# Complexity in continuous optimization (2/6)
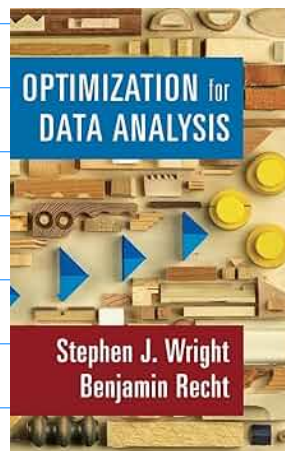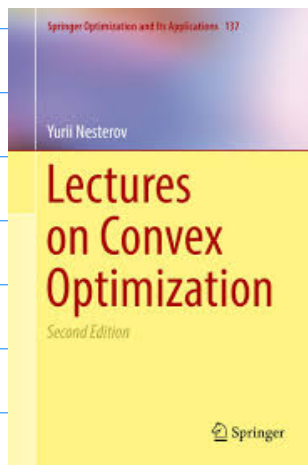
February 6, 2024

Today: Accelerated gradient from convex to nonconvex

References:



Yurii Nesterov — Lectures on Convex Optimization, Second Edition (Springer)

Stephen J. Wright, Benjamin Recht — OPTIMIZATION for DATA ANALYSIS

---

## "Convex Until Proven Guilty": Dimension-Free Acceleration of Gradient Descent on Non-Convex Functions

ICML 2017

Yair Carmon   John C. Duchi   Oliver Hinder   Aaron Sidford [1]

---

### Lower bounds for finding stationary points II: first-order methods

Yair Carmon[1] · John C. Duchi[2] · Oliver Hinder[3] · Aaron Sidford[3]

# Where we stand:

minimize $f(x)$ , $f \in C_L^{1,1}$ ( $\nabla f$ exists at every $x \in \mathbb{R}^n$
$x \in \mathbb{R}^n$
and
$$\|\nabla f(x) - \nabla f(y)\| \le L\|x-y\|)$$

## Goal:
Bound the number of iterations/gradient evaluations/function evaluations to find a point such that $\|\nabla f(x)\| \le \varepsilon$

## One algorithm: Gradient descent

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) \quad , \quad \alpha_k > 0$$

(e.g. $\alpha_k = \frac{1}{L}$)

Complexity: $\|\nabla f(x_k)\| \le \varepsilon$ after at most $O(\varepsilon^{-2})$ iterations

(valid $\forall \varepsilon > 0$, implies $\liminf_{k \to \infty} \|\nabla f(x_k)\| = 0$)

Upper bound is $O(\varepsilon^{-2})$

Lower bound matches the upper bound: There exists a function $f$ $C_L^{1,1}$ such that GD takes exactly $\varepsilon^{-2}$ iterations to reach a point such that $\|\nabla f(x)\| \le \varepsilon$.

$\Rightarrow$ Sharp analysis (Lower and upper bound match)

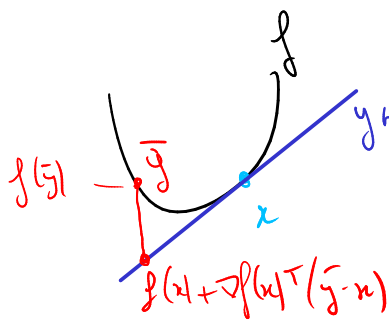$\Rightarrow$ How can we improve these results? Look at subclasses of $C_L^{1,1}$ functions

Today: • convex

• $C_L^{1,1} + C^2$ (nonconvex)

# ① Convex optimization

In this part, we suppose that $f$ is $C_L^{1,1}$ and convex or strongly convex.

