

# Optimization for Machine Learning

September 30, 2024

Today: Lecture on advanced gradient descent

Tomorrow: Tutorial / lab (bring laptops if you can)

Aftwards: A long break until December

⇒ Part 1 of homework assignment

⇒ Updated lecture notes with exercises

# ADVANCED ASPECTS OF GRADIENT DESCENT

Setup: minimize  $f(w)$   $f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $w \in \mathbb{R}^d$

Assumptions: .  $f$  is bounded below

$$(\exists \bar{f} \in \mathbb{R}, f(w) \geq \bar{f} \text{ for any } w \in \mathbb{R}^d)$$

  $\bar{f}$  is not necessarily equal to  $\min_{w \in \mathbb{R}^d} f(w)$

.  $f$  is  $C_L^{1,1}$  (the gradient of  $f$  exists  $\forall w \in \mathbb{R}^d$  and it is L-Lipschitz continuous)

$$\nabla f(w, v) \in (\mathbb{R}^d)^2, \quad \|\nabla f(v) - \nabla f(w)\| \leq L \|v - w\|$$

("If  $v$  is close to  $w$ , then  $\nabla f(v)$  is close to  $\nabla f(w)$ ")

Ex) Linear regression, logistic regression

Remark: In general, functions are only  $C_L^{1,1}$  on subsets of  $\mathbb{R}^d$   
 $\Rightarrow$  The analysis of gradient descent for  $C_L^{1,1}$  functions extends to more general classes of functions

Key inequality: If  $f$  is  $C^1_L$ , then for any  $(w, v) \in (\mathbb{R}^d)^2$ ,

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

function we want to minimize (as a function of  $v$ )

For fixed  $w$ , this a function of  $v$  that is quadratic, convex, easy to minimize

Given  $w \in \mathbb{R}^d$ , if we find  $v$  such that

$$f(w) + \nabla f(w)^T(v-w) + \frac{L}{2} \|v-w\|^2 < f(w), \text{ then}$$

the inequality guarantees  $f(v) < f(w)$  without having to compute  $f(v)$

The function  $v \mapsto f(w) + \nabla f(w)^T(v-w) + \frac{L}{2} \|v-w\|^2$  is minimized at  $v^* = w - \frac{1}{L} \nabla f(w)$

$\Rightarrow$  This is precisely gradient descent with stepsize  $\frac{1}{L}$ !

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

$$\begin{aligned} u^T(\alpha v) &= \alpha(u^Tv) \\ \|\alpha v\|^2 &= \alpha^2 \|v\|^2 \end{aligned}$$

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

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

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

$$= f(w) - \frac{1}{L} \|\nabla f(w)\|^2 + \frac{1}{2L} \|\nabla f(w)\|^2 = f(w) + \underbrace{\left[ -\frac{1}{L} + \frac{1}{2L} \right] \|\nabla f(w)\|^2}_{-\frac{1}{2L}}$$

$$\left( \nabla f(w)^T \nabla f(w) = \|\nabla f(w)\|^2 \right)$$

$$= f(w) - \frac{1}{2L} \|\nabla f(w)\|^2 \leq f(w) \quad \text{if } \nabla f(w) \neq 0_{R^d}$$


---

minimize  $\underbrace{f(w) + \nabla f(w)^T(v-w) + \frac{L}{2} \|v-w\|^2}_{g(v)}$  for fixed  $w \in \mathbb{R}^d$

$g$  is convex (sum of convex functions: linear + norm<sup>2</sup>)

$$\begin{aligned} v^* \in \underset{v \in \mathbb{R}^d}{\operatorname{argmin}} g(v) &\quad (\Rightarrow \nabla g(v^*) = 0) \\ &\quad \Leftrightarrow \nabla f(w) + L(v^*-w) = 0 \\ &\quad \Leftrightarrow Lv^* = Lw - \nabla f(w) \\ &\quad \Leftrightarrow v^* = w - \frac{1}{L} \nabla f(w) \end{aligned}$$

