# Tutorial on Gradient methods for non-convex problems

### Part 1

Guillaume Garrigos – November 28th – ENS





• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?

• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$x^{k+1} = x^k - \lambda_k \nabla f(x^k)$$

#### Proposition

- Let  $0 \ll \lambda_k \ll 2/L$ , then:
- 1)  $f(x^k)$  is decreasing
- 2) if  $x^{k_n} \to x_\infty$  then  $\nabla f(x_\infty) = 0$
- 3) **Isolated** local minima are attractive

[Pro 1.2.3, 1.2.5 & Ex. 1.2.18] Bertsekas, Nonlinear Programming, 1999.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$x^{k+1} = x^k - \lambda_k \nabla f(x^k)$$

### Proposition

- Let  $0 \ll \lambda_k \ll 2/L$ , then:
- 1)  $f(x^k)$  is decreasing

2) if 
$$x^{k_n} \to x_\infty$$
 then  $\nabla f(x_\infty) = 0$ 

3) Isolated local minima are attractive

 $x^k$  can have no limit !!

No convergence  $\neq$  Lack of regularity, but rather wildness

[Pro 1.2.3, 1.2.5 & Ex. 1.2.18] Bertsekas, Nonlinear Programming, 1999.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

 $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 



[Ex. 3] Palis, de Melo, Geometric Theory of Dynamical Systems: An Introduction, 1982. H.B.Curry, The method of steepest descent for nonlinear minimization problems, 1944.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 



 $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$  $\dot{x}(t) + \nabla f\bigl(x(t)\bigr) = 0$ 

[Ex. 3] Palis, de Melo, Geometric Theory of Dynamical Systems: An Introduction, 1982. H.B.Curry, The method of steepest descent for nonlinear minimization problems, 1944.

• A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$ 

- A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$ 
  - It is a classic result that `Finite Length' implies convergence
  - Converse is not true (but tricky):

$$x^n \coloneqq \sum_{k=1}^n \frac{(-1)^k}{k} \to -\log(2) \text{ but } \sum \left| x^{n+1} - x^n \right| = \sum \frac{1}{n}$$

- A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$
- Length is invariant up to a reparametrization in time

- A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$
- Length is invariant up to a reparametrization in time
- We have a natural diffeomorphism  $f \circ x : [0, \infty[ \rightarrow] s_{\infty}, s_0]$  where  $s_0 = f(x_0)$  and  $s_{\infty} = \lim_{\infty} f(x(t))$



- A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$
- Length is invariant up to a reparametrization in time
- We have a natural diffeomorphism  $f \circ x : [0, \infty[ \rightarrow] s_{\infty}, s_0]$  where  $s_0 = f(x_0)$  and  $s_{\infty} = \lim_{\infty} f(x(t))$
- With s = f(x(t)) we can define  $y(s) = x((f \circ x)^{-1}(s))$  s.t.

$$\dot{y}(s) = \nabla f(y(s)) \|\nabla f(y(s))\|^{-2}$$



- A sufficient condition for x(t) to converge is  $\int_0^\infty \|\dot{x}(t)\| dt < \infty$
- Length is invariant up to a reparametrization in time
- We have a natural diffeomorphism  $f \circ x : [0, \infty[ \rightarrow] s_{\infty}, s_0]$  where  $s_0 = f(x_0)$  and  $s_{\infty} = \lim_{\infty} f(x(t))$
- With s = f(x(t)) we can define  $y(s) = x((f \circ x)^{-1}(s))$  s.t.  $\dot{y}(s) = \nabla f(y(s)) \| \nabla f(y(s)) \|^{-2}$
- So the length becomes  $\int_{s_{\infty}}^{s_{0}} \frac{1}{\|\nabla f(y(s))\|} ds$ Finite interval ! Ignore  $\nabla f(y(s)) = 0$

• How to upper bound  $\int_0^\infty \|\dot{x}(t)\| dt = \int_{s_\infty}^{s_0} \frac{1}{\|\nabla f(y(s))\|} ds$  ?