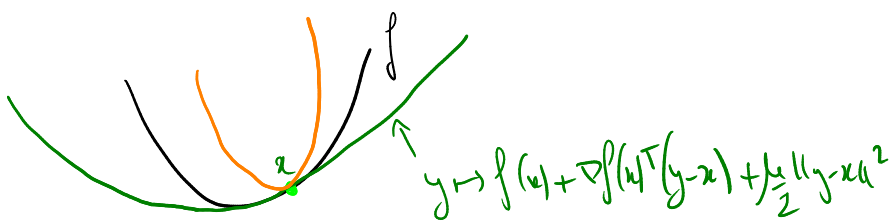$f \; C_L^{1,1}$ is convex $\iff$ $\forall (x,y) \in (\mathbb{R}^n)^2$,
$$f(y) \geq f(x) + \underbrace{\nabla f(x)^T (y-x)}_{\text{linear function of } y}$$



$y \mapsto f(x) + \nabla f(x)^T (y-x)$

$f(\bar{y}) - \bar{y}$

$x$

$f(x) + \nabla f(x)^T (\bar{y}-x)$

Property: If $f$ is $C^1$ convex, then $\nabla f(\bar{x}) = 0 \iff \bar{x} \in \arg\min_x f(x)$

($\bar{x}$ global minimum of $f$)

• $f \; C_L^{1,1}$ is $\mu$-strongly convex $\quad \mu > 0$
$$\iff \forall (x,y) \in (\mathbb{R}^n)^2, \quad f(y) \geq \underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2}_{\text{quadratic function of } y}$$



$y \mapsto f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$

$\left( \text{NB: } f \; C_L^{1,1} \implies \forall (x,y) \in (\mathbb{R}^n)^2, \quad f(y) \leq f(x) + \nabla f(x)^T (y-x) + \underbrace{\frac{L}{2} \|y-x\|^2}_{} \right)$

Property: If $f \; C_L^{1,1}$ and $\mu$-strongly convex, then it has a unique minimum which is the unique solution of $\nabla f(x) = 0_{\mathbb{R}^n}$

$f$ is $\mu$-strongly convex $\iff$ $x \mapsto \underbrace{f(x) - \frac{\mu}{2}\|x\|^2}_{\bar{f}}$ is convex

minimize $\quad \bar{f}(y) - \bar{f}(x) - \nabla\bar{f}(x)^T(y-x)$
$x,y$

$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = x_1^2$

$\quad\quad$ Gradient descent on convex problems

$\to x_{k+1} = x_k - \alpha_k \nabla f(x_k) = x_k - \frac{1}{L}\nabla f(x_k)$

$\to f$ $C_L^{1,1}$ and convex, and has at least 1 minimum

$\quad\quad$ Then, GD reaches a point $x_K$ such that

$\quad\quad$ <mark>$f(x_K) - \min\limits_{x \in \mathbb{R}^d} f(x) \le \varepsilon$</mark> $\quad$ in at most $O(\varepsilon^{-1})$ iterations

$\quad\quad\quad\quad ( \implies \|\nabla f(x_K)\| \le \varepsilon$ in at most $O(\varepsilon^{-1})$ iterations)

Typical criterion for complexity in the convex setting
(replaces $\|\nabla f(x)\| \le \varepsilon$ used in the nonconvex setting)
$\Rightarrow$ Indicates how far the current function value is from the minimum value

Proof: $\bullet$ Since $f$ is $C_L^{1,1}$, $\forall k \in \mathbb{N}$,

$f\left(\underbrace{x_k - \frac{1}{L}\nabla f(x_k)}_{x_{k+1}}\right) \le f(x_k) + \nabla f(x_k)^T\left(x_k - \frac{1}{L}\nabla f(x_k) - x_k\right) + \frac{L}{2}\|x_k - \frac{1}{L}\nabla f(x_k) - x_k\|^2$

$\le f(x_k) - \frac{1}{2L}\|\nabla f(x_k)\|^2$

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2$$

Let $x^* \in \operatorname*{argmin}_x f(x)$ and $f^* = f(x^\circ) = \min_x f(x)$

By convexity, $\quad f(x^*) \geq f(x_k) + \nabla f(x_k)^T (x^* - x_k)$

$$\Longleftrightarrow \quad f(x_k) \leq f(x^*) + \nabla f(x_k)^T (x_k - x^*)$$

Hence, $\quad f(x_{k+1}) \leq f(x^*) + \nabla f(x_k)^T (x_k - x^*) - \frac{1}{2L} \|\nabla f(x_k)\|^2$

