

# Advanced gradient descent and convex optimization

I Convexity

II Gradient descent (GD) on convex problems

III Acceleration

## I] Setup

Pb: minimize  $f(w)$   $f: \mathbb{R}^d \rightarrow \mathbb{R}$   
problem variables  $w \in \mathbb{R}^d$  objective function

$\operatorname{argmin}_{w \in \mathbb{R}^d} f(w) \subseteq \mathbb{R}^d$ : set of solutions

$\min_{w \in \mathbb{R}^d} f(w)$ : minimum value of  $f$  / of the problem

$\mathbb{R} \cup \{-\infty, +\infty\}$   
 $f$  unbounded  $f$  not defined on  $\mathbb{R}^d$

$\hookrightarrow$  We will assume that  $f$  is convex and continuously differentiable.

Continuously differentiable ( $C^1$ ):  $\nabla f$  exists at every point and  $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is continuous

If  $w^* \in \mathbb{R}^d$  is a local minimum, then  $\|\nabla f(w^*)\| = 0$

Convexity:  $f \in C^1$  is convex iff

$$\forall (u, v) \in (\mathbb{R}^d)^2, f(u) \geq f(v) + \nabla f(v)^T (u - v)$$

Property: If  $f$  is  $C^1$  and convex,

$$[w^* \in \mathbb{R}^d \text{ global minimum of } f] \Leftrightarrow \|\nabla f(w^*)\| = 0$$

Def:  $f \in C^1$  is  $\mu$ -strongly convex (for  $\mu > 0$ )  
 if  $\forall (u, v) \in (\mathbb{R}^d)^2$ ,  
 $f(u) \geq f(v) + \nabla f(v)^T (u-v) + \frac{\mu}{2} \|u-v\|^2$

Property: If  $f \in C^1$  and  $\mu$ -strongly convex, it has a unique global minimum.

Question: Given  $f \in C^1$  and (at least) convex, can we design a procedure to compute a global minimum of  $f$ ?

Goals: ① Design an algorithm that provably converges towards a solution

For convex problems, this may mean:

$\Leftrightarrow$  conv towards global min thanks to convexity

a) Generate  $\{w_k\}_k$  such that  $\|\nabla f(w_k)\| \xrightarrow[k \rightarrow \infty]{} 0$

b)  $\{w_k\}_k$  such that  $f(w_k) \xrightarrow[k \rightarrow +\infty]{} \min_{w \in \mathbb{R}^d} f(w)$

c)  $\{w_k\}_k$  such that  $w_k \xrightarrow[k \rightarrow +\infty]{} w^*$  where  $w^* \in \arg \min_{w \in \mathbb{R}^d} f(w)$

② Efficiently compute an approximate solution to the problem

$\Rightarrow$  Popular metrics: complexity / convergence rate

\* (Work-case) complexity: Given a tolerance / level of accuracy, how much does it cost to reach this level of accuracy

Ex)  $\forall \varepsilon > 0$ ,  $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$

in at most  $O(\varepsilon^{-1})$  iterations

$O(A) = \text{constant} \times A$   
 does not depend on  $A$

\* Convergence rate: Given a budget (of iterations, e.g.), what is the level of accuracy that we can expect?  
 Ex)  $\forall k \geq 1, f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{k}\right)$

## II Gradient descent for convex and strongly convex problems

Idea:  $w^*$  global minimum  $\Leftrightarrow \nabla f(w^*) = 0$

If  $\nabla f(w) \neq 0$ ,  $w$  not global (nor local) minimum, and there must exist a better point nearby  
 $\Rightarrow$  To find such a point, we look into the direction of steepest descent ( $-\nabla f(w)$ )

GD iteration:  $k=0$  Pick  $w_0 \in \mathbb{R}^d$   
 $k \geq 0$   $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$   
 $\alpha_k > 0$  stepsize\*

\* learning rate

How to choose  $\alpha_k$ ?