Gradient formula (generic rules)

- The gradient of  $w \mapsto x^T w + y$  is  $x$
  - The gradient of  $w \mapsto \alpha \|w-u\|^2$  is  $2\alpha(w-u)$
- 

For  $C_L^{1,1}$  functions, applying gradient descent (with  $\frac{1}{L}$ ) corresponds to solving a sequence of convex, quadratic

## optimization problems

$$\text{GD: } w_0 \in \mathbb{R}^d, \quad w_{k+1} = w_k - \frac{1}{L} \nabla f(w_k) \quad k=0, 1, \dots$$

$$w_0 \rightarrow \text{Pick } w_1 \in \underset{v \in \mathbb{R}^d}{\arg \min} \quad f(w_0) + \nabla f(w_0)^T(v - w_0) + \frac{L}{2} \|v - w_0\|^2$$

$$\Rightarrow w_1 = w_0 - \frac{1}{L} \nabla f(w_0)$$

$$w_1 \rightarrow \text{Pick } w_2 \in \underset{v \in \mathbb{R}^d}{\arg \min} \quad f(w_1) + \nabla f(w_1)^T(v - w_1) + \frac{L}{2} \|v - w_1\|^2$$

$$\Rightarrow w_2 = w_1 - \frac{1}{L} \nabla f(w_1)$$

$$w_2 \rightarrow \dots$$

Q: What can we guarantee after K iterations of gradient descent?

For any  $k=0, \dots, K-1$ , we have shown that

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

Therefore

$$f(w_K) \leq f(w_{K-1}) - \frac{1}{2L} \|\nabla f(w_{K-1})\|^2$$

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

$$f(w_K) \leq f(w_0) - \frac{1}{2L} \sum_{k=0}^{K-1} \|\nabla f(w_k)\|^2$$

Since  $f$  is bounded below,  $f(w_K) \geq \bar{f}$

Hence  $\bar{f} \leq f(w_0) - \frac{1}{2L} \sum_{k=0}^{K-1} \|\nabla f(w_k)\|^2$

$$\Leftrightarrow \frac{1}{2L} \sum_{k=0}^{K-1} \|\nabla f(w_k)\|^2 \leq f(w_0) - \bar{f}$$

$$\Leftrightarrow \sum_{k=0}^{K-1} \|\nabla f(w_k)\|^2 \leq 2L(f(w_0) - \bar{f})$$

Using  $\sum_{k=0}^{K-1} \|\nabla f(w_k)\|^2 \geq K \left( \min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \right)^2$

$$\left( \sum_{k=0}^{K-1} a_k^2 \geq \sum_{k=0}^{K-1} (\min a_k)^2 = K (\min_k a_k)^2 \right)$$

Overall, we obtain that

$$K \left( \min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \right)^2 \leq 2L(f(w_0) - \bar{f})$$

$$\min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \leq \sqrt{\frac{2L(f(w_0) - \bar{f})}{K}}$$

Theorem Apply gradient descent to minimize  $f(w)$  with  $w \in \mathbb{R}^d$ ,  $f \in C_L^{1,1}$  and the stepsize for gradient descent chosen as  $\frac{1}{L}$ .

Then, for any  $K \geq 1$ ,

$$\min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \leq \sqrt{\frac{2L(f(w_0) - \bar{f})}{K}} = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

"A constant times  $1/\sqrt{K}$ "

⇒ We say that the convergence rate for gradient descent on nonconvex functions is  $\frac{1}{\sqrt{K}}$

⇒ As  $K \rightarrow \infty$ ,  $\min_{0 \leq k \leq K-1} \|\nabla f(w_k)\| \rightarrow 0$ : the method converges to a point with zero gradient

(Recall that  $w^*$  is a minimum of  $f(w) \Rightarrow \nabla f(w^*) = 0_{\text{nd}}$ )