Then $\quad f(x_{k+1}) - f(x^*) \leq \nabla f(x_k)^T (x_k - x^*) - \frac{1}{2L} \|\nabla f(x_k)\|^2$

$$= \frac{L}{2} \left( \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2 \right)$$

$$\underset{\uparrow}{\phantom{x}}$$

$$x_k - \frac{1}{L} \nabla f(x_k)$$

Suppose that $\quad f(x_k) - f(x^*) > \varepsilon \quad \forall k = 0, \ldots, K$

Then $\quad \sum_{k=0}^{K-1} f(x_{k+1}) - f(x^\circ) \leq \frac{L}{2} \sum_{k=0}^{K-1} \left( \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2 \right)$

$$= \frac{L}{2} \left( \|x_0 - x^*\|^2 - \|x_K - x^*\|^2 \right)$$

$$K\varepsilon < \sum_{k=0}^{K-1} \underbrace{\left( f(x_{k+1}) - f(x^\circ) \right)}_{> \varepsilon} \leq \boxed{\frac{L}{2} \|x_0 - x^*\|^2}$$

$$\Rightarrow \quad K < \frac{L}{2} \| x_0 - x^\bullet \|^2 \, \varepsilon^{-1} = O(\varepsilon^{-1})$$

For $\mu$-strongly convex case, can show

$$f(x_k) - f(x^\bullet) \le \left(1 - \frac{\mu}{L}\right)^k \left(f(x_0) - f(x^\bullet)\right)$$

( Behind this: $\underline{\| \nabla f(x) \|^2 \ge 2\mu \left(f(x) - f(x^\bullet)\right)}$ , nontrivial

consequence of strong convexity inequality )

$$f C_L^{1,1} \Rightarrow \quad f(x_{k+1}) \le f(x_k) - \frac{1}{2L} \| \nabla f(x_k) \|^2$$

$$\le f(x_k) - \frac{\mu}{L} \left(f(x_k) - f(x^\bullet)\right)$$

$$f(x_{k+1}) - f(x^\bullet) \le f(x_k) - f(x^\bullet) - \frac{\mu}{L} \left(f(x_k) - f(x^\bullet)\right) = \left(1 - \frac{\mu}{L}\right)\left(f(x_k) - f(x^\bullet)\right)$$

$\rightarrow$ The proof in the strongly convex setting is (in some way)
easier than the proof in the convex setting

$$\text{Complexity of GD}: \quad O\left(\frac{L}{\mu} \ln\left(\varepsilon^{-1}\right)\right)$$

$\uparrow$
constant in terms of $\varepsilon$ but
important for the complexity
$L/\mu$ : "condition number"

Same GD algorithm $f C_L^{1,1}$

$f$ nonconvex    $\| \nabla f(x) \| \le \varepsilon$    $O(\varepsilon^{-2})$

$f$ convex    $f(x) - \min_x f(x) \le \varepsilon$    $O(\varepsilon^{-1})$

$f$ $\mu$-strongly convex    $\sim$    $O\left(\frac{L}{\mu} \ln(\varepsilon^{-1})\right)$

<u>Pb</u>: In the convex/strongly convex settings, the upper bonds for GD do not match the lower bonds

Nemirovski & Yudin 1983

Existence result:
here exists an algorithm that uses only one gradient per iteration

→ If $f$ $C_L^{1,1}$ convex, there exits an algorithm with complexity $O(\varepsilon^{-1/2})$

→ If $f$ $C_L^{1,1}$ $\mu$-strongly convex, $\exists$ algorithm with complexity $O\left(\sqrt{\frac{L}{\mu}} \ln(\varepsilon^{-1})\right)$

