# Session 12: Gradient descent

Optimization and Computational Linear Algebra for Data Science

Léo Miolane

# Contents

# Gradient descent

# Gradient descent algorithm

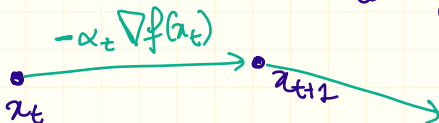**Goal:** minimize a differentiable function $f : \mathbb{R}^n \to \mathbb{R}$.

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

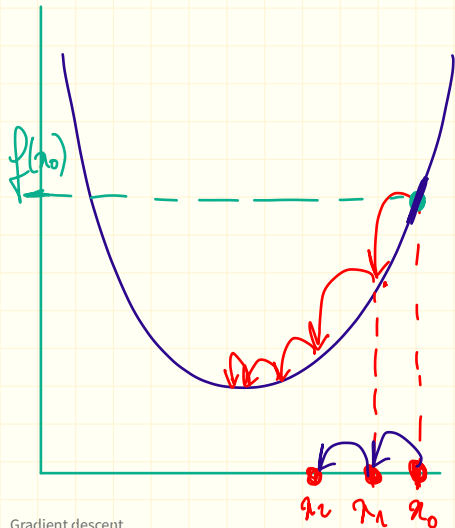$$x_{t+1} = x_t - \alpha_t \nabla f(x_t).$$

"step size"

"learning rate"

$$-\alpha_t \nabla f(x_t)$$

$x_t$     $x_{t+1}$

IDEA: $f(x_t + h) \simeq f(x_t) + h \cdot \nabla f(x_t)$

$$f(x_{t+1}) \simeq f(x_t) - \alpha_t \|\nabla f(x_t)\|^2$$

$\leq 0$

Cauchy (~1850)

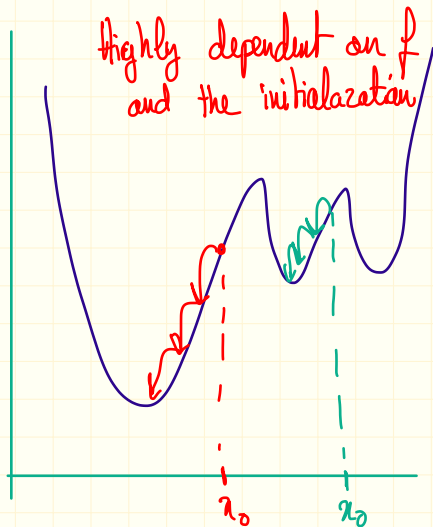for $\alpha_t$ small enough.

# Convex vs non-convex

# Numerical observations

- 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.

# 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

▶ $L$-smooth if for all $x \in \mathbb{R}^n$, $\lambda_{\max}(H_f(x)) \leq L$.

▶ $\mu$-strongly convex if for all $x \in \mathbb{R}^n$, $\lambda_{\min}(H_f(x)) \geq \mu$.

Remark: if $f$ is $\begin{vmatrix} L\text{-smooth} \\ \mu\text{- strongly convex} \end{vmatrix}$ then:

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

## 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:

$$f(x_t) - f(x^\star) \leq \frac{2L\|x_0 - x^\star\|^2}{t + 4} \cdot \subseteq \frac{\text{Constant}}{t}$$

Why step size $\alpha_t = \frac{1}{L}$ ?

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

this is minimal for $\boxed{h = -\frac{1}{L}\nabla f(x_t)}$

$$\rightarrow x_{t+1} = x_t + h = x_t - \frac{1}{L}\nabla f(x_t)$$

# $L$-smooth + $\mu$-strongly cvx functions

**Theorem**

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

$$f(x_t) - f(x^\star) \leq \underbrace{\left(1 - \frac{\mu}{L}\right)^t}_{\leq\, e^{-\frac{\mu}{L}t}} \underbrace{(f(x_0) - f(x^\star))}_{\text{Constant}}.$$

<u>Remark</u>: • GD with step size $\alpha_t = \dfrac{1}{L}$ is "adaptive" to strong convexity"

• The quantity $K = \dfrac{L}{\mu} \geqslant 1$ is called the "condition number"