Remark: Sometimes  $K$  is predetermined (budget criterion)  
But it can also be selected to achieve some accuracy/  
some convergence guarantee

For example, if we want to compute  $w_K$  with  $\|\nabla f(w_K)\| \leq \varepsilon$ ,  
then choosing  $K \geq 2L(f(w_0) - \bar{f})\varepsilon^{-2}$  guarantees that GD  
will compute such a point

→ If  $L$  and  $\bar{f}$  are known, can use this value in practice  
→ Otherwise, this value is used as a guidance on the  
performance of the method

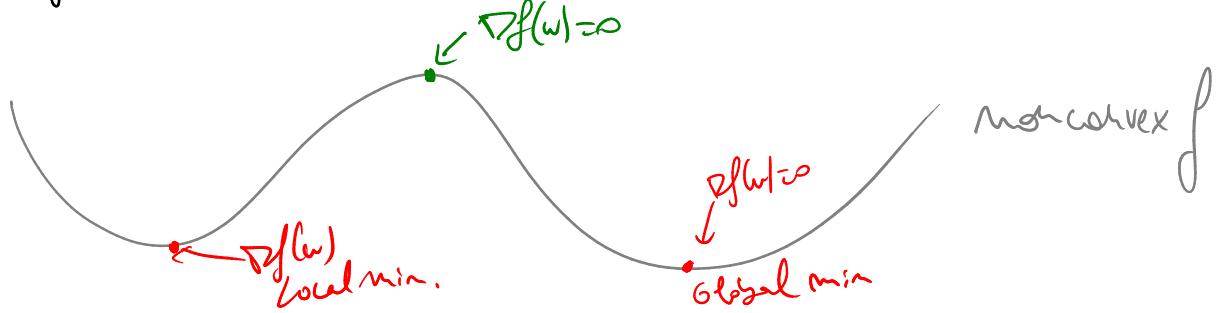
$$\varepsilon = 10^{-6} \rightarrow 2L(f(w_0) - \bar{f})\varepsilon^{-2} \text{ iterations}$$

$$\frac{\varepsilon}{10} = 10^{-7} \rightarrow 2L(f(w_0) - \bar{f})\varepsilon^{-2} \times 100 \text{ iterations}$$

Remark: ( $\|\nabla f(w_k)\| \rightarrow 0$  and nonconvexity)

↪ The theorem shows that GD converges to a point with zero gradient.

In theory, since  $f$  is not assumed to be convex, a point with zero gradient is not necessarily a minimum!



so GD can converge to points that are not minima in theory. In practice, however, this almost never happens, i.e. GD tends to converge to local or even global minima.

⇒ Explained very recently (through sophisticated math arguments)

Theorem (Lee et al. 2015) !INFORMAL VERSION

Suppose that we run GD with  $w_0$  chosen randomly in  $\mathbb{R}^d$ . Then, with probability 1, GD converges to a local minimum.

- Result actually holds for a subclass of nonconvex functions
- This subclass includes classical problems in ML
  - . Low-rank matrix completion
  - , Shallow neural networks
  - . Phase retrieval

→ See the homework!

# Convergence rates for gradient descent

- ↳ GD can be applied to any  $C^1$  function, regardless of the function being convex or nonconvex  
=> All you need is a gradient and a stepsize
- ↳ when the function is convex, GD has stronger theoretical guarantees.
- ↳ when the function is strongly convex, GD has even better guarantees.

