

OPTIMIZATION FOR MACHINE LEARNING

September 16, 2024

Today (Lecture 1/16)

- Logistics
- Intro to optimization (Pt 1)

LOGISTICS

Course google doc (key resource)

<https://bit.ly/3XwecET>

Course webpage on my website

<https://www.lamsade.dauphine.fr/~croyer/teachOIM.html>

INTRODUCTION TO OPTIMIZATION (IN ML)

Formal definition: Optimization is a field of study concerned with making the best decision out of a set of alternatives

L> As a field of computational mathematics, optimization matured between 1980s and 2000s
=> Success stories: Finance, Vaccine dispatch, Scheduling, Moving plan for Dauphine's offices

→ Model-driven approach

L> 2000s-2020s: Shift of paradigm from model-driven to data-driven

* changed the problem formulations people care about

* changed the state-of-the-art algorithms

(1990s: Newton's method
2024: Stochastic gradient and its variants)

Typical ML optimization problem

↳ Data: $\mathcal{D} = \{ (a_i, y_i) \}_{i=1 \dots n}$ n samples

↑ Input (vector, matrix, ...)
↑ Output (scalar, vector, matrix, tensor, ...)

↳ Goal: * Learn a mapping from a_i s to y_i s
i.e. a function h such that $h(a_i) \approx y_i$

* In practice, consider that h is parameterized by a vector $x \in \mathbb{R}^d$
 $h(\cdot; x) : a \mapsto h(a; x)$

⇒ The goal becomes to find the value of x that yields the best mapping in terms of matching $h(a_i; x)$ with y_i

Problem formulation

Alternative:
minimize
 $x \in \mathbb{R}^d$

→ **min** $x \in \mathbb{R}^d$ $\{ \underbrace{f_{\mathcal{D}}(x)}_{\text{objective function } f(x)} + \underbrace{\lambda \sum r(x)}_{\lambda \sum r(x)} \}$

"Minimize
 $f_{\mathcal{D}}(x) + \lambda r(x)$
over $x \in \mathbb{R}^d$ "

x : decision variables / parameters / weights
 \Rightarrow Define mapping h (and thus how it fits the data)

J_D : data-fitting term, defined using h and $D = \{(a_i, y_i)\}_{i=1..m}$

Typical structure:

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

Average over all data points

loss function: measures agreement between $h(a_i; x)$ and y_i

mapping parameterized by x

$\in \mathbb{R}$

$\pi(x)$: regularization term

* typically doesn't depend on the data

* used to penalize values of x that do not satisfy desired properties (e.g. overfitting, density)

$\lambda > 0$: regularization parameter

Expresses trade-off between regularization and data-fitting

$\lambda \gg 1$: more weight on the regularization

$\lambda \ll 1$: data-fitting term

EXAMPLES

① Linear least squares (aka linear regression)

Data $D = \{(a_i, y_i)\}_{i=1 \dots m}$ $a_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$

\Rightarrow Representation

$$A = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}$$

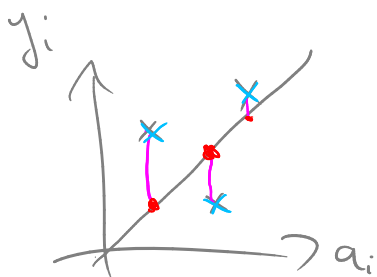
Set of matrices with
real coefficients
m rows
d columns

$$a \in \mathbb{R}^d \quad a = \begin{bmatrix} \\ \\ \end{bmatrix} \uparrow d$$

$$a^T = \begin{bmatrix} & & \end{bmatrix} \uparrow d$$

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

Goal: Compute a linear model of the data, that is a vector $x \in \mathbb{R}^d$ such that $a_i^T x \approx y_i \quad \forall i=1 \dots m$



$$\forall (u, v) \in (\mathbb{R}^d)^2, \quad u^T v = \sum_{j=1}^d u_j v_j$$

$$u^T v = \begin{bmatrix} u^T \\ \end{bmatrix} \begin{bmatrix} v \\ \end{bmatrix}$$

Mapping $a \mapsto a^T x$

Problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{2m} \|Ax - y\|^2 = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{2} (a_i^T x - y_i)^2 \right\}$$

loss (l2)
 $l(h, y) = \frac{1}{2} (h - y)^2$

where $\|Ax - y\|^2 = \sum_{i=1}^m (a_i^T x - y_i)^2$

$$\forall v \in \mathbb{R}^m, \quad \|v\|^2 = \sum_{j=1}^m v_j^2$$

Variants

* Ridge regression:

$$\min_{x \in \mathbb{R}^d} \frac{1}{2m} \|Ax - y\|^2 + \lambda \|x\|^2$$

data fitting regularization

* LASSO:

$$\min_{x \in \mathbb{R}^d} \frac{1}{2m} \|Ax - y\|^2 + \lambda \|x\|_1$$

\uparrow
 $\|x\|_1 = \sum_{j=1}^d |x_j|$

\Rightarrow Produces vector x with more zero coordinates than the original formulation