- $\hookrightarrow$  Predefined sequence (constant, decreasing)
- $\hookrightarrow$  Adaptive way (e.g. via line search)
- $\ominus$  Induces additional cost in general

$\Rightarrow$  Choice of  $\alpha_k$  is generally guided by own assumptions on  $f$

Here, we will assume that  $f$  is continuously differentiable with  $L$ -Lipschitz continuous gradient  
 (Other names:  $L$ -smooth  
 $L$ -Lipschitz continuously differentiable  
 $C_{L,1}^{1,1}$ )

$\nabla f$   $L$ -Lipschitz continuous  $\|\nabla f(u) - \nabla f(v)\| \leq L\|u - v\|$

Prop: If  $f \in C_{L,1}^{1,1}$ , then  $\forall (u, v) \in (\mathbb{R}^d)^2$

$$(*) \quad f(u) \leq f(v) + \nabla f(v)^T (u - v) + \frac{L}{2} \|u - v\|^2$$

( $\neq$  convexity inequality  
 $f(u) \geq f(v) + \nabla f(v)^T (u - v) + \frac{\mu}{2} \|u - v\|^2$ )

Proposition (Descent property).

Suppose that  $f \in C_{L,1}^{1,1}$ .

Run gradient descent with  $\alpha_k = \frac{1}{L}$ .

Then,  $\forall k \in \mathbb{N}$ ,

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

Proof: Use the inequality (\*) with

$$u = w_{k+1} = w_k - \alpha_k \nabla f(w_k) = w_k - \frac{1}{L} \nabla f(w_k)$$

$$v = w_k$$

$$f(w_{k+1}) \leq f(w_k) + \nabla f(w_k)^T \left(-\frac{1}{L} \nabla f(w_k)\right) + \frac{L}{2} \left\|-\frac{1}{L} \nabla f(w_k)\right\|^2$$

$$= f(w_k) - \frac{1}{L} \nabla f(w_k)^T \nabla f(w_k) + \frac{L}{2} \times \frac{1}{L^2} \|\nabla f(w_k)\|^2$$

$$= f(w_k) + \left(-\frac{1}{L} + \frac{1}{2L}\right) \|\nabla f(w_k)\|^2$$

$$= f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2$$

NB:  $w_k - \frac{1}{L} \nabla f(w_k) \in \underset{w \in \mathbb{R}^d}{\text{argmin}} \left\{ \underbrace{f(w_k) + \nabla f(w_k)^T (w - w_k) + \frac{L}{2} \|w - w_k\|^2}_{(*) \text{ with } v = w_k} \right\}$

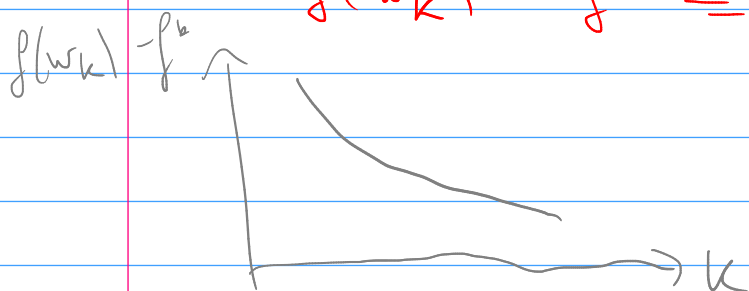
## Th (Convergence rate of GD on convex functions)

Suppose that  $f$  is  $C^{1,1}$  and convex.

Suppose that  $\text{argmin}_{w \in \mathbb{R}^d} f(w) \neq \emptyset$  and  
 let  $w^* \in \text{argmin}_{w \in \mathbb{R}^d} f(w)$  and  $f^* = f(w^*) = \min_{w \in \mathbb{R}^d} f(w)$

Run GD with  $\alpha_k = \frac{1}{L}$  starting with  $w_0 \in \mathbb{R}^d$ .  
 Then,  $\forall k \geq 1$ ,

$$f(w_k) - f^* \leq \frac{L}{2} \|w_0 - w^*\|^2 \times \frac{1}{k} = O\left(\frac{1}{k}\right)$$



If you increase  $k$  to  $10k$ ,  
 you are 10 times  
 more accurate

## Th (Convergence rate of GD for strongly convex functions)

Let  $f \in C^{1,1}$  and  $\mu$ -strongly convex.

Let  $w^* \in \text{argmin}_{w \in \mathbb{R}^d} f(w)$  and  $f^* = f(w^*)$

Run GD with  $\alpha_k = \frac{1}{L}$  starting from  $w_0 \in \mathbb{R}^d$ .  
 Then,  $\forall k \in \mathbb{N}$ ,

$$f(w_k) - f^* \leq \left(1 - \sqrt{\frac{\mu}{L}}\right)^k (f(w_0) - f^*)$$

$$\in (0, 1)$$

$$\left(1 - \frac{\mu}{L}\right)^k \xrightarrow{k \rightarrow +\infty} 0$$

$\hookrightarrow \left(1 - \frac{\mu}{L}\right)^k$  goes  
 to 0 faster than  $\frac{1}{k}$

"Solving strongly convex problems is easier than  
 convex problems."