$\uparrow$ $\mu \leq L$  $\sqrt{\frac{L}{\mu}} \leq \frac{L}{\mu}$

→ Algorithm was unknown until Yurii Nesterov discovered it in 1983

<span style="color:red">Accelerated gradient/Nesterov's method $(x_0 \in \mathbb{R}^n)$</span>

$$x_{k+1} = x_k - \alpha_k \nabla f\left(x_k + \underbrace{\beta_k(x_k - x_{k-1})}_{\text{Momentum term}}\right) + \beta_k(x_k - x_{k-1})$$

$\alpha_k > 0$ , $\beta_k > 0$   $(\beta_0 = 0, x_{-1} = x_0)$

Equivalent formulation:
$$\begin{cases} y_k = x_k + \beta_k(x_k - x_{k-1}) \\ x_{k+1} = y_k - \nabla f(y_k) \end{cases}$$

$\beta_0 = 0$
$x_{-1} = x_0$
$y_0 = x_0$

<u>Note</u>: Momentum methods in ML (SGD with momentum, Adam)
$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) + \beta_k(x_k - x_{k-1})$$

# Analyzing Nesterov's method

$$\begin{cases} y_k = x_k + \beta_k (x_k - x_{k-1}) \\ x_{k+1} = y_k - \alpha_k \nabla f(y_k) \end{cases} \qquad y_0 = x_0, \; \beta_0 = 0, \; x_{-1} = x_0$$

a) $f$ is $\mu$-strongly convex

$$\rightarrow \alpha_k = \frac{1}{L}, \; \beta_k \underset{k \geq 1}{=} \frac{\sqrt{K}-1}{\sqrt{K}+1} \quad \text{where} \quad K = \frac{L}{\mu}$$

## Key difference with GD analysis                    $f^* = \min f(x)$

- In GD, we look at $f(x_k) - f^*$ and we show

$$f(x_k) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k (f(x_0) - f^*)$$

$$= \left(1 - \frac{1}{K}\right)^k (f(x_0) - f^*)$$

"Lyapunov function for GD"

- In AG (Accelerated Gradient)'s analysis, we use

$$V_k := f(x_k) - f^* + \frac{L}{2} \| x_k - x^* - \rho^2 (x_{k-1} - x^*) \|^2$$

$$\text{with} \quad \rho^2 = \left(1 - \frac{1}{\sqrt{K}}\right) = 1 - \sqrt{\frac{\mu}{L}}$$

$$\Rightarrow \text{the analysis shows } V_k \leq \rho^k V_0$$

$$V_0 = f(x_0) \cdot f^* + \frac{L}{2} \| (1 - \rho^2)(x_0 - x^*) \|^2$$
$$= f(x_0) - f^* + \frac{\mu}{2} \| x_0 - x^* \|^2$$

$\Rightarrow$ Given the $O\left( \sqrt{\frac{L}{\mu}} \ln(\varepsilon^{-1}) \right)$ bound

vs $O\left( \frac{L}{\mu} \ln(\varepsilon^{-1}) \right)$ for gradient descent

b) <u>Convex</u> case ( <u>Not</u> strongly convex)

Nesterov's implementation of accelerated gradient

$$\begin{cases} y_k = x_k + \beta_k (x_k - x_{k-1}) \\ x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k) \end{cases}$$

$\beta_0 = 0$

$\beta_k = \rho_k \, \rho_{k-1} \quad k \geq 1$

$\rho_0 = \rho_{-1} = 0$

<span style="color:red">Most non-intuitive part of the algorithm: $\{\rho_k\}$ and $\{\beta_k\}$ is defined independently of the problem ($f$) and $x_0$</span>

$\rho_{k+1}$ is the positive root

of $1 + \rho_{k+1} (\rho_k^2 - 1) - \rho_{k+1}^2 = 0$

$$\rho_{k+1} = \frac{1 + \sqrt{1 + 4\rho_k^2}}{2}$$

$\longrightarrow$ Analysis relies on the Lyapunov function

$$W_k = f(x_k) - f^* + \frac{L}{2} \| (x_k - x^*) - \rho_{k-1}^2 (x_{k-1} - x^*) \|^2$$