- How to upper bound  $\int_0^\infty \|\dot{x}(t)\| dt = \int_{s_\infty}^{s_0} \frac{1}{\|\nabla f(y(s))\|} ds$ ?
- ``Naive'' hypothesis:  $\|\nabla f(y)\| \ge C$  i.e. sharpness

- How to upper bound  $\int_0^\infty \|\dot{x}(t)\| dt = \int_{s_\infty}^{s_0} \frac{1}{\|\nabla f(y(s))\|} ds$ ?
- ``Naive'' hypothesis:  $\|\nabla f(y)\| \ge C$  i.e. sharpness



- How to upper bound  $\int_0^\infty \|\dot{x}(t)\| dt = \int_{s_\infty}^{s_0} \frac{1}{\|\nabla f(y(s))\|} ds$ ?
- ``Naive'' hypothesis:  $\|\nabla f(y)\| \ge C$  i.e. sharpness
- ``Smart'' hypothesis:  $\frac{1}{\|\nabla f(y(s))\|} \le \varphi'(s)$  with  $\varphi \ge 0$ ,  $\varphi \uparrow$

so the length is  $\leq \varphi(s_0) - \varphi(s_\infty) \leq \varphi(s_0)$ 

- How to upper bound  $\int_0^\infty \|\dot{x}(t)\| dt = \int_{s_\infty}^{s_0} \frac{1}{\|\nabla f(y(s))\|} ds$  ?
- ``Naive'' hypothesis:  $\|\nabla f(y)\| \ge C$  i.e. sharpness
- ``Smart'' hypothesis:  $\frac{1}{\|\nabla f(y(s))\|} \leq \varphi'(s)$  with  $\varphi \geq 0$ ,  $\varphi \uparrow$  so the length is  $\leq \varphi(s_0) \varphi(s_\infty) \leq \varphi(s_0)$
- In other words  $\varphi'(f(x(t))) \| \nabla f(x(t)) \| \ge 1$  i.e.  $\varphi \circ f$  is sharp:



 $\|\nabla(\varphi \circ f)(x)\| \ge 1$ 

# The Łojasiewicz property

#### Definition

We say that f is **Lojasiewicz** at a critical point  $x^*$  if

$$\varphi'(f(x) - f(x^*)) \|\nabla f(x)\| \ge 1,$$

- with  $\varphi: [0, \infty[ \rightarrow [0, \infty[ \text{ s.t. } \varphi(0) = 0, \varphi \uparrow, \varphi \text{ concave}$
- for all  $x \in \{ x' \in \mathbb{B}(x^*, \delta) \mid f(x^*) < f(x') < f(x^*) + r \}$

#### Definition

- f is Łojasiewicz if it is Łojasiewicz at every critical point
- f is p-Łojasiewicz if it is Łojasiewicz at every critical point with  $\varphi(s) \simeq s^{1/p}$ :

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $C_L^{1,1}$ 

 $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

### Theorem (convergence)

Let f be **Lojasiewicz** and  $\lambda_k \in [0,2/L[.$ 

If  $x^k$  is **bounded**, then it converges to some critical point  $x^{\infty}$ .

#### Theorem (capture)

Let f be **Łojasiewicz** and  $\lambda_k \in [0,2/L[$ . For every  $x^* \in argmin f$ , if  $x^0 \sim x^*$  then  $x^k$  converges to  $x^{\infty} \in argmin f$ .

Łojasiewicz. Sur les trajectoires du gradient d'une fonction analytique, 1984.

Absil, Mahony, Andrews. Convergence of the Iterates of Descent Methods for Analytic Cost Functions, 2005.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L$ 

 $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

$$\varphi\left(f(x^k) - f(x^*)\right) - \varphi\left(f(x^{k+1}) - f(x^*)\right)\right)$$

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

$$\varphi \left( f(x^k) - f(x^*) \right) - \varphi \left( f(x^{k+1}) - f(x^*) \right)$$
  
 
$$\ge \varphi'(f(x^k) - f(x^*))(f(x^k) - f(x^{k+1})) \text{ because } \varphi \text{ concave}$$

 $f: \mathbb{R}^n \to \mathbb{R} \text{ is of class } \mathbb{C}^{1,1}_L$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

$$\varphi\left(f(x^{k}) - f(x^{*})\right) - \varphi\left(f(x^{k+1}) - f(x^{*})\right)\right)$$

$$\geq \varphi'(f(x^{k}) - f(x^{*}))(f(x^{k}) - f(x^{k+1})) \quad \text{because } \varphi \text{ concave}$$

$$\geq \varphi'\left(f(x^{k}) - f(x^{*})\right)c_{\lambda,L} \left\|x^{k+1} - x^{k}\right\|^{2} \quad \text{with Descent Lemma}$$

