# Session 12: Gradient Descent

Optimization and Computational Linear Algebra for Data Science

BYOD  & VANDENBHERGE - CHAP9

FRANCIS BACH ⟶ CHAP5 .

Marylou Gabrié (based on material by Léo Miolane)

# Contents

# Clarification about saddle points

- A critical is always either a local minimum or a local maximum, a saddle point.

- **Definitions:**
  - A critical point $x^*$ is a local extrema for a small $\delta > 0$ for any $x \in B(x^*, \delta)$, $f(x)$ is bigger/smaller than $f(x^*)$.
  - If a critical point not a local extrema, then it is a saddle point.

- **Caracterizations ( sufficient but not necessary conditions):** Examine Hessian $H_f(x^*)$:
  - is positive definite $\Rightarrow$ local minimum.
  - has strictly positive and strictly negative eigenvalues $\Rightarrow$ saddle
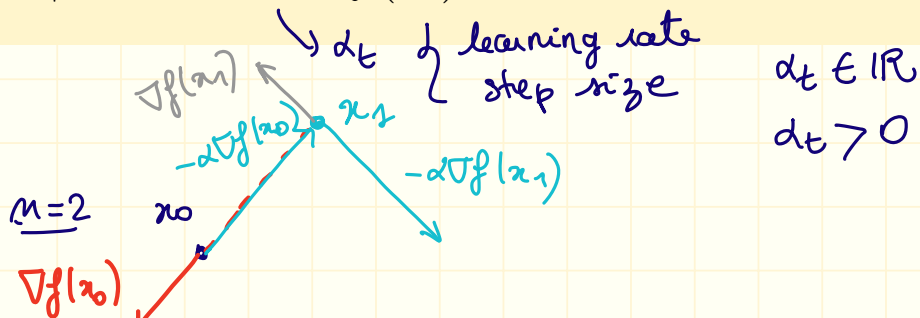
# 1. Gradient descent

# Gradient descent algorithm

**Goal:** minimize a differentiable function $f : \mathbb{R}^n \to \mathbb{R}$.

Starting from a point $x_0 \in \mathbb{R}^n$, perform the updates:

$$x_{t+1} = x_t - \alpha_t \nabla f(x_t).$$

CAUCHY ($\sim$1850)

$\nabla f(x_1)$

$\searrow \alpha_t$ } learning rate
step size

$\alpha_t \in \mathbb{R}$
$\alpha_t > 0$

$-\alpha \nabla f(x_0)$ $x_1$

$-\alpha \nabla f(x_1)$

$n = 2$   $x_0$

$\nabla f(x_0)$

IDEA: $f(x_t + h) = f(x_t) + h \cdot \nabla f(x_t)$
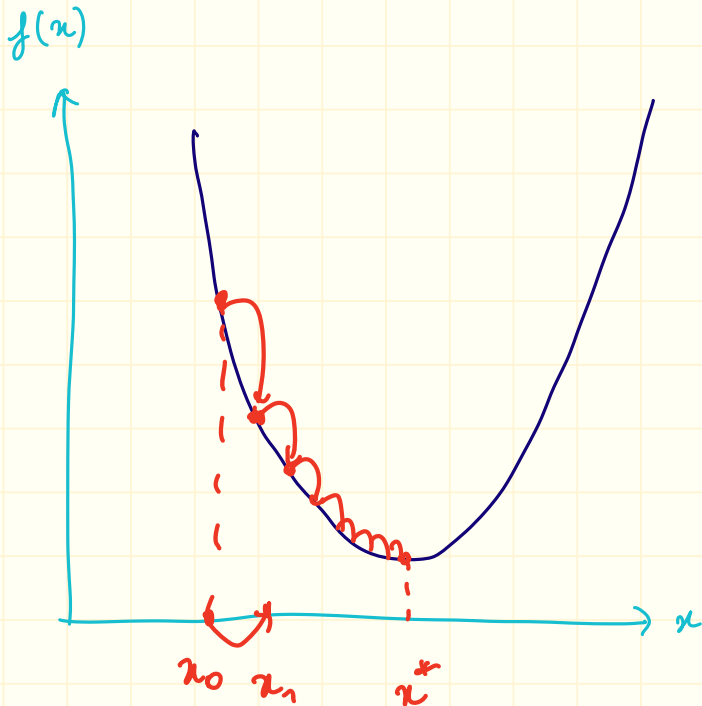
$h = -\alpha \nabla f(x_t)$