$\longrightarrow$ slow $\quad W_k \leq \rho_{k-1}^2 \cdots \rho_1^2 \, W_1 = (1 - \rho_{k-1}^2)^2 W_1$

and $W_1 \leq \boxed{\frac{L}{2} \| x_0 - x^* \|^2}$

$$\longrightarrow \quad 1 - \rho_k^2 \leq \frac{2}{k+2} \qquad \forall k \quad \text{(by definition of } \rho_k \text{)}$$

Overall, prove

$$f(x_k) - f^\circ \leq V_k \leq \frac{2L}{(k+1)^2} \|x_0 - x^\circ\|^2$$

$$\Rightarrow \quad f(x_k) - f^\circ \leq \varepsilon \quad \text{after at most} \quad O(\varepsilon^{-1/2})$$

iterations.

$\longrightarrow$ Nesterov's method attains the upper bound $O(\varepsilon^{-1/2})$

$\longrightarrow$ Nesterov also showed a lower bound for the convex setting

$$f(x) = \frac{1}{2} x^T A x - e_1^T x \qquad \text{convex quadratic}$$

$$= \frac{1}{2}(x_1 - 1)^2 + \frac{1}{2}\sum_{i=1}^{n-1}(x_i - x_{i+1})^2$$

$$A = \begin{bmatrix} 2 & -1 & & 0 \\ -1 & & & \\ & & & -1 \\ 0 & & -1 & 2 \end{bmatrix}, \quad e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \text{unique solution} \quad x^\circ = \begin{bmatrix} 1 - \tfrac{1}{m+1} \\ \vdots \\ 1 - \tfrac{m}{m+1} \end{bmatrix}$$

Any algorithm starting from $x_0 = 0$ and using only 1 gradient per iteration will satisfy

$$f(x_k) - f^\circ \geq \frac{3}{8(k+1)^2} \|x_0 - x^*\|^2 \qquad \forall k \leq \frac{m}{2} - 1$$

Nesterov's example: dimension "$\varepsilon^{-1/2}$" but quadratic convex

Cartis, Gould Toint examples: dimension 1 or 2, but very nonlinear nonquadratic

## ② Acceleration in the noncovex setting

$\longrightarrow$ If $f$ is just $C_L^{1,1}$ and noncovex, cannot do better than GD!

$\longrightarrow$ If $f$ is $C_L^{1,1}$, noncovex and $C^{2,2}$, then you can do better than GD, and you can do that with a variant of Nesterov's method  ( Carmon et al '17 )

Idea:  "Convex until proven guilty"

- Run AG as if the function were strongly convex
- Two cases:  — Either the method works as in the strongly convex setting
 — Or it doesn't, and take a negative curvature step instead of doing 1 iteration of AG

## Negative curvature

$f$ noncovex $C^2$ :     $\bar{x} \in \arg\min_x f(x) \Rightarrow \begin{cases} \nabla f(\bar{x}) = 0 \\ \nabla^2 f(\bar{x}) \succeq 0 \end{cases}$

$\forall v \in \mathbb{R}^n, \ v^\top \nabla^2 f(\bar{x}) v \geq 0$

When $f$ is convex $C^2$, $\nabla^2 f(x) \succeq 0$ true $\forall x$!

For nonconvex $f$, if $\nabla^2 f(x) \succeq 0$ does not hold, the function decreases from $x$ along any direction $v$ such that $v^T \nabla^2 f(x) v < 0$

$\underset{\uparrow}{}$
Negative curvature direction

(and $v^T \nabla f(x) \leq 0$)

Negative curvature step: $x + \alpha v$
$\alpha > 0$

Algorithm $\left( x_0, L, \varepsilon, \varepsilon_H \right)$      $L$ Lipschitz constant for the gradient