$K \nearrow$ the convergence speed $\searrow$

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

Apply it for $x = x_t$, $h = -\frac{1}{L} \nabla f(x_t)$

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

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

Combining the two inequalities:

$$f(x_{t+1}) - f(x^*) \leq f(x_t) - f(x^*) - \frac{\mu}{L} \left( f(x_t) - f(x^*) \right)$$
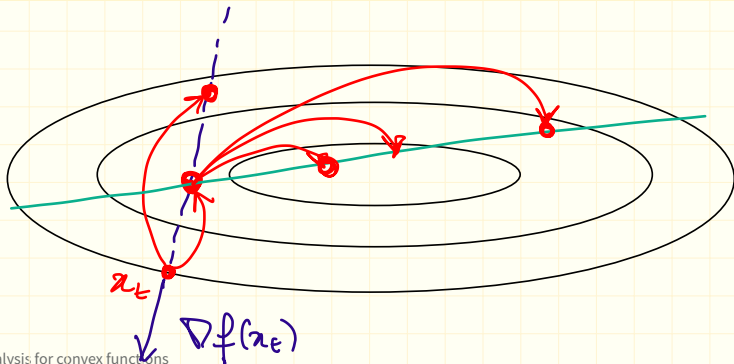
$$= \left( f(x_t) - f(x^*) \right) \left( 1 - \frac{\mu}{L} \right)$$

# Choosing the step size

Backtracking line search

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 $\alpha = 0.8\alpha$.
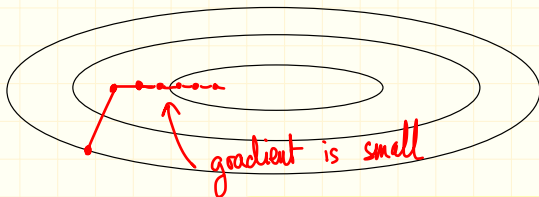


$x_t$

$\nabla f(x_t)$

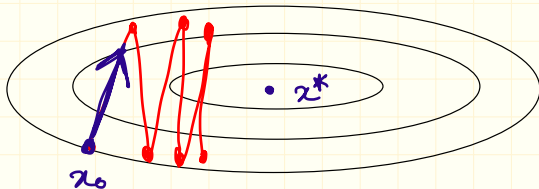# Improvements

# Issues with gradient descent

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

1. the norm $\|\nabla f(x)\|$ is sometimes too small.
   $\rightarrow$ gradient descent steps are too small.



gradient is small

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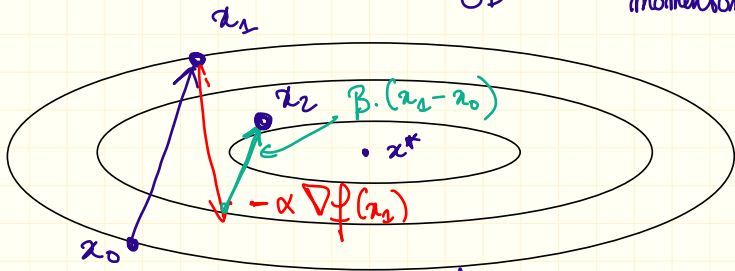   $\rightarrow$ gradient descent oscilates.



$x^\star$

$x_0$

# Gradient descent + momentum

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

$$x_{t+1} = x_t + v_t \qquad \text{where} \quad v_t = \underbrace{-\alpha_t \nabla f(x_t)}_{GD} + \underbrace{\boxed{\beta_t} v_{t-1}}_{\text{momentum}}.$$

$$x_t - x_{t-1}$$



$x_1$

$x_2$    $\beta \cdot (x_1 - x_0)$

$\cdot \; x^*$

$-\alpha \nabla f(x_1)$

$x_0$

Momentum damps the oscillations + accumulate momentum in the horizontal direction.

# Newton's method

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

Newton's method perform the updates:

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

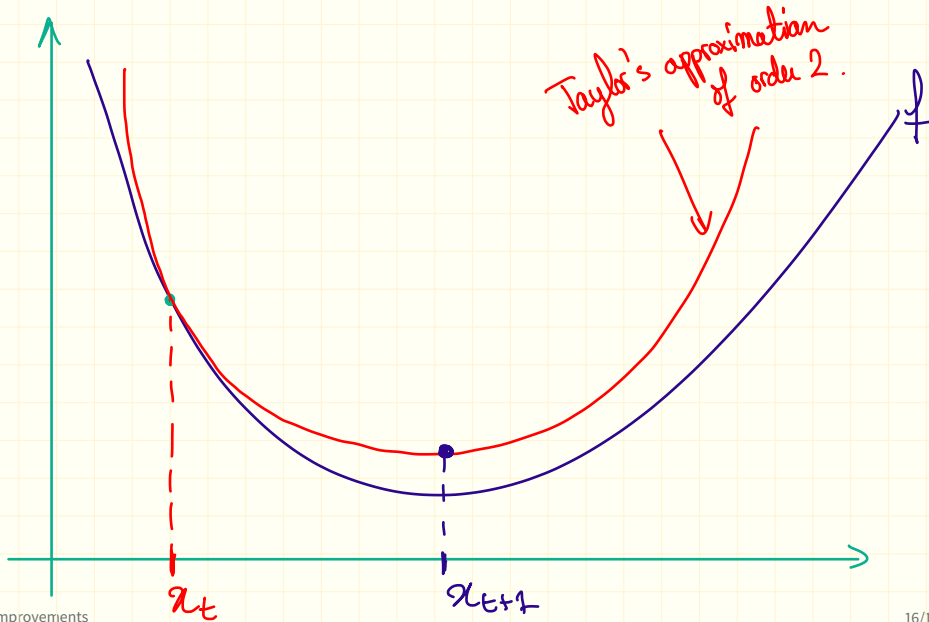*eigenvalues* $\geq \mu > 0$

*hence invertible!*

**IDEA:** $f(x_t + h) \simeq f(x_t) + h \cdot \nabla f(x_t) + \frac{1}{2} h^T H_f(x_t) h$.

$\overset{def}{=} Q(h)$, minimal for $h = -H_f(x_t)^{-1}\nabla f(x_t)$

Proof

- $Q$ is convex. $\left( H_Q(h) = H_f(x_t) \leftarrow \text{PSD} \right)$

- let's solve $\nabla Q(h) = 0$ : $\nabla f(x_t) + H_f(x_t) h = 0$

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

Taylor's approximation of order 2.

$f$

$x_t$

$x_{t+1}$

# Advantages and drawbacks

- Extremly fast there exists $C, \rho > 0$ such that
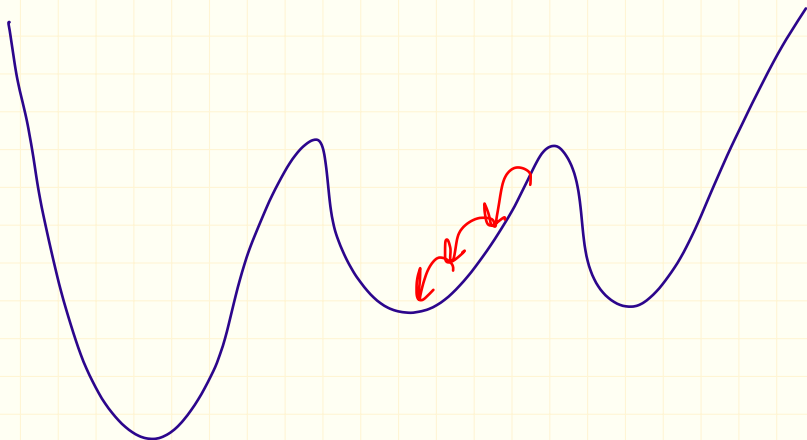
$$\|x_t - x^\star\|^2 \leq Ce^{-\rho 2^t}.$$

$$e^{-\rho 2^t}$$

- 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.

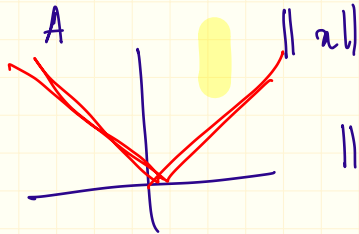$$x_{t+1} = x_t - B_t \nabla f(x_t)$$

# Questions?

$$A = U \Sigma V^T$$

$$\min_{x \in \mathbb{R}^n} \|Ax - y\|^2 = \min_{u \in \text{Im}(A)} \|u - y\|^2$$

$$\rightarrow x^* = A^+ y$$

$$u^* = A A^+ y$$

$$u^* = \begin{pmatrix} | & & | \\ u_1 & \cdots & u_r \\ | & & | \end{pmatrix} \begin{pmatrix} - u_1 - \\ - u_r - \end{pmatrix} y$$

$A$      $\|x\|$

$$\|x\| = |x)$$