(2) Matrix completion

\hookrightarrow Data: $\{(A_i, y_i)\}_{i=1..m}$

$$A = \begin{bmatrix} A_{11} & \dots & A_{1p} \\ \vdots & & \vdots \\ A_{q1} & \dots & A_{qp} \end{bmatrix} \in \mathbb{R}^{q \times p}$$

$$A^T = \begin{bmatrix} A_{11} & \dots & A_{q1} \\ \vdots & & \vdots \\ A_{1p} & \dots & A_{qp} \end{bmatrix} \in \mathbb{R}^{p \times q}$$

$y_i \in \mathbb{R}$

$$A_i = A_i^T \in \mathbb{R}^{d \times d}$$

Symmetric matrix

$$\begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} \checkmark \quad \begin{bmatrix} 1 & 2 \\ 3 & 1 \end{bmatrix} \times$$

\hookrightarrow Underlying ground truth $\underbrace{\text{matrix}}_{\text{symmetric}} M^* \in \mathbb{R}^{d \times d}$

$\rightarrow A_i$ measurement operators

$$\rightarrow y_i = \langle A_i, M^* \rangle = \text{trace}(A_i^T M^*)$$

Goal: Find a good approximation for M^* that corresponds to $\{(A_i, y_i)\}_{i=1..m}$

Problem

$$\min_{X \in \mathbb{R}^{d \times d}} \frac{1}{2m} \sum_{i=1}^m (\langle A_i, X \rangle - y_i)^2$$

(similar to linear least squares problem)

\hookrightarrow If you know that M^* is (symmetric) positive semidefinite

consider an alternate formulation

$A \in \mathbb{R}^{d \times d}$ is
positive semi-definite
if symmetric
- $\forall v \in \mathbb{R}^d$,
 $v^T A v \geq 0$

$$\min_{U \in \mathbb{R}^{d \times d}} \frac{1}{2m} \sum_{i=1}^m (\langle A_i, U U^T \rangle - y_i)^2$$

$$(U U^T)^T = U U^T$$

$$v^T U U^T v = \|U^T v\|^2 \geq 0$$

\Rightarrow Nonlinear least squares

\hookrightarrow If you further know that M^* has low rank r
(can be represented by $< d \times d$ coefficients), can

$$\begin{matrix} v^T A v & & \\ \underbrace{\quad}_{1 \times a} & \underbrace{\quad}_{d \times d} & \underbrace{\quad}_{a \times 1} \\ & & \underbrace{\quad}_{1 \times 1} \end{matrix}$$

consider

$$\min_{U \in \mathbb{R}^{d \times r}} \frac{1}{2m} \sum_{i=1}^m (\langle \overbrace{A_i}^{d \times d}, \underbrace{U U^T}_{\substack{d \times r \quad r \times d \\ d \times d}} \rangle - y_i)^2$$

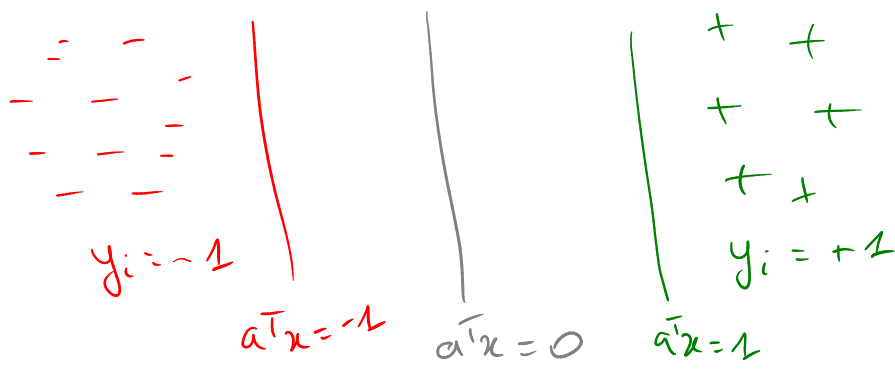
③ Support vector machine (SVM)

