

Optimization for Machine Learning

Lecture 1: Basics of optimization

Clément W. Royer

M2 IASD - 2025/2026

September 18, 2025



- 1 Optimization theory
- 2 Exercises
- 3 Bonus

- 1 Optimization theory
- 2 Exercises
- 3 Bonus

Formulation of an (unconstrained) optimization problem

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

Formulation of an (unconstrained) optimization problem

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

- \boldsymbol{w} represents the optimization variable(s);
- d is the dimension of the problem (we will assume $d \geq 1$);
- $f(\cdot)$ is the **objective/cost/loss** function.

Formulation of an (unconstrained) optimization problem

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

- \boldsymbol{w} represents the optimization variable(s);
- d is the dimension of the problem (we will assume $d \geq 1$);
- $f(\cdot)$ is the **objective/cost/loss** function.

Maximizing f is equivalent to minimizing $-f$!

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

- $\operatorname{argmin}_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{w})$: Set of solutions (can be empty).
- $\min_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{w})$: Optimal value (can be infinite).

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

- $\operatorname{argmin}_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{w})$: Set of solutions (can be empty).
- $\min_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{w})$: Optimal value (can be infinite).

Global and local minima

- \boldsymbol{w}^* is a solution or a **global minimum** of f if $f(\boldsymbol{w}^*) \leq f(\boldsymbol{w}) \forall \boldsymbol{w} \in \mathbb{R}^d$.
- \boldsymbol{w}^* is a **local minimum** of f if $f(\boldsymbol{w}^*) \leq f(\boldsymbol{w}) \forall \boldsymbol{w}, \|\boldsymbol{w} - \boldsymbol{w}^*\|_2 \leq \epsilon$ for some $\epsilon > 0$.

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

- $\text{argmin}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$: Set of solutions (can be empty).
- $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$: Optimal value (can be infinite).

Global and local minima

- \mathbf{w}^* is a solution or a **global minimum** of f if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w} \in \mathbb{R}^d$.
 - \mathbf{w}^* is a **local minimum** of f if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}, \|\mathbf{w} - \mathbf{w}^*\|_2 \leq \epsilon$ for some $\epsilon > 0$.
-
- Finding global/local minima is hard in general!
 - Regularity of f is needed.

First notion of regularity: Smoothness

Class of \mathcal{C}^1 functions

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable/ \mathcal{C}^1 if

- For any $\mathbf{w} \in \mathbb{R}^d$, the **gradient** $\nabla f(\mathbf{w})$ exists.
- $\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is continuous.

First notion of regularity: Smoothness

Class of \mathcal{C}^1 functions

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable/ \mathcal{C}^1 if

- For any $\mathbf{w} \in \mathbb{R}^d$, the **gradient** $\nabla f(\mathbf{w})$ exists.
- $\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is continuous.

Class of $\mathcal{C}_L^{1,1}$ functions ($L > 0$)

f is $\mathcal{C}_L^{1,1}$ if it is \mathcal{C}^1 and ∇f is L -Lipschitz continuous, i.e.

$$\forall (\mathbf{v}, \mathbf{w}) \in (\mathbb{R}^d)^2, \quad \|\nabla f(\mathbf{v}) - \nabla f(\mathbf{w})\| \leq L \|\mathbf{v} - \mathbf{w}\|.$$

Ex) Linear regression, logistic regression, etc.

Important for today

Function $f(\mathbf{w}) \in \mathbb{R}$ Gradient $\nabla f(\mathbf{w}) \in \mathbb{R}^d$

$$\mathbf{a}^T \mathbf{w} + \mathbf{b}$$
$$\frac{1}{2} \|\mathbf{w} + \mathbf{b}\|_2^2$$

$$\mathbf{a}$$
$$\mathbf{w} + \mathbf{b}$$

Important for today

Function $f(\mathbf{w}) \in \mathbb{R}$ Gradient $\nabla f(\mathbf{w}) \in \mathbb{R}^d$

$$\mathbf{a}^T \mathbf{w} + b$$
$$\frac{1}{2} \|\mathbf{w} + \mathbf{b}\|_2^2$$

$$\mathbf{a}$$
$$\mathbf{w} + \mathbf{b}$$

Next week How to compute derivatives in ML (bring laptops!).

Smoothness and optimality conditions

Problem: minimize $w \in \mathbb{R}^d$ $f(w)$, $f \in \mathcal{C}^1$.