 $f: \mathbb{R}^n \to \mathbb{R} \text{ is of class } C_L^{1,1}$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

Sketch of proof : show that  $\varphi'(s) \ge \|\dot{x}(t)\|$ 

$$\begin{split} \varphi\left(f(x^{k}) - f(x^{*})\right) &- \varphi\left(f(x^{k+1}) - f(x^{*})\right)\right) \\ \geq \varphi'(f(x^{k}) - f(x^{*}))(f(x^{k}) - f(x^{k+1})) \quad \text{because } \varphi \text{ concave} \\ \geq \varphi'\left(f(x^{k}) - f(x^{*})\right)c_{\lambda,L} \left\|x^{k+1} - x^{k}\right\|^{2} \quad \text{with Descent Lemma} \\ &= \varphi'\left(f(x^{k}) - f(x^{*})\right)C_{\lambda,L} \left\|x^{k+1} - x^{k}\right\| \left\|\nabla f(x^{k})\right\| \end{split}$$

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

Sketch of proof : show that  $\varphi'(s) \ge \|\dot{x}(t)\|$ 

$$\begin{split} \varphi\left(f(x^{k}) - f(x^{*})\right) &- \varphi\left(f(x^{k+1}) - f(x^{*})\right)\right) \\ \geq \varphi'(f(x^{k}) - f(x^{*}))(f(x^{k}) - f(x^{k+1})) \quad \text{because } \varphi \text{ concave} \\ \geq \varphi'\left(f(x^{k}) - f(x^{*})\right)c_{\lambda,L} \left\|x^{k+1} - x^{k}\right\|^{2} \quad \text{with Descent Lemma} \\ &= \varphi'\left(f(x^{k}) - f(x^{*})\right)C_{\lambda,L} \left\|x^{k+1} - x^{k}\right\| \left\|\nabla f(x^{k})\right\| \end{split}$$

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$   $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

Sketch of proof : show that  $\varphi'(s) \ge \|\dot{x}(t)\|$ 

 $\varphi\left(f(x^{k}) - f(x^{*})\right) - \varphi\left(f(x^{k+1}) - f(x^{*})\right)\right)$   $\geq \varphi'(f(x^{k}) - f(x^{*}))(f(x^{k}) - f(x^{k+1})) \quad \text{because } \varphi \text{ concave}$   $\geq \varphi'\left(f(x^{k}) - f(x^{*})\right)c_{\lambda,L} \|x^{k+1} - x^{k}\|^{2} \quad \text{with Descent Lemma}$   $= \varphi'\left(f(x^{k}) - f(x^{*})\right)C_{\lambda,L} \|x^{k+1} - x^{k}\| \|\nabla f(x^{k})\|$   $\geq 1 \cdot C_{\lambda,L} \|x^{k+1} - x^{k}\| \quad \text{with:}$ 



$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $C_L^{1,1}$ 

 $x^{k+1} = x^k - \lambda_k \nabla f(x^k)$ 

### Theorem (convergence)

Let f be **Lojasiewicz** and  $\lambda_k \in [0,2/L[.$ 

If  $x^k$  is **bounded**, then it converges to some critical point  $x^{\infty}$ .

#### Theorem (capture)

Let f be **Łojasiewicz** and  $\lambda_k \in [0,2/L[$ . For every  $x^* \in argmin f$ , if  $x^0 \sim x^*$  then  $x^k$  converges to  $x^{\infty} \in argmin f$ .

Łojasiewicz. Sur les trajectoires du gradient d'une fonction analytique, 1984.

Absil, Mahony, Andrews. Convergence of the Iterates of Descent Methods for Analytic Cost Functions, 2005.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

### Examples

- If f is  $\mu$ -strongly convex, then it is 2-Łojasiewicz with  $\mu = \mu$
- If f convex and  $\mu d(x, \operatorname{argmin} f)^p \le f(x) \inf f$

#### Counter-example

• There is a convex function  $f: \mathbb{R}^2 \to \mathbb{R}$  which is not Łojasiewicz

Bolte, Daniilidis, Ley, Mazet, Characterizations of Łojasiewicz inequalities [...], 2010.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

#### Theorem

Any **analytic** function is p-Łojasiewicz at its critical points.