$f(x_{t+1}) \simeq f(x_t) - \alpha \underbrace{\nabla f(x_t) \cdot \nabla f(x_t)}_{\|\nabla f(x_t)\|^2} \leqslant 0$
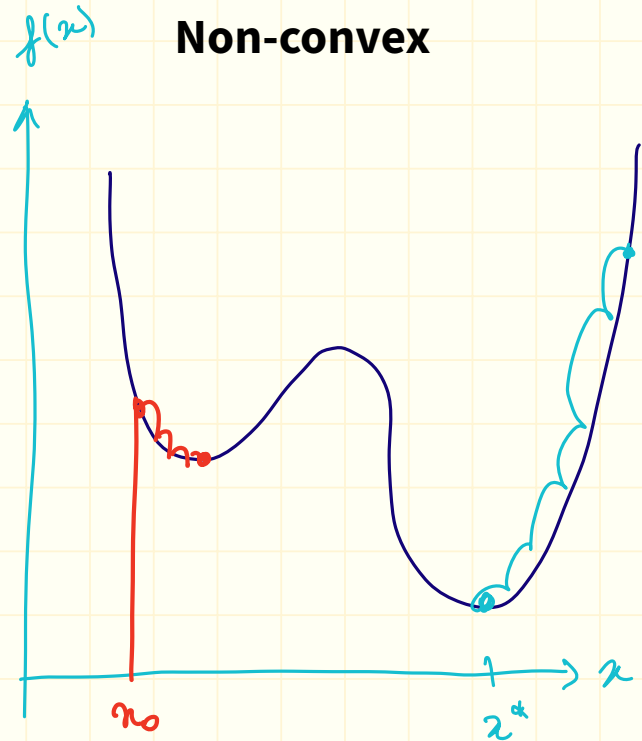
for $\alpha$ small.
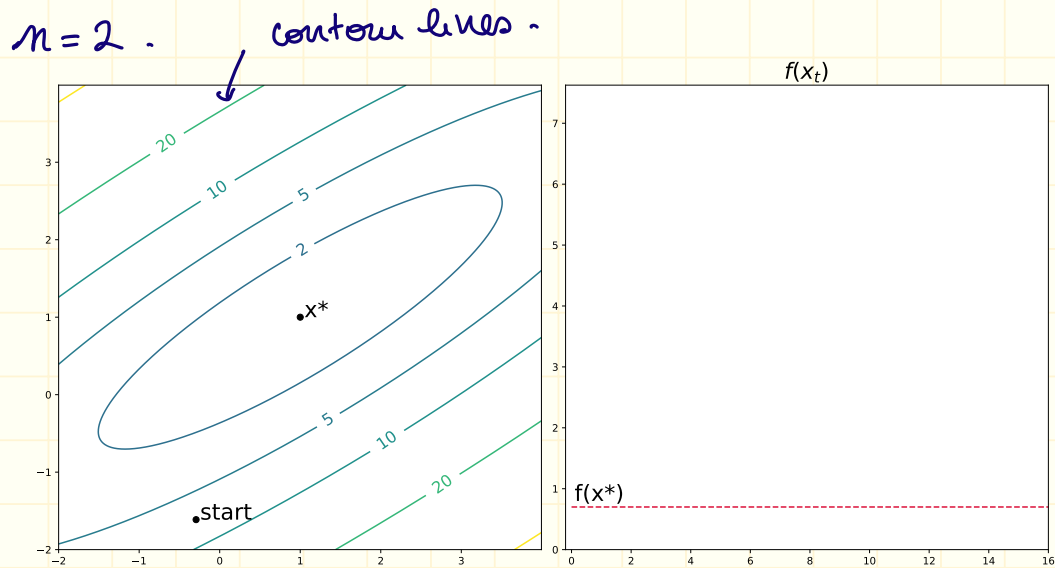
# Convex vs non-convex
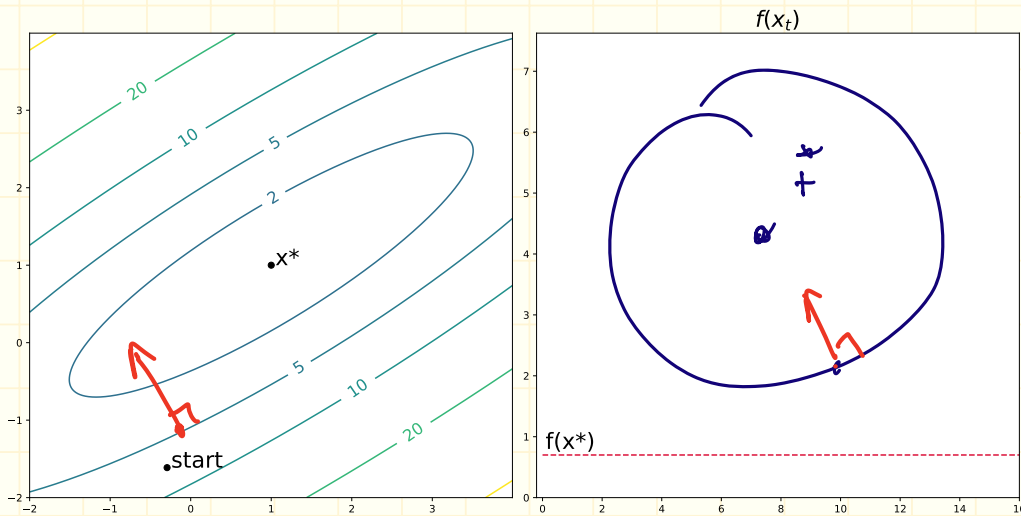


