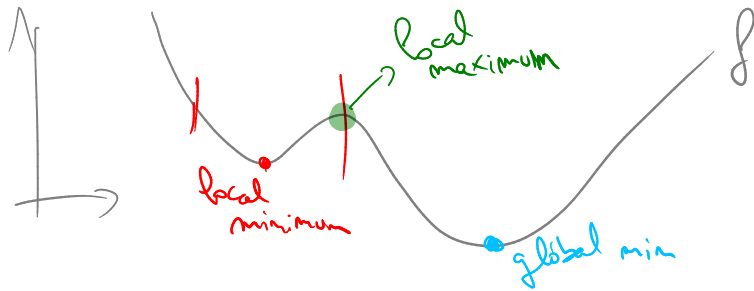
Two problems are equivalent if the solution of one can be readily obtained by solving the other

⚠ In general, it is hard (computationally hard/intractable) to compute the argmin directly
↳ simply by function comparison

⇒ In practice, we exploit other properties of the objective to build algorithms that compute approximate solutions

Def. $\bar{w} \in \operatorname{argmin}_w f(w)$ is called a global minimum

$\bar{w} \in \mathbb{R}^d$ is called a local minimum of f if
 $\forall w \in \mathbb{R}^d$ such that $\|w - \bar{w}\|$ is small enough $\Rightarrow f(\bar{w}) \leq f(w)$



\bar{w} local maximum
 $f(\bar{w}) \geq f(w)$ for
 $\|w - \bar{w}\|$ small enough

Local minima are (on principle) easier to converge to than global minima but they are still hard to compute and they may correspond to values of f that are far from the optimum

\Rightarrow In some data science tasks, it is possible to quantify how far local minima are from global ones

\hookrightarrow The notions of convexity and differentiability allow to characterize local minima (or approximations) in finite time in an "easy" way

a) Differentiability

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is C^1 (continuously differentiable)

if its derivative exists $\forall w \in \mathbb{R}^d$ and is a continuous mapping.

$\Rightarrow \forall w \in \mathbb{R}^d, \exists \nabla f(w) \in \mathbb{R}^d$ such that
 \uparrow
 gradient of f at w

$\forall v \in \mathbb{R}^d, f(v) \approx f(w) + \underbrace{\nabla f(w)^T (v-w)}_{\text{linear function of } v}$ when $\|v-w\|$ is small enough

\Rightarrow Approximation of the change in function values around w

\Rightarrow Local property that can be connected to local optimality.

Theorem (First-order **necessary** conditions)

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ and $\bar{w} \in \mathbb{R}^d$

[\bar{w} is a local minimum of f]

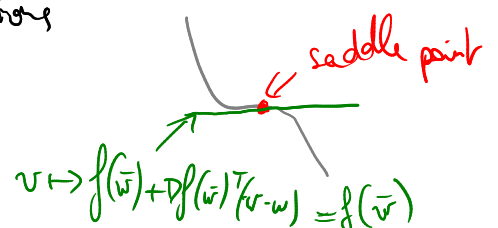
$\Rightarrow \nabla f(\bar{w}) = 0_{\mathbb{R}^d}$ ($\|\nabla f(\bar{w})\| = 0$)

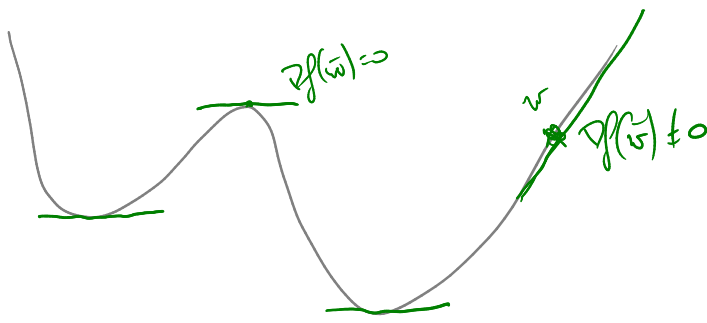
Necessary: the converse is not true

In fact, [\bar{w} is a local maximum of f] $\Rightarrow \nabla f(\bar{w}) = 0_{\mathbb{R}^d}$

[\bar{w} is a saddle point of f] $\Rightarrow \nabla f(\bar{w}) = 0_{\mathbb{R}^d}$

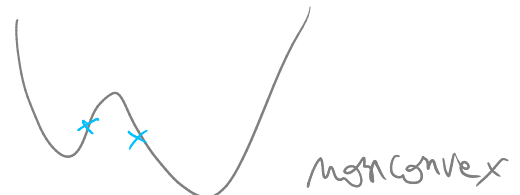
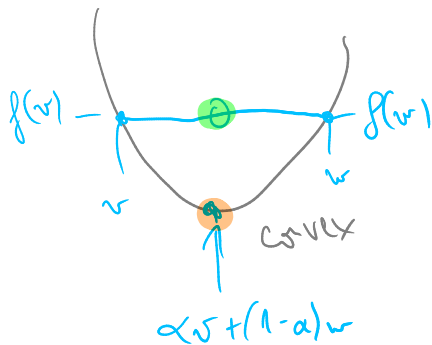
Saddle point: local min in some directions
 local max in others





Ⓟ Convexity

A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is called convex if
 $\forall (w, v) \in (\mathbb{R}^d)^2, \forall \alpha \in [0, 1],$
 $f(\alpha v + (1-\alpha)w) \leq \alpha f(v) + (1-\alpha)f(w)$



Theorem: Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function. Then,

- i) any local minimum of f is a global minimum
- ii) $\exists!$ $f \in C^1, \forall \bar{w} \in \mathbb{R}^d$

$$[\bar{w} \text{ is a global minimum of } f] \Leftrightarrow [\nabla f(\bar{w}) = 0]$$

↑
 Necessary and sufficient
 condition for global optimality

$\hookrightarrow C^1$ convex functions are easy to minimize
(ex: linear least squares, SVM)

\Rightarrow Need to find \bar{w} such that $\nabla f(\bar{w}) = 0$