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

- Any  $\mu$ -strongly convex function is convex
- The converse is not true:  $w \mapsto \bar{x}^Tw + \gamma$  is convex but not strongly convex



$\rightarrow$  If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C^1_L$  and  $\mu$ -strongly convex, then  $\nabla f(v, w) \in \mathbb{R}^{nd}$

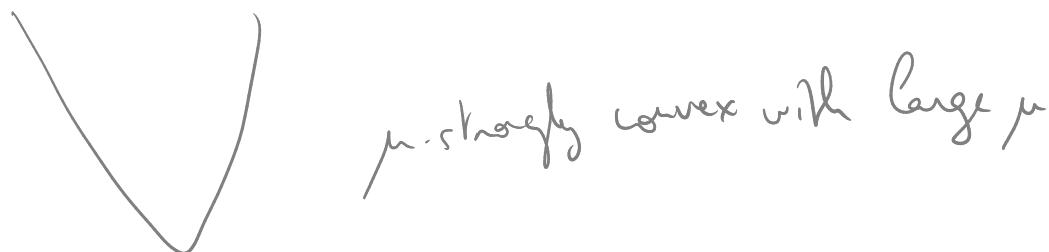
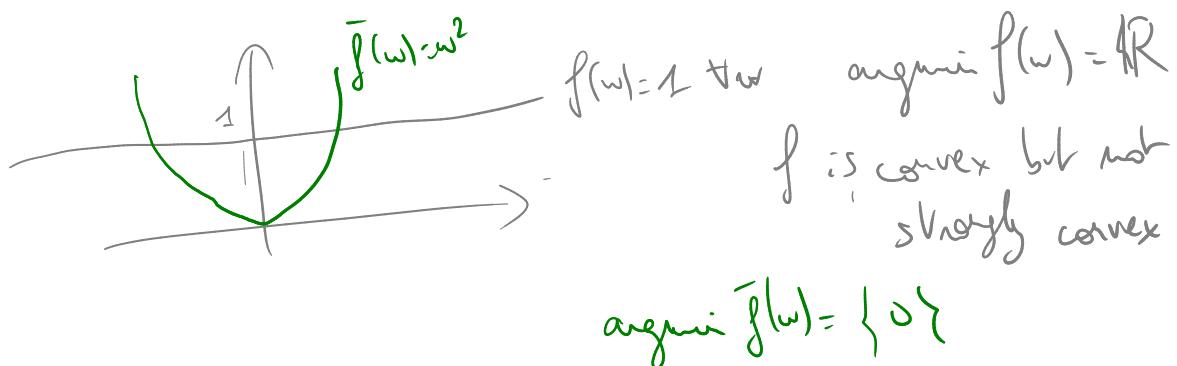
$$f(w) + \nabla f(w)^T(v-w) + \frac{\mu}{2}\|v-w\|^2 \leq f(v) \leq f(w) + \nabla f(w)^T(v-w) + \frac{L}{2}\|v-w\|^2$$

$\hookrightarrow$  At every  $w$ , you can lower and upper bound  $f$  by quadratic functions

$\hookrightarrow$  Implies  $L \geq \mu$

Proposition: If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C^1$  and  $\mu$ -strongly convex, it has a unique minimum, which is the only point  $w^* \in \mathbb{R}^d$  such that  $\nabla f(w^*) = 0_{\mathbb{R}^d}$

$\rightarrow$  Convex functions can have infinitely many global minima (ex: constant functions)



NB: If  $f$  ( $\mu$ -strongly convex) had two global minima  $w_1$  and  $w_2$ , we would have

$$f(w_2) \geq f(w_1) + \underbrace{\nabla f(w_1)^T(w_2 - w_1)}_{> 0} + \underbrace{\frac{\mu}{2} \|w_2 - w_1\|^2}_{\text{("} w_1 \text{ minimum")}}$$

$$f(w_2) \geq f(w_1) + \underbrace{\frac{\mu}{2} \|w_2 - w_1\|^2}_{> 0 \text{ if } w_2 \neq w_1} > f(w_1) \quad \text{if } w_1 \neq w_2$$

This contradicts the fact that  $w_2 \in \arg\min_w f(w)$

Q) What is the convergence rate of GD on convex functions and strongly convex functions?

Setup: minimize  $f(w)$ ,  $f \in \mathcal{L}^1$ , run GD with stepsize  $\frac{1}{L}$  for  $K \geq 1$  iterations.