$M = 1$

highly dependent of the initialization —

**Convex**

$f(x)$

$x_0$  $x_1$  $x^*$  $x$

**Non-convex**

$f(x)$

$x_0$  $x^*$  $x$

# Numerical observations



n = 2.  contour lines.

x*

start

f(x_t)

f(x*)

1. Gradient descent

- If the step size $\alpha$ is small enough, gradient descent converges to $x^\star$ **but** this may take a while.

- If the step size $\alpha$ is large, gradient descent moves faster **but** it may oscilate or even diverge.

- The convergence is faster when the eigenvalues of the Hessian $H_f$ are of close to each other.

# 2. Convergence analysis for convex functions

.

## Definition

Given $L, \mu > 0$, we say that a twice-differentiable convex function
$f : \mathbb{R}^n \to \mathbb{R}$ is

↳ function not to be infinitely steep.

- $L$-smooth if for all $x \in \mathbb{R}^n$, $\lambda_{\max}(H_f(x)) \leq L$.
- $\mu$-strongly convex if for all $x \in \mathbb{R}^n$, $\lambda_{\min}(H_f(x)) \geq \mu$.   Hw 9 or 10.

Remark: if $f$ convex if $\begin{cases} L \text{ smooth} \\ \mu \text{ strongly convex} \end{cases}$, then for

$h$ small:

$$\frac{1}{2} h^T H_f(x) h$$

$$f(x) + \nabla f(x) \cdot h + \frac{\mu}{2} \|h\|^2 \leq f(x+h) \leq f(x) + \nabla f(x) \cdot h + \frac{L}{2} \|h\|^2$$

## Proposition

Assume that $f$ is convex, $L$-smooth and admits a global minimizer $x^\star \in \mathbb{R}^n$. Then, gradient descent with constant step size $\alpha_t = 1/L$ verifies:

$\rightarrow$ initial distance to the solution

$$f(x_t) - f(x^\star) \le \frac{2L\|x_0 - x^\star\|^2}{t + 4}. \qquad = O\left(\frac{1}{t}\right)$$

how close in terms of function value we are after $t$ step of GD

Why step $\alpha_t = \frac{1}{L}$: $\quad f(x_t + h) \le f(x_t) + \nabla f(x_t) \cdot h + \frac{L}{2} \|h\|^2$