#### Theorem

- Any **semi-algebraic** function is p-Łojasiewicz at its critical points.
- Any **o-minimal** function is Łojasiewicz.

Łojasiewicz, Ensembles semi-analytiques, 1965.

Kurdyka, On gradients of functions definable in o-minimal structures, 1998.

Bolte, Daniilidis, Lewis, Shiota, Clarke Subgradients of Stratifiable Functions, 2007.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Examples of semi-algebraic functions

• Polynomials by parts

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Examples of semi-algebraic functions

• Polynomials by parts

Theorem (``Tarski-Seidenberg'')

The class of **semi-algebraic** functions is stable under:

- addition, multiplication, division, sup, inf
- restriction, composition, inverse  $f^{-1}$
- derivative

Coste, An Introduction to O-minimal Geometry, 2000.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

#### Examples of semi-algebraic functions

- Polynomials by parts
- $\alpha \|x\|_0 + \|Ax b\|^2, \|x\|_*, \dots$

#### Theorem (``Tarski-Seidenberg'')

The class of **semi-algebraic** functions is stable under:

- addition, multiplication, division, sup, inf
- restriction, composition, inverse  $f^{-1}$
- derivative

Coste, An Introduction to O-minimal Geometry, 2000.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Counter-examples of semi-algebraic functions

Exponential/Logarithmic stuff

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

#### Counter-examples of semi-algebraic functions

Exponential/Logarithmic stuff

#### Theorem

There exists a class of functions (o-minimal structure) which:

- Includes the semi-algebraic structure
- Contains the exponential function
- Has the same stability property than the semi-algebraics
- Is also stable by integration (and resolution of 1st order ODEs)

Speissegger, The Pfaffian closure of an o-minimal structure, 1999.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Take-home message:

Virtually any function you can think about is Łojasiewicz, as long as it does not involve

 $\mathbb{R} \to \mathbb{R}$  $x \mapsto \sin(x)$ 

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Take-home message:

Virtually any function you can think about is Łojasiewicz, as long as it does not involve

 $\mathbb{R} \to \mathbb{R}$  $x \mapsto \sin(x)$ 





 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Take-home message:

Virtually any function you can think about is Łojasiewicz, as long as it does not involve

 $\mathbb{R} \to \mathbb{R}$  $x \mapsto \sin(x)$ 

So gradient descent ``always converges'' to a critical point

# The Łojasiewicz property : rates for free

$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

#### Theorem (p=2)

Let f be globally 2-Łojasiewicz,  $\lambda \in ]]0,2/L[[, and <math>x^k \rightarrow x^*]$ . Then we have **linear** convergence :

$$f(x^{k+1}) - f(x^*) \le \theta\left(f(x^k) - f(x^*)\right)$$

where  $\theta \in [0,1[$ , and  $\theta = 1 - \mu/L$  if  $\lambda = 1/L$ .

- If f strongly convex we have  $\theta = (1 \kappa)^2$  for  $\lambda = 1/L$ .
- If f is [any weak s. convex notion] we have a better  $\theta$ .
- Rates become asymptotic if local Łojasiewicz only.

Polyak, Gradient methods for the minimisation of functionals, 1963.

# The Łojasiewicz property : rates for free

$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

#### Theorem (p>2)

Let *f* be globally **p**-Łojasiewicz,  $\lambda \in ]]0,2/L[[, and <math>x^k \to x^*]$ . Then we have **sublinear** convergence :

$$f(x^k) - f(x^*) = O\left(k^{\frac{-p}{p-2}}\right)$$

• 
$$\frac{p}{p-2} \to +\infty$$
 when  $p \downarrow 2$ ;  $\frac{p}{p-2} \to 1$  when  $p \uparrow \infty$ 

• Rates are matched for 
$$f(x) = x^p$$

Attouch, Bolte, On the convergence of the proximal algorithm for nonsmooth functions [...], 2009. Chouzenoux, Pesquet, Repetti, A block coordinate variable metric forward-backward algorithm, 2014.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

Take-home message:

Virtually any function you can think about is Łojasiewicz, as long as it does not involve

 $\mathbb{R} \to \mathbb{R}$  $x \mapsto \sin(x)$ 

So gradient descent ``always converges'' to a critical point

What about **other methods**?