$\underset{\uparrow}{}$
$\varepsilon > 0$ for complexity: used in the algorithm

For $k = 0, 1, \text{---}$

$\rightarrow$ Define $\hat{f}_k(z) = f(z) + \varepsilon_H \|z - x_k\|^2 \simeq$ Lyapunov function

$\rightarrow$ Run AG on $\hat{f}_k$ to find a point $z_k$ such that

$\|\nabla \hat{f}_k(z_k)\| \leq \varepsilon/10$

assuming that the function $\hat{f}_k$ is $\varepsilon_H$-strongly convex

true if $f$ is convex, false otherwise

$f(u) \geq \hat{f}(v) + \nabla \hat{f}(v)^T(u-v)$
$+ \frac{\varepsilon_H}{2}\|u-v\|^2$:

· At every iteration of AG, check whether the strong convexity inequality holds between every pair of points that is computed, and stop if the inequality is violated.

- If AG stops with $\|\nabla \hat{f}_k(z_k)\| \leq \frac{\varepsilon}{10}$, define $x_{k+1} = z_k$.

- Otherwise, get a pair $(u_k, v_k)$ that violates the strong convexity inequality. and use $u_k - v_k$ as a negative curvature direction.

$$x_{k+1} = x_k + \eta\,(u_k \cdot v_k)$$

The function is nonconvex "enough" around $x_k$ for AG to fail

---

Analysis: At every iteration $k$,

$$f(x_k) - f(x_{k+1}) \geq O\left( \frac{\|\nabla f(x_k)\|^2}{\varepsilon_H} \right) \geq O\left( \frac{\varepsilon^2}{\varepsilon_H} \right)$$

if $\|\nabla f(x_k)\| \geq \varepsilon$

if AG succeeds

- $f(x_k) - f(x_{k+1}) \geq O\left( \varepsilon_H^3 \right)$

if AG fails

$\Rightarrow$ The method computes $x_k$ such that $\|\nabla f(x_k)\| \leq \varepsilon$ in at most $O\left( \max\left( \varepsilon^{-2}\,\varepsilon_H,\ \varepsilon_H^{-3} \right) \right)$ iterations

$\Rightarrow$ With $\varepsilon_H = \varepsilon^{1/2}$, get $O\left( \varepsilon^{-3/2} \right)$ outer iterations

Better than $O(\varepsilon^{-2})$ for GD !

$\rightarrow$ Every call to AG terminates (with success/failure)
after at most $O\left(\sqrt{\frac{L+2\varepsilon_H}{\varepsilon_H}}\, \ln\left(\frac{1}{\varepsilon}\right)\right) = O\left(\sqrt{\frac{1}{\varepsilon_H}}\, \ln\left(\frac{1}{\varepsilon}\right)\right)$

$\Rightarrow$ Total number of iterations:

$$O\left(\frac{1}{\sqrt{\varepsilon_H}}\, \ln\left(\frac{1}{\varepsilon}\right)\right) \times O\left(\max\left(\varepsilon^{-2}\varepsilon_H,\ \varepsilon_H^{-3}\right)\right)$$

$\varepsilon_H = \varepsilon^{1/2}$ $\quad\curvearrowright$

$$O\left(\varepsilon^{-7/4}\, \ln\left(1/\varepsilon\right)\right) = \widetilde{O}\left(\varepsilon^{-7/4}\right)$$

$\Rightarrow$ Still better (for small $\varepsilon$) than $O(\varepsilon^{-2})$

$\varepsilon^{-7/4}$ : Best known upper bound for algorithms that use only gradient information to optimize a $C^{1,1} \cap C^{2,2}$ function

Lower bound (for this class of methods/functions): $O\left(\varepsilon^{-12/7}\right)$

$\rightarrow$ Attained for an example with dimension $O(\varepsilon^{-12/7})$

Takeaways
- AG better than GD for convex $\left.\phantom{\begin{array}{c}\\ \text{strongly convex}\end{array}}\right\}$versions
  strongly convex
- Analysis: Lyapunov functions + convexity inequalities
- Extension to nonconvex (+ negative curvature)