$\underbrace{\qquad\qquad\qquad\qquad\qquad}$ min w.r.t $h$

$$h^* = -\frac{1}{L} \nabla f(x_t)$$

$$= \alpha$$

## Theorem

Assume that $f$ is convex, $L$-smooth and $\mu$-strongly convex. Then, gradient descent with constant step size $\alpha_t = 1/L$ verifies:

$$f(x_t) - f(x^\star) \leq \left(1 - \frac{\mu}{L}\right)^t (f(x_0) - f(x^\star)). = O\left(e^{-\frac{\mu}{L}t}\right)$$

$1/K$

distance to solution at initialization.

Remark: $K = \dfrac{L}{\mu} \gg \dfrac{\max_x \lambda_{max} H_f(x)}{\min_x \lambda_{min} H_f(x)} \gg 1$. CONDITION NUMBER

$\searrow$ speed of convergence if $K \nearrow$

# Proof

Recall that $f(x+h) \leq f(x) + \nabla f(x) \cdot h + \frac{L}{2} \|h\|^2$

Apply this for: $x = x_t$ and $h = -\frac{1}{L} \nabla f(x_t)$

$\Rightarrow f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2$ ①

exercice

By strong convexity: $f(x_t) - f(x^*) \leq \frac{1}{2\mu} \|\nabla f(x_t)\|^2$ ②

↗ global minimum

① + ②: $f(x_{t+1}) - f(x^*) \leq \overbrace{f(x_t) - f(x^*)}^{① - f(x^*)} - \frac{1}{2L} \|\nabla f(x_t)\|^2$

② $\|\nabla f(x_t)\|^2 \geq 2\mu (f(x_t) - f(x^*))$

$-\frac{1}{2L} \|\nabla f(x_t)\|^2 \leq -\frac{\mu}{L} (f(x_t) - f(x^*))$

$+$ ②

$\leq (f(x_t) - f(x^*)) \left(1 - \frac{\mu}{L}\right)^2$

④ FORMAL WRITING OF THE INDUCTION

Backtracking line search

decrease by at least $\|\alpha/2 \nabla f(x_t)\|^2$
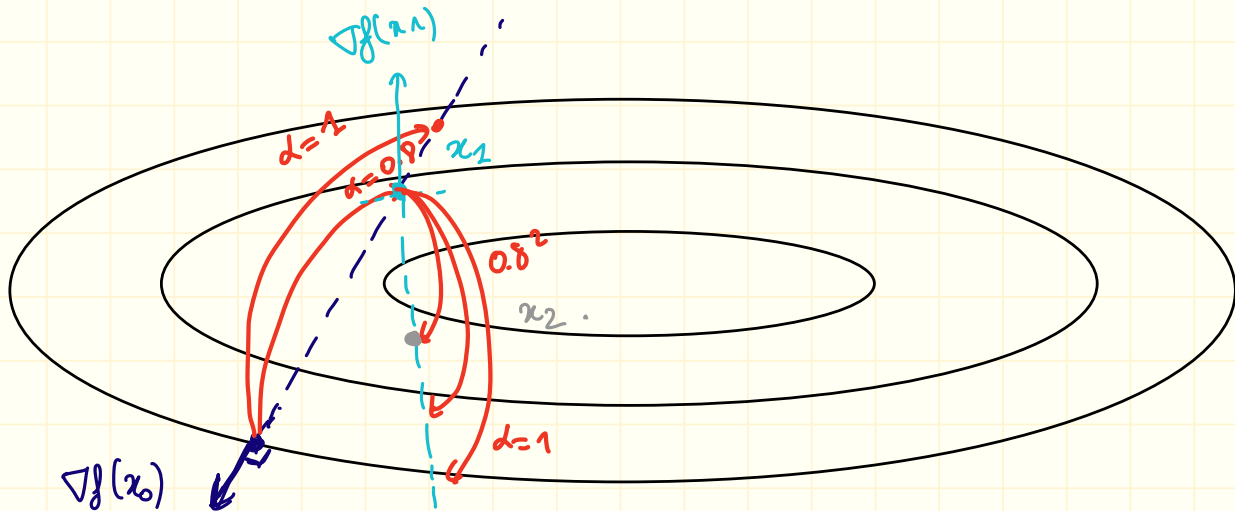
Start with $\alpha = 1$ and while