First-order necessary condition

If w^* is a local minimum of the problem, then

$$\|\nabla f(w^*)\|_2 = 0.$$

- This condition is only necessary;
- A point such that $\|\nabla f(w^*)\|_2 = 0$ can also be a local maximum or a saddle point.

Smoothness and optimality conditions

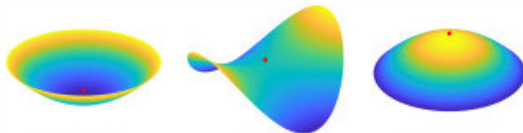
Problem: minimize $w \in \mathbb{R}^d$ $f(w)$, $f \in \mathcal{C}^1$.

First-order necessary condition

If w^* is a local minimum of the problem, then

$$\|\nabla f(w^*)\|_2 = 0.$$

- This condition is only necessary;
- A point such that $\|\nabla f(w^*)\|_2 = 0$ can also be a local maximum or a saddle point.



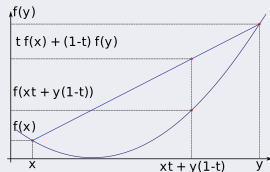
Picture from (Wright and Ma '22).

Another notion of regularity: Convexity

Generic definition (+Wikicommons picture)

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

$$\forall(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \\ f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t) f(\mathbf{v}).$$

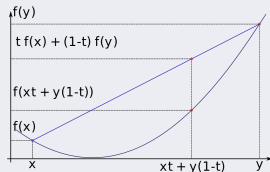


Another notion of regularity: Convexity

Generic definition (+Wikicommons picture)

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

$$\forall(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \\ f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t) f(\mathbf{v}).$$



Examples in ML

- Linear function $\mathbf{w} \mapsto \mathbf{a}^T \mathbf{w} + b$.
- ℓ_2 loss $\|\mathbf{w}\|_2^2 = \sum_{j=1}^d w_j^2$.
- Logistic loss.

Showing convexity (from Hardt and Recht '21)

Showing convexity with more than two variables is hard.

Showing convexity (from Hardt and Recht '21)

Showing convexity with more than two variables is hard.

Basic blocks

- All norms (and $\|\mathbf{w}\|_2^2$) are convex.
- All linear functions $\mathbf{w} \mapsto \mathbf{A}\mathbf{w} + \mathbf{b}$ are convex.
- f convex $\Rightarrow \alpha f$ convex $\forall \alpha \geq 0$.
- f, g convex $\Rightarrow f + g$ convex.
- f, g convex $\Rightarrow \max(f, g)$ convex.
- f convex $\Rightarrow \mathbf{w} \mapsto f(\mathbf{A}\mathbf{w} + \mathbf{b})$ convex.

Convexity and gradient

A continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}).$$

Convexity and gradient

A continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}).$$

A key inequality in optimization.

Convex optimization problem

$$\underset{\boldsymbol{w} \in \mathbb{R}^d}{\text{minimize}} \ f(\boldsymbol{w}), \ f \text{ convex.}$$

Convex optimization problem

$$\underset{\boldsymbol{w} \in \mathbb{R}^d}{\text{minimize}} \ f(\boldsymbol{w}), \ f \text{ convex.}$$

Theorem

Every local minimum of f is a global minimum.

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \ f(\mathbf{w}), \ f \text{ convex.}$$

Theorem

Every local minimum of f is a global minimum.

Corollary

If f is \mathcal{C}^1 ,

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \ f(\mathbf{w}) = \{ \bar{\mathbf{w}} \mid \|\nabla f(\bar{\mathbf{w}})\|_2 = 0 \}.$$

Any point with a zero gradient is a global minimum!

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in \mathcal{C}^1 is μ -strongly convex (or *strongly convex of modulus $\mu > 0$*) if for all $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ and $t \in [0, 1]$,

$$f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq t f(\mathbf{u}) + (1 - t)f(\mathbf{v}) - \frac{\mu}{2}t(1 - t)\|\mathbf{v} - \mathbf{u}\|_2^2.$$

Strong convexity

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in \mathcal{C}^1 is μ -strongly convex (or *strongly convex of modulus $\mu > 0$*) if for all $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ and $t \in [0, 1]$,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|_2^2.$$

Theorem

Any strongly convex function in \mathcal{C}^1 has a unique global minimizer.

Gradient and strong convexity

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $f \in \mathcal{C}^1$. Then,

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}) + \frac{\mu}{2}\|\mathbf{v} - \mathbf{u}\|_2^2.$$

Key rules

- For any $\mu > 0$ and $\mathbf{w}_0 \in \mathbb{R}^d$, $\mathbf{w} \mapsto \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}_0\|_2^2$ is μ -strongly convex.