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

• Nonsmooth 1st-order methods : works the same

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

- Nonsmooth 1st-order methods : works the same
  - Projected gradient
  - Forward Backward
  - Douglas Rachford
  - ADMM
  - Adapts to Maximal Monotone theory (saddle point problems)

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

- Nonsmooth 1st-order methods : works the same
- Inertial (2<sup>nd</sup> order in time) methods :

$$\ddot{x}(t) + \alpha(t)\dot{x}(t) + \nabla f(x(t)) = 0$$
$$x^{k+1} = y^k - \lambda \nabla f(y^k)$$
$$y^k = x^k + \frac{1}{1+\alpha_k} (x^k - x^{k-1})$$

Bégout, Bolte, Jendoubi, On damped second-order gradient systems, 2015. + refs within! Li et al., Convergence Analysis of Proximal Gradient with Momentum for Nonconvex Optimization, 2017.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

- Nonsmooth 1st-order methods : works the same
- Inertial (2<sup>nd</sup> order in time) methods :
  - Heavy-Ball  $(\alpha(t) \equiv \alpha)$  ok
  - Nesterov ( $\alpha(t) \sim \alpha/t$ ) + Monotone OK pour p=2 global

 $\ddot{x}(t) + \alpha(t)\dot{x}(t) + \nabla f(x(t)) = 0$ 

$$x^{k+1} = y^k - \lambda \nabla f(y^k)$$
$$y^k = x^k + \frac{1}{1+\alpha_k} (x^k - x^{k-1})$$

Bégout, Bolte, Jendoubi, On damped second-order gradient systems, 2015. + refs within! Li et al., Convergence Analysis of Proximal Gradient with Momentum for Nonconvex Optimization, 2017.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

- Nonsmooth 1st-order methods : works the same
- Inertial (2<sup>nd</sup> order in time) methods : ok
- Newton-like (2<sup>nd</sup> order in space) methods: some results for trustregion methods, Landweber iterations. IDK for Newton/BFGS.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $C_L^{1,1}$ 

$$\mu(f(x) - f(x^*))^{p-1} \le \|\nabla f(x)\|^p$$

- Nonsmooth 1st-order methods : works the same
- Inertial (2<sup>nd</sup> order in time) methods : ok for Heavy-Ball
- Newton-like (2<sup>nd</sup> order in space) methods: some results for trustregion methods, Landweber iterations. IDK for Newton/BFGS.
- Stochastic methods: (under global 2-Łojasiewicz)
  - SAGA, SVRG: linear rates
  - SGD: rates O(1/k), and linear if vanishing variance
  - SVRG + monotone Nesterov: linear rates

Reddi, Hefny, Sra, Poczos, Smola, Stochastic variance reduction for nonconvex optimization, 2016. Karimi, Nutini, Schmidt, Linear Convergence of Gradient and Proximal-Gradient [...], 2016. Lei et al., Stochastic Gradient Descent for Nonconvex Learning without Bounded Gradient Assumptions, 2019.

• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?

• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?



• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

A symmetric operator

 $\dot{x}(t) + Ax(t) = 0$ 

A symmetric operator

 $\dot{x}(t) + Ax(t) = 0$ 



Stolen from Robert Ghrist' Twitter account @robertghrist



- Positive eigs. are attractive, negative eigs. are repulsive.
- Converging to the saddle point requires starting from  $E_{-1}(A)$

A symmetric operator

 $\dot{x}(t) + Ax(t) = 0$ 

#### Definition

Let  $\bar{x}$  be an equilibrium of the system. We define

$$W(\bar{x}) = \{x \mid x(t) \to \bar{x} \text{ with } x(0) = x\}$$

#### Theorem

$$W(\bar{x}) \simeq \bigoplus_{\lambda > 0} E_{\lambda}(A)$$

#### Corollary

If  $\lambda_{min}(A) < 0$ , then  $W(\bar{x})$  has Lebesgue measure 0.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^2$ 

 $\dot{x}(t) + \nabla f(x(t)) = 0$ 

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^2$ 

 $\dot{x}(t) + \nabla f(x(t)) = 0$ 

### Theorem (Stable Manifold Lemma)

 $W(\bar{x})$  is a submanifold of dimension smaller than the one of

 $\bigoplus_{\lambda>0} E_{\lambda}(\nabla^2 f(\bar{x}))$ 

### Corollary

If  $\lambda_{min}(\nabla^2 f(\bar{x})) < 0$ , then  $W(\bar{x})$  has Lebesgue measure 0.