$$f(x_t - \alpha \nabla f(x_t)) \geq f(x_t) - \frac{\alpha}{2}\|\nabla f(x_t)\|^2,$$

update let's say $\alpha = 0.8\alpha$.

# 3. Improvements

# Issues with gradient descent

When the condition number $\kappa = L/\mu$ is large:

$\rightarrow$ or learning rate

1.  the norm $\|\nabla f(x)\|$ is sometimes too small.

    $\rightarrow$ gradient descent steps are too small.

2.  The vector $-\nabla f(x)$ does « not really » points towards the minimizer $x^\star$.
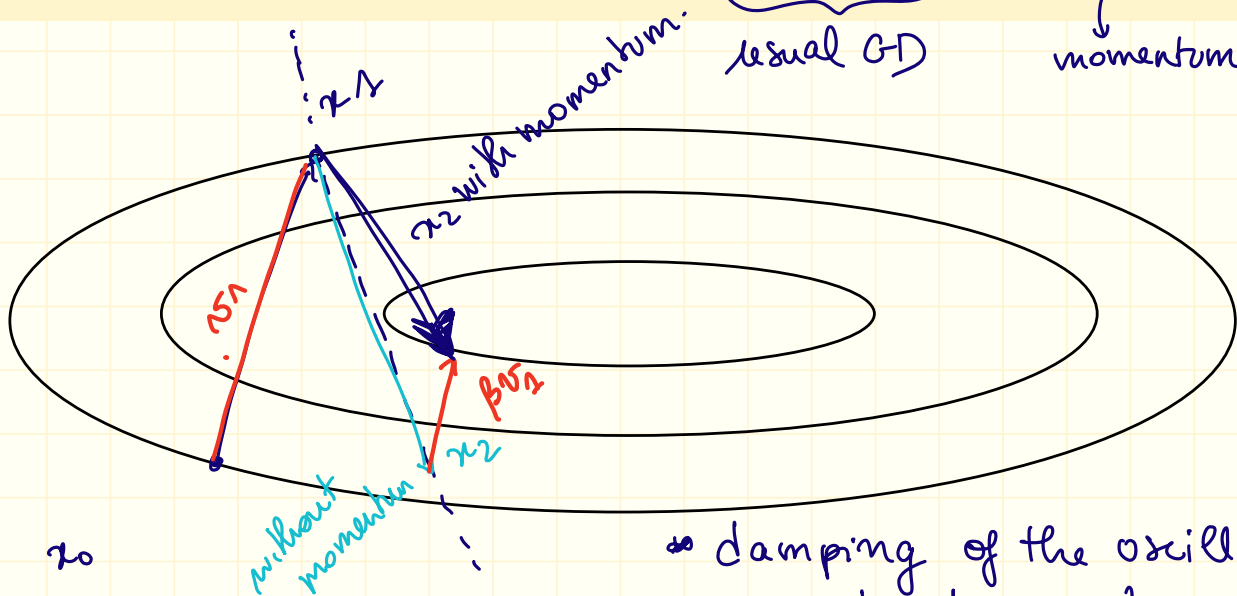
    $\rightarrow$ gradient descent oscilates.

**Idea:** mimic the trajectory of an « heavy ball » that goes down the slope:

$$x_{t+1} = x_t + v_t \qquad \text{where} \qquad v_t = -\alpha_t \nabla f(x_t) + \beta_t v_{t-1}.$$

parameter $\beta$.

$\underbrace{-\alpha_t \nabla f(x_t)}_{\text{usual GD}}$

momentum

$x_1$

$x_2$ with momentum.