Key rules

- For any $\mu > 0$ and $\mathbf{w}_0 \in \mathbb{R}^d$, $\mathbf{w} \mapsto \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}_0\|_2^2$ is μ -strongly convex.
- If f is μ -strongly convex and g is convex, $f + g$ is μ -strongly convex.

- M. Hardt and B. Recht, *Patterns, Predictions and Actions*, Princeton University Press, 2021.
- J. Wright and Y. Ma, *High-Dimensional Data Analysis with Low-Dimensional Models*, Cambridge University Press, 2022.
- S. J. Wright and B. Recht, *Optimization for Data Analysis*, Cambridge University Press, 2022.

- 1 Optimization theory
- 2 Exercises
- 3 Bonus

Show that the SVM objective

$$\mathbf{w} \in \mathbb{R}^d \longmapsto \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i \mathbf{x}_i^T \mathbf{w}, 0\} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

is a convex function for any $\lambda \geq 0$.

Exercise 1.b - Strong convexity

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be \mathcal{C}^1 and μ -strongly convex, and denote by \mathbf{w}^* the minimum of f .

- ❶ For any $\mathbf{w} \in \mathbb{R}^d$, show that the function

$$\varphi_{\mathbf{w}} : \mathbf{z} \mapsto f(\mathbf{w}) + \nabla f(\mathbf{w})^T(\mathbf{z} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{z} - \mathbf{w}\|^2$$

is strongly convex.

- ❷ Compute $\min_{\mathbf{z}} \varphi_{\mathbf{w}}(\mathbf{z})$ and $\operatorname{argmin}_{\mathbf{z}} \varphi_{\mathbf{w}}(\mathbf{z})$.

- ❸ Show that

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*)).$$

Exercise 1.c - Least-squares

Let $\mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\|_2 \neq 0$ and $\mathbf{y} \in \mathbb{R}^d$.

- 1 Consider the problem

$$\underset{w \in \mathbb{R}}{\text{minimize}} \quad \frac{1}{2} \|w\mathbf{x} - \mathbf{y}\|_2^2.$$

Is it convex? Is the minimum value 0?

- 2 Consider now the problem

$$\mathbf{W} \in \mathbb{R}^{d \times d} \longmapsto \frac{1}{2} \|\mathbf{W}\mathbf{x} - \mathbf{y}\|_2^2.$$

Is this a convex problem?

- 3 Justify that

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times d}} \frac{1}{2} \|\mathbf{W}\mathbf{x} - \mathbf{y}\|_2^2 = 0,$$

and find a global minimum. Is the minimum unique?

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be $\mathcal{C}_L^{1,1}$ and convex. Suppose that $\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$ and let $f^* = f(\mathbf{w}^*)$.

- ❶ Let $\mathbf{w} \in \mathbb{R}^d$. Show that $f(\mathbf{w}) - f(\mathbf{w}^*) \geq \frac{1}{2L} \|\nabla f(\mathbf{w})\|_2^2$.
- ❷ Let $(\mathbf{w}, \mathbf{v}) \in (\mathbb{R}^d)^2$. Show that

$$(\nabla f(\mathbf{v}) - \nabla f(\mathbf{w}))^\top (\mathbf{v} - \mathbf{w}) \geq \frac{1}{L} \|\nabla f(\mathbf{v}) - \nabla f(\mathbf{w})\|_2^2.$$

Consider $z \mapsto f(z) - \nabla f(\mathbf{v})^\top z$ and $z \mapsto f(z) - \nabla f(\mathbf{w})^\top z$.

- 1 Optimization theory
- 2 Exercises
- 3 Bonus

- Convex problems: All local minima are global!
- Nonconvex problems: May have local, non-global (aka spurious minima).

Landscape analysis

Identify classes of **nonconvex problems** for which there are no spurious minima (and possibly more).

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_L}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- $\mathbf{W}_i \in \mathbb{R}^{d_{i+1} \times d_i}$.
- $\mathbf{X} \in \mathbb{R}^{d_1 \times d_0}$, $\mathbf{Y} \in \mathbb{R}^{d_{L+1} \times d_0}$.
- $\|\mathbf{A}\|_F^2 = \sum_i \sum_j A_{ij}^2$.