Perron, Die stabilitätsfrage bei differentialgleichungen, 1930. Smale, Differentiable dynamical systems, 1967.



The 3 kinds of critical points:

- The local minima (e.g.  $\lambda_{min} (\nabla^2 f(\bar{x})) > 0)$  , attractive
- The strict saddles  $\lambda_{min}(\nabla^2 f(\bar{x})) < 0$ , repulsive
- The degenerated ones (they have  $\lambda_{min}(\nabla^2 f(\bar{x})) = 0$ ), ???

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^2$ 

 $\ddot{x}(t) + \alpha \dot{x}(t) + \nabla f(x(t)) = 0$ 

Goudou, Munier, The gradient and heavy ball with friction dynamical systems: the quasiconvex case, 2007.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^2$ 

 $\ddot{x}(t) + \alpha \dot{x}(t) + \nabla f(x(t)) = 0$ 

### Theorem (Stable Manifold Lemma)

 $W(\bar{x})$  is a submanifold of dimension smaller than the one of  $\bigoplus_{\lambda>0} E_{\lambda}(\nabla^2 f(\bar{x}))$ 

### Corollary

If  $\lambda_{min}(\nabla^2 f(\bar{x})) < 0$ , then  $W(\bar{x})$  has Lebesgue measure 0.

Goudou, Munier, The gradient and heavy ball with friction dynamical systems: the quasiconvex case, 2007.

$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $\mathbb{C}^{1,1}_L \cap \mathbb{C}^2$ 

$$x^{k+1} = x^k - \lambda \nabla f(x^k)$$

#### Theorem

Let  $\lambda \in ]0,1/L[$ . Then  $W(\bar{x})$  is a submanifold of dimension smaller than the one of  $\bigoplus_{\lambda>0} E_{\lambda}(\nabla^2 f(\bar{x}))$ 

### Corollary

If  $\lambda_{min}(\nabla^2 f(\bar{x})) < 0$ , then  $W(\bar{x})$  has Lebesgue measure 0.

### Corollary

If f has no degenerated critical points and is Łojasiewicz, then  $x^k$  converges a.s. to a local minima with random initialization.

Lee, Simchowitz, Jordan, Recht, Gradient Descent Converges to Minimizers, 2016.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L \cap \mathbb{C}^2$ 

 $x^{k+1} = x^k - \lambda \nabla f(x^k)$ 

It is time now for <u>examples</u>.



$$f: \mathbb{R}^n \to \mathbb{R}$$
 is of class  $\mathbb{C}^{1,1}_L \cap \mathbb{C}^2$ 

$$x^{k+1} = x^k - \lambda \nabla f(x^k)$$

#### Corollary

If f has no degenerated critical points and is Łojasiewicz, then  $x^k$  converges a.s. to a local minima with random initialization.

Some problems have no degenrated critical points, like the matrix factorization problem a.k.a. two-layer-linear-neural-network

$$\min_{X \in \mathbb{R}^{d \times r}} f(X) = \left\| X^T X - A \right\|_F^2$$

Li et al., Symmetry, saddle points, and global geometry of nonconvex matrix factorization, 2016.

 $f: \mathbb{R}^n \to \mathbb{R}$  is of class  $\mathbb{C}^{1,1}_L \cap \mathbb{C}^2$ 

$$x^{k+1} = x^k - \lambda \nabla f(x^k) + \xi_k$$

#### Corollary

If f has no degenerated critical points and is Łojasiewicz, then  $x^k$  converges a.s. to a local minima with random initialization.

#### Corollary

The above result remains true for the noisy gradient method.

The noise here must be isotropic!

Not the case for SGD (proportional to eigenvalues), but proof can be adapted for RKHS learning with a loss s.t.  $|\ell''| = O(|\ell'|)$ .

Daneshmand et al., Escaping saddles with stochastic gradients, 2018.

• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?



• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?



• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?



• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?

• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

Depends strongly on :

- What your problem is
- how you initialize

Yes, this bold statement will be my conclusion





• Does my algorithm converge?  $x_{\infty} \coloneqq \lim_{k \to +\infty} x_k$  exists?



• What is the nature of the limit  $x_{\infty}$ ? Global/Local minima? Saddle?

Depends strongly on :

- What your problem is
- how you initialize

Yes, this bold statement will be my conclusion





### Any questions ?