## Proof of the convergence rate of GD in the convex case

Setup:  $f \in C_L^{1,1}$  and convex

$$\forall (u, v) \in \mathbb{R}^d,$$

$$(*) \quad f(u) \leq f(v) + \nabla f(v)^T (u-v) + \frac{L}{2} \|u-v\|^2$$

$$(**) \quad f(u) \geq f(v) + \nabla f(v)^T (u-v)$$

GD run with  $\alpha_k = 1/L$  starting from  $w_0$

$\hookrightarrow$  Because of  $(*)$  with  $\alpha_k = 1/L$ , the descent property holds, i.e.

$$(1) \quad \forall k \geq 0, \quad f(w_{k+1}) \leq f(w_k) - \frac{1}{2L} \|\nabla f(w_k)\|^2 \leq f(w_k)$$

$\hookrightarrow$  Because of convexity,  $(**)$  applied with  $v = w_k$  and  $u = w^*$  ( $\in \arg\min_w f(w)$ )

$$f(w^*) \geq f(w_k) + \nabla f(w_k)^T (w^* - w_k)$$

$\Leftrightarrow$

$$(2) \quad f(w_k) \leq f(w^*) + \nabla f(w_k)^T (w_k - w^*)$$

For every  $k \in \mathbb{N}$ ,  $(1) + (2)$  gives

$$f(w_{k+1}) \leq f(w^*) + \nabla f(w_k)^T (w_k - w^*) - \frac{1}{2L} \|\nabla f(w_k)\|^2$$

"Trick": Notice that

$$\nabla f(w_k)^T (w_k - w^*) - \frac{1}{2L} \|\nabla f(w_k)\|^2 = \frac{L}{2} (\|w_k - w^*\|^2 - \|w_{k+1} - w^*\|^2)$$

Hence,  $\forall k \geq 0$

$$(*) \quad f(w_{k+1}) \leq f(w^*) + \frac{L}{2} (\|w_k - w^*\|^2 - \|w_{k+1} - w^*\|^2)$$

Our goal is to bound  $f(w_k) - f(w^*)$  for some  $k \geq 1$ .

We sum all inequalities of the form  $(*)$  over  $k \in \{0, \dots, k-1\}$

$$\sum_{k=0}^{k-1} f(w_{k+1}) \leq k f(w^*) + \frac{L}{2} \sum_{k=0}^{k-1} (\|w_k - w^*\|^2 - \|w_{k+1} - w^*\|^2)$$

$\sum_{k=0}^{k-1} (a_k - a_{k+1}) = a_0 - a_k$   
Telescopic sum

$$\sum_{k=0}^{k-1} f(w_{k+1}) \leq k f(w^*) + \frac{L}{2} (\|w_0 - w^*\|^2 - \underbrace{\|w_k - w^*\|^2}_{\leq 0})$$

$$\sum_{k=0}^{k-1} f(w_{k+1}) \leq k f(w^*) + \frac{L}{2} \|w_0 - w^*\|^2$$

$$\sum_{k=0}^{k-1} (f(w_{k+1}) - f(w^*)) \leq \frac{L}{2} \|w_0 - w^*\|^2$$

Finally, per the descent property, we have

$$f(w_0) \geq f(w_1) \geq \dots \geq f(w_{k-1}) \geq f(w_k)$$

$$\text{Therefore, } \sum_{k=0}^{k-1} (f(w_{k+1}) - f(w^*)) \geq k (f(w_k) - f(w^*))$$

and we obtain

$$k [f(w_k) - f(w^*)] \leq \frac{L}{2} \|w_0 - w^*\|^2$$

$$f(w_k) - \underbrace{f(w^*)}_{f^*} \leq \frac{L}{2} \|w_0 - w^*\|^2 \times \frac{1}{k}$$

Q.E.D.