$v_1$

$\beta v_1$

$x_2$

without momentum

$x_0$

※ damping of the oscillations
* promotes direction towards minimum-

# Newton's method

Assume that $f$ is $\mu$-strongly convex and $L$-smooth.

Newton's method perform the updates:

$$x_{t+1} = x_t - H_f(x_t)^{-1}\nabla f(x_t).$$

IDEA: Optimizing the learning rate by considering the second order Taylor expansion.

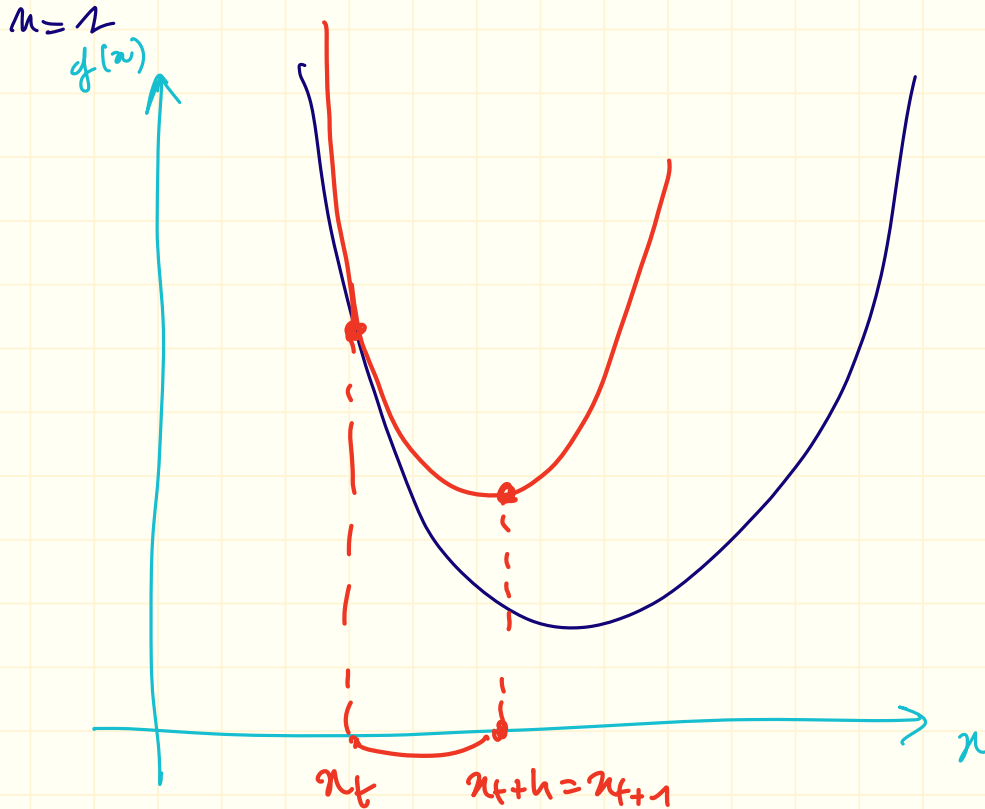$$f(x_{t+1}) = f(x_t + h) = f(x_t) + h \cdot \nabla f(x_t) + \frac{1}{2} h^T H_f(x_t) h.$$
$$= Q(h)$$

• $Q$ is convex    $H_Q(h) = H_f(x_t)$    $H_Q(h)$ is PSD.