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

Goal: Perform binary classification, i.e. learn how to map a_i to y_i

Approach: Consider a linear model parameterized by $x \in \mathbb{R}^d$ $a \mapsto a^T x$
 $\mathbb{R}^d \quad \mathbb{R}$

- when $y_i = +1$, want $a^T x \geq 1$ ← "Margin of 1"
- when $y_i = -1$, want $a^T x \leq -1$



Problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max(1 - y_i a^T x, 0) + \lambda \|x\|^2$$

\downarrow
 $l(h, y) = \max(1 - yh, 0)$
 "Hinge loss"

↳ Variants:

- * Add other regularization (l_1)
- * Go beyond linear models (kernel SVM)
 $a \mapsto a^T x$ $a \mapsto \psi(a)^T x$
 with ψ a nonlinear vector function

④ Multi class classification using neural networks

↳ Basic feedforward NN

$$a^0 \in \mathbb{R}^{d_0} \mapsto a^1 \in \mathbb{R}^{d_1} \mapsto a^2 \in \mathbb{R}^{d_2} \mapsto \dots \mapsto a^L \in \mathbb{R}^{d_L}$$

L : number of layers

$$\forall l=1..L, \quad a^l = \sigma(W^l a^{l-1} + b^l)$$

(Typically Nonlinear) activation function

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

bias vector
 \mathbb{R}^{d_l}

from \mathbb{R} to \mathbb{R} ReLU $\sigma(t) = \max(t, 0)$ applied componentwise

$$\forall v \in \mathbb{R}^m, \sigma(v) = \begin{bmatrix} \sigma(v_1) \\ \vdots \\ \sigma(v_m) \end{bmatrix} \in \mathbb{R}^m$$

$$(x) \quad \begin{matrix} \sigma(t) = t \\ \sigma(t) = \tanh(t) \end{matrix}$$

Mapping: $a \mapsto NN(a; x)$, where $NN(a; x)$ is the output of the last layer (a^L) when given a as the input (a^0) and using x as parameters

concatenation
of $w^1, b^1, \dots, w^L, b^L$
into one (long!) vector

$$x \in \mathbb{R}^d \quad d = d_1 d_0 + d_1 + \dots + d_L d_{L-1} + d_L$$

Data: $\{(a_i, y_i)\}_{i=1..n} \quad a_i \in \mathbb{R}^{d_0}$

$y_i \in \mathbb{R}^K$ K number of classes

y_i 1-hot encoding of the class

$$y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow k : \text{class of } a_i$$

Goal: Use NN to predict y_i from a_i

$NN(a_i; x)$ returns a vector that will be transformed into a vector of probabilities in \mathbb{R}^K

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_k \end{bmatrix} \quad \begin{matrix} \sum_{k=1}^K v_k = 1 \\ \forall k \quad v_k \geq 0 \end{matrix}$$

For simplicity, assume $NN(a_i; x) \in \mathbb{R}^K$

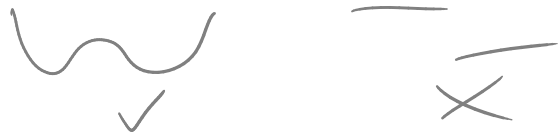
Optimization problem: Minimizing the negative log-likelihood of $NN(a_i; x)$ given y_i + Average over the samples

\Leftrightarrow Maximize the probability of correct classification

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left\{ \log \left(\sum_{k=1}^K \exp([NN(a_i, x)]_k) \right) - \sum_{k=1}^K [y_i]_k [NN(a_i, x)]_k \right\}$$

Takeaways

- All examples involve an average of data points (finite-sum structure)
- All involve a continuous objective function



- Additional structure: loss, mapping + possibly regularization

\Rightarrow We want to build algorithms that exploit this structure

Note: • S. J. Wright and B. Recht, Optimization for Data Analysis (2022, Cambridge University Press)

BASICS OF OPTIMIZATION

General setup (no visible dependency data)

$$(P) \quad \min_{x \in \mathbb{R}^d} f(x) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

Def: A solution of (P) aka a global minimum of (P) is a vector $x^* \in \mathbb{R}^d$ such that $f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^d$

Notation: The set of global minima is denoted by $\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \subseteq \mathbb{R}^d$

$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$ can be empty ($x \mapsto x$ in \mathbb{R})