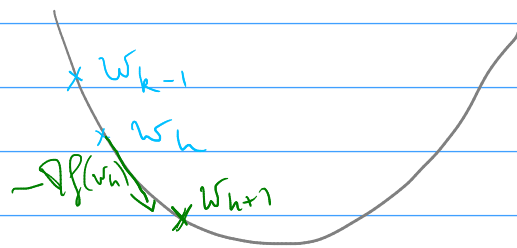
### III Acceleration

Question: Given a class of optimization problems ( $C^{1,1}$  and convex/strongly convex) what is the fastest algorithm in terms of convergence rate that relies only on 1 gradient evaluation per iteration?  
comparable to GD

Answer: Accelerated gradient techniques  
Gradient methods with momentum

Idea: At every iteration, we continue moving in the direction of the previous step, and we combine this move with one in the (negative) gradient direction.

- \* Locally,  $-\nabla P(\cdot)$  is a good direction
- \* The step you just took (in a negative gradient direction) may still be good for your current iterate



### III - 1) Heavy-ball method (Polyak, 1964)

Iteration  $w_0 \in \mathbb{R}^d$ ,  $w_{-1} = w_0$

$$(HB) \quad \forall k \geq 0, \quad w_{k+1} = \underbrace{w_k - \alpha_k \nabla P(w_k)}_{\text{GD step}} + \underbrace{\beta_k (w_k - w_{k-1})}_{\text{Momentum step}}$$

Guarantees: If  $f$  is a  $\mu$ -strongly convex quadratic and  $C_{L,1}^{1,1}$ ,  
 then (HB) with

$$\alpha_k = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2} \quad \text{and} \quad \beta_k = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$$

then  $\forall k \geq 0$ ,

$$f(w_k) - f^* \leq O\left(\left(1 - \sqrt{\frac{\mu}{L}}\right)^k\right)$$

Optimal CV rate for strongly convex quadratics

(Recall: For GD,  $f(w_k) - f^* \leq O\left(\left(1 - \frac{\mu}{L}\right)^k\right)$ )

$$\mu \leq L \quad 1 - \frac{\mu}{L} \geq 1 - \sqrt{\frac{\mu}{L}}$$

so  $\left(1 - \sqrt{\frac{\mu}{L}}\right)^k \rightarrow 0$  faster than  $\left(1 - \frac{\mu}{L}\right)^k$

⚠ The analysis does not extend to general strongly convex functions

⚠ In heavy ball, you can have  $f(w_{k+1}) > f(w_k)$

III 2.) Nesterov's method (strongly convex)

1983

Yurii Nesterov

Iteration  $w_0 \in \mathbb{R}^d, w_{-1} = w_0$

$$\forall k \geq 0, w_{k+1} = w_k - \alpha_k \nabla f(w_k + \beta_k(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1})$$

Idea: Do the momentum step before the gradient step  
 (Heavy ball: gradient step before momentum)

Convergence rate:  $f \in C^{1,2}$ ,  $\mu$ -strongly convex (not necessarily quadratic)

Run Nesterov's method (aka Nesterov's accelerated gradient) with  $\alpha_k = \frac{1}{L}$  (like GD) and  $\beta_k = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$  (like HB)

Then,  $\forall k \geq 0$ ,

$$f(w_k) - f^* \leq \frac{L + \mu}{2} \|w_0 - w^*\|^2 \left(1 - \sqrt{\frac{\mu}{L}}\right)^k \\ = O\left(\left(1 - \sqrt{\frac{\mu}{L}}\right)^k\right)$$

optimal for strongly convex functions

Proof: Technical but the idea is to work with a sequence of quadratic functions that approximate  $f$  more closely as  $h \rightarrow +\infty$

III - 3 } Nesterov's method for convex functions

Same as before but  $\{\beta_k\}_k$  is predefined independently of  $f$ !

$$t_0 = 0, \quad t_{k+1} = \frac{1}{2} (1 + \sqrt{1 + 4t_k^2}), \quad \beta_k = \frac{t_k - 1}{t_{k+1}}$$

CV rate: If  $f \in C^{1,2}$  convex, run Nesterov's method with the above choice for  $\beta_k$  and  $\alpha_k = \frac{1}{L}$ ,

then  $\forall k \geq 0$ ,

$$f(w_k) - f^* \leq \frac{2L \|w_0 - w^*\|^2}{(k+1)^2} = O\left(\frac{1}{k^2}\right)$$

(GD:  $O(1/k)$ )

optimal for convex functions