(1) For nonconvex functions, we showed above that

$$\min_{0 \leq k \leq K_1} \|\nabla f(w_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

$\mathcal{O}(A)$ : a constant times A

(2) For convex functions, we can show that

$$f(w_K) - \min_{w \in \mathbb{R}^d} f(w) \leq \mathcal{O}\left(\frac{1}{K}\right)$$

This rate is better than the rate for nonconvex functions!

•  $\frac{1}{K}$  decreases faster than  $\frac{1}{\sqrt{K}}$

•  $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \rightarrow 0$  is a stronger guarantee than  $\|\nabla f(w_k)\| \rightarrow 0$

③ If  $f$  is  $\mu$ -strongly convex, we can also show that

$$f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left((1 - \frac{\mu}{L})^k\right)$$

- Same guarantee than convex functions:  $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \rightarrow 0$
- But much better rate:  $(1 - \frac{\mu}{L})^k$  vs  $\frac{1}{K}$

Ex)  $\mu=1, L=2, K=100$

$$(1 - \frac{\mu}{L})^k = \left(\frac{1}{2}\right)^{100} = \frac{1}{2^{100}} \quad \text{vs} \quad \frac{1}{K} = \frac{1}{100}$$

- Even stronger guarantee:

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

where  $\{w^*\} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w)$

$w_k - w^* \rightarrow 0$ : The iterates of GD converge to the minimum of  $f$ !

## Takeaways

→ Strongly convex functions are easier to minimize than convex functions, and convex functions are easier to minimize than nonconvex ones!

→ For the same algorithm (GD in our case), the difference between these 3 classes of functions can be seen in the convergence rates

→ Convergence rates are a modern tool to compare optimization algorithms and choose the fastest method for a given class of functions.

## The fastest methods for minimizing $C_L^{1+}$ functions

- If  $f$  is nonconvex, the fastest rate is  $\frac{1}{\sqrt{K}}$ , which is the rate of gradient descent.
- If  $f$  is convex, the fastest rate is  $\frac{1}{K^2}$ , which is faster than the  $\frac{1}{K}$  rate of gradient descent.
- If  $f$  is  $\mu$ -strongly convex, the fastest rate is  $(1 - \sqrt{\frac{\mu}{L}})^K$ , which is faster than the  $(1 - \frac{\mu}{L})^K$  rate of gradient descent.

$$(L \geq \mu \Rightarrow \frac{\mu}{L} \in (0, 1] \Leftrightarrow \sqrt{\frac{\mu}{L}} \geq \frac{\mu}{L} \Leftrightarrow 1 - \sqrt{\frac{\mu}{L}} \leq 1 - \frac{\mu}{L})$$

→ GD is not the fastest method for convex/strongly convex  $C_L^{1+}$  functions

⇒ This result was shown before the fastest methods were discovered.

⇒ For a while, there did not exist an algorithm with the best convergence rates

⇒ Today, we know of such an algorithm: accelerated gradient, aka Nesterov's method

↳ History of accelerated gradient

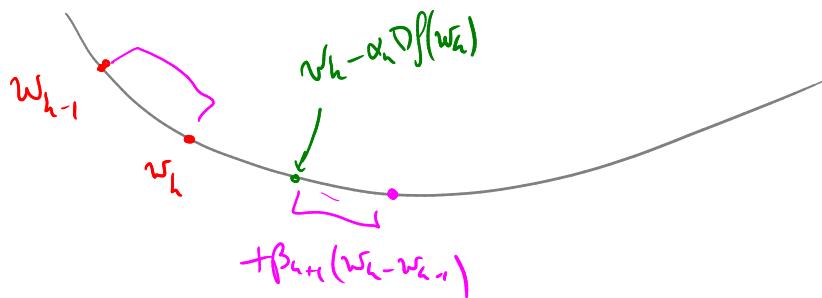
\* 1964: Polyak's method (Heavy ball / Gradient descent with momentum)

$$GD: \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad (\text{typically } \alpha_k = \frac{1}{L})$$