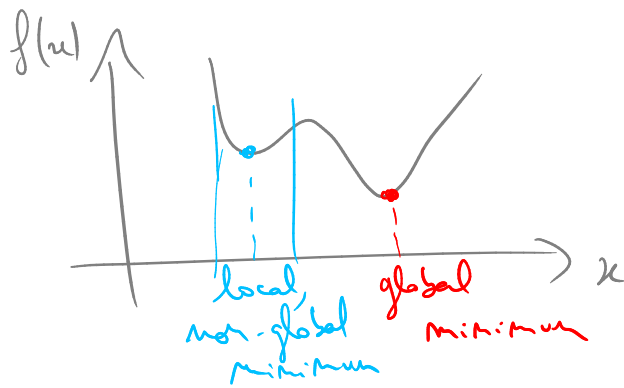
$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$ can be finite ($x \mapsto x^2$ in \mathbb{R})
for

$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$ can be infinite ($x \mapsto 0$ in \mathbb{R})

↳ In general, global minima are not tractable (hard or even impossible to compute in finite time)

Def: A vector $\bar{x} \in \mathbb{R}^d$ is a local minimum of (P) if $\forall x \in \mathbb{R}^d$, $\|x - \bar{x}\|$ sufficiently small, $f(\bar{x}) \leq f(x)$
 $\exists \varepsilon > 0, \forall x \in \mathbb{R}^d, (\|x - \bar{x}\| \leq \varepsilon) \Rightarrow f(\bar{x}) \leq f(x)$

NB: Any global minimum is a local minimum, but the converse is false in general



↳ Local minima still hard to compute in general
 ⇒ In ML, there are many problem instances for which we can find local minima or approximations of local minima in finite time
 ⇒ Local minima are easier to identify than global ones

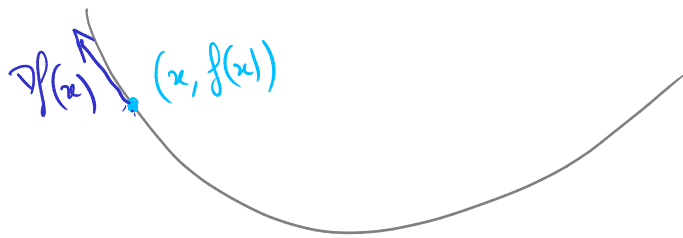
↳ We will assume properties on the objective function f that guarantee that local minima can be identified, such as differentiability and convexity

Two important properties

- C^1 (aka smooth for some authors)

Def: f is C^1 (" C^1 /continuously differentiable") if
 $\forall x \in \mathbb{R}^d$, there exists a vector $\nabla f(x) \in \mathbb{R}^d$
 (called the gradient of f at x) that represents the
 derivative of f at x and $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is

continuous



Theorem: (First-order necessary condition)

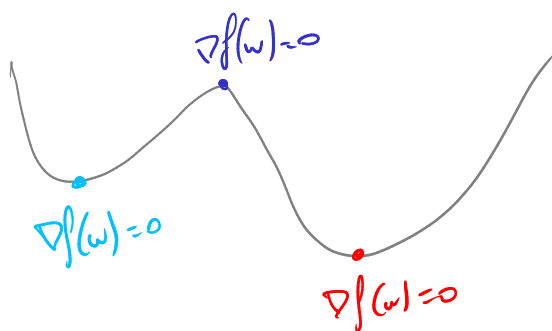
Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be C^1 .

(IF) $\bar{w} \in \mathbb{R}^d$ is a local minimum of f
 (Then) $\nabla f(\bar{w}) = 0$ ($\equiv \|\nabla f(\bar{w})\| = 0$)

→ the gradient characterizes local changes in the function

→ Necessary condition: (IF) ... (THEN) ...

- Does not help you identify local minima
- Helps you identifying points that are not local minima



• = local maximum
 (yet its gradient is 0)

Any point other than the red/light blue/dark blue ones has a non-zero gradient

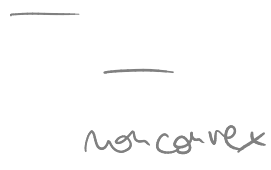
* Convex functions

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

$$\forall (x, y) \in (\mathbb{R}^d)^2, \forall \alpha \in [0, 1], f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

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

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



l_2 loss, hinge loss are convex

Theorem: Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ C^1 convex and consider (P).

(i) Any local minimum of (P) is a global minimum of (P)
(and thus local = global)

(ii) $[\bar{x} \in \mathbb{R}^d \text{ local minimum of (P)}] \Leftrightarrow [\nabla f(\bar{x}) = 0]$

↑
Necessary and sufficient condition

Takeaway:

- With C^1 + convexity, can find global minima by checking ∇f !
- With C^1 , the gradient still helps to identify points that are not local minima.

Exercise:

Find a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ C^1 nonconvex
such that $\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) = \{x \in \mathbb{R}^d \mid \nabla f(x) = 0\}$