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_L}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- $\mathbf{W}_i \in \mathbb{R}^{d_{i+1} \times d_i}$.
 - $\mathbf{X} \in \mathbb{R}^{d_1 \times d_0}$, $\mathbf{Y} \in \mathbb{R}^{d_{L+1} \times d_0}$.
 - $\|\mathbf{A}\|_F^2 = \sum_i \sum_j A_{ij}^2$.
-
- Also called deep matrix factorization.
 - Initially used to better understand neural networks.
 - Numerous landscape results, especially between 2016-2022.

Case $L = 1$ (One-layer)

$$\underset{\mathbf{W}_1}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- Convex problem!
- Explicit form of a solution (often costly to compute).

Case $L = 1$ (One-layer)

$$\underset{\mathbf{W}_1}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- Convex problem!
- Explicit form of a solution (often costly to compute).

Case $L = 2$ (two-layer network)

$$\underset{\substack{\mathbf{W}_1 \in \mathbb{R}^{d_2 \times d_1} \\ \mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}}}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- If $\mathbf{X} \mathbf{X}^T$ full rank, there are no spurious minima.
- If $d_2 \geq \max\{d_1, d_3\}$, the optimal value is 0!

Bad example for $L = 3$

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{1 \times 2}, \mathbf{W}_2 \in \mathbb{R}, \mathbf{W}_3 \in \mathbb{R}^{2 \times 1}}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\|_F^2$$

Bad example for $L = 3$

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{1 \times 2}, \mathbf{W}_2 \in \mathbb{R}, \mathbf{W}_3 \in \mathbb{R}^{2 \times 1}}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\|_F^2$$

- The point $\left(\begin{bmatrix} 1 & 0 \end{bmatrix}, 0, \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right)$ is a local, non-global minimum!
- Due to intermediate dimensions.

Bad example for $L = 3$

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{1 \times 2}, \mathbf{W}_2 \in \mathbb{R}, \mathbf{W}_3 \in \mathbb{R}^{2 \times 1}}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\|_F^2$$

- The point $\left(\begin{bmatrix} 1 & 0 \end{bmatrix}, 0, \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right)$ is a local, non-global minimum!
- Due to intermediate dimensions.

A positive result (informal)

$$\underset{\mathbf{W}_1 \in \mathbb{R}^{d \times 2}, \mathbf{W}_2 \in \mathbb{R}^{d \times d}, \mathbf{W}_3 \in \mathbb{R}^{2 \times d}}{\text{minimize}} \quad \frac{1}{2} \left\| \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\|_F^2$$

If $d \geq 2$ (overparameterized regime), no spurious minima!

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_L}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- $\mathbf{W}_i \in \mathbb{R}^{d_{i+1} \times d_i}$.
- $\mathbf{X} \in \mathbb{R}^{d_1 \times d_0}$, $\mathbf{Y} \in \mathbb{R}^{d_{L+1} \times d_0}$.
- $\|\mathbf{A}\|_F^2 = \sum_i \sum_j A_{ij}^2$.

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_L}{\text{minimize}} \frac{1}{2} \|\mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2$$

- $\mathbf{W}_i \in \mathbb{R}^{d_{i+1} \times d_i}$.
 - $\mathbf{X} \in \mathbb{R}^{d_1 \times d_0}$, $\mathbf{Y} \in \mathbb{R}^{d_{L+1} \times d_0}$.
 - $\|\mathbf{A}\|_F^2 = \sum_i \sum_j A_{ij}^2$.
-
- Full characterization of the landscape possible (Achour et al '22).
 - IF all dimensions are equal and $\mathbf{X}\mathbf{X}^T$ full rank, no spurious local minima!

Basic block in optimization

- Derivatives (more on that next week).
- Convexity and strong convexity.

Both help characterize solutions of a problem!

Towards the nonconvex case

- Challenge: Presence of spurious minima.
- Overparameterization helps (often the case in ML)!
- Still a lot to be understood (optional course, internships?).

Basic block in optimization

- Derivatives (more on that next week).
- Convexity and strong convexity.

Both help characterize solutions of a problem!

Towards the nonconvex case

- Challenge: Presence of spurious minima.
- Overparameterization helps (often the case in ML)!
- Still a lot to be understood (optional course, internships?).

Summary

Basic block in optimization

- Derivatives (more on that next week).
- Convexity and strong convexity.

Both help characterize solutions of a problem!

Towards the nonconvex case

- Challenge: Presence of spurious minima.
- Overparameterization helps (often the case in ML)!
- Still a lot to be understood (optional course, internships?).

For now

- Material available online by tomorrow (with corrections if needed).
- Questions are always welcome.