• $\nabla Q(h) = 0$  $\Rightarrow$  $\nabla f(x_t) + H_f(x_t) h = 0$

   $\Rightarrow$  $h = - H_f^{-1}(x_t) \nabla f(x_t).$

$m = 1$

$f(x)$

$x$

$x_t$   $x_t + h = x_{t+1}$

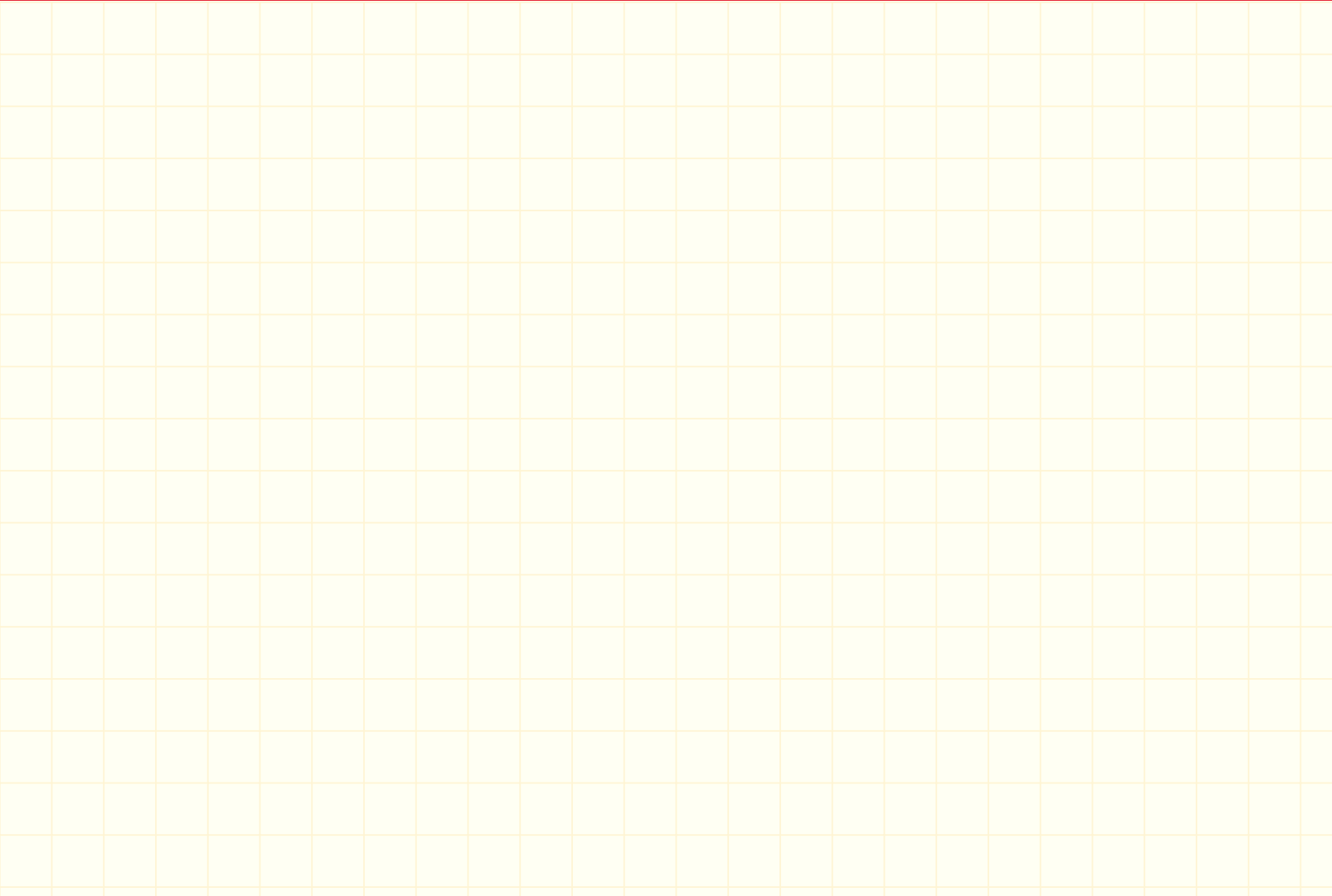# Advantages and drawbacks

▰ Extremly fast there exists $C, \rho > 0$ such that

$$\|x_t - x^\star\|^2 \le C e^{-\rho 2^t}.$$

▰ Computationally expensive: requires $\sim n^3$ operations to compute the inverse of the $n \times n$ matrix $H_f(x_t)$.

▰ In non-convex setting, Newton's method gets attracted by any critical points (which could be saddle points/maximas...).

**Quasi-Newton methods**: try to approximate $H_f(x_t)^{-1}$ by matrices $B_t$ that are easier to compute.

# Questions?

# Questions?