Heavy ball:  $w_{h+1} = w_h - \alpha_h \nabla f(w_h) + \beta_{h+1} (w_h - w_{h-1})$

with  $\beta_{h+1} \geq 0, \quad w_0 = w_0$

Idea: Incorporate information from the previous iteration ( $w_{h-1} \rightarrow w_h$ ) into the calculation of  $w_{h+1}$   
 $\Rightarrow$  Tilted Momentum term  $\beta_{h+1} (w_h - w_{h-1})$



↳ Polyak showed that this method is the fastest for  $C_2^{1,1}$ , strongly convex quadratic functions (special case)

→ Unfortunately, Polyak's method can fail on other convex functions.

\* 1983: Yurii Nesterov introduce the accelerated gradient method

$$w_{h+1} = w_h - \alpha_h \nabla f(w_h + \beta_{h+1} (w_h - w_{h-1})) + \beta_{h+1} (w_h - w_{h-1})$$

- Still combine gradient steps and momentum like in heavy ball
- But you add the momentum before computing the gradient  
 $\Rightarrow$  Nesterov's method can be written as

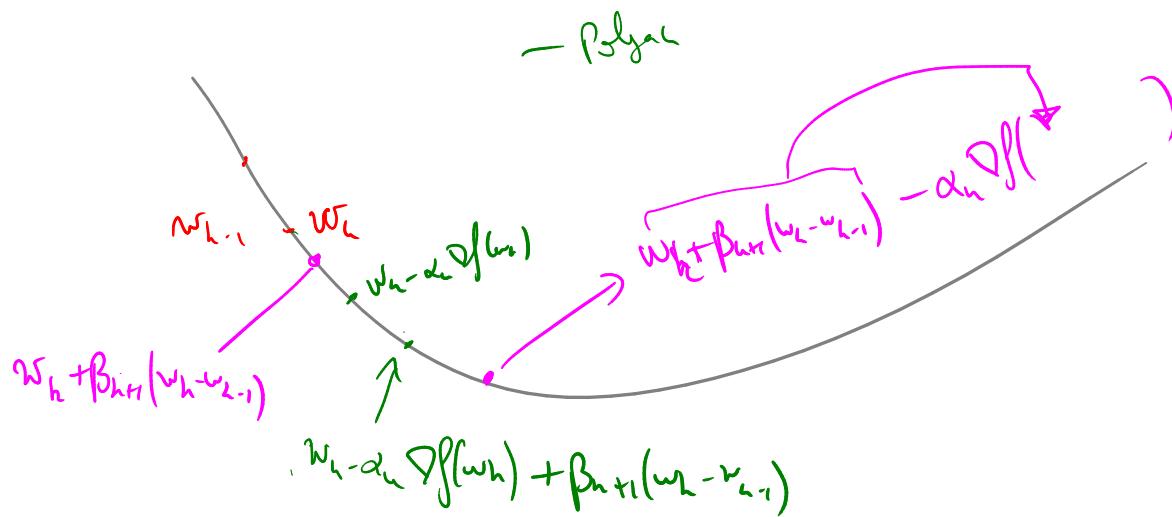
$$\bar{z}_{k+1} = w_k + \beta_{k+1} (w_k - w_{k-1})$$

$$w_{k+1} = \bar{z}_{k+1} - \alpha_k \nabla f(\bar{z}_{k+1})$$

$\hookrightarrow$  Nesterov showed that this method has the fastest convergence rate possible for convex and strongly convex functions.

$\triangle$   $\alpha_k$  and  $\beta_{k+1}$  need to be chosen carefully

$$\mu\text{-Strongly convex: } \alpha_k = \frac{1}{L} \quad \beta_{k+1} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$$



convex case:  $\alpha_k = \frac{1}{L}$ ,  $\beta_{k+1}$  chosen independently of  $f$  and  $L$ !