

The subgradient method

Acknowledgement: slides are based on Prof. Lieven Vandenberghes.

- subgradient method
- convergence analysis
- optimal step size when f^* is known
- alternating projections
- optimality

Subgradient method

to minimize a nondifferentiable convex function f : choose $x^{(0)}$ and repeat

$$x^{(k)} = x^{(k-1)} - t_k g^{(k-1)}, \quad k = 1, 2, \dots \quad \forall g^{(k-1)} \in \partial f(x^{(k-1)})$$

$g^{(k-1)}$ is any subgradient of f at $x^{(k-1)}$

$$\begin{cases} x^{k+1} = x^k + t_k \underline{g^k} \\ \|\nabla f(x^k) - g^k\| \leq \underline{\varepsilon_k} \end{cases}$$

Step size rules

- fixed step: t_k constant
- fixed length: $t_k \|g^{(k-1)}\|_2 = \|x^{(k)} - x^{(k-1)}\|_2 \stackrel{=s}{\text{is constant}}$
- diminishing: $t_k \rightarrow 0, \sum_{k=1}^{\infty} t_k = \infty$ (Learning rate η_t)

Inexact Gradient descent.

Assumptions

- f has finite optimal value f^* , minimizer x^*

- f is convex, $\text{dom } f = \mathbf{R}^n$

(∇f)

- f is Lipschitz continuous with constant $G > 0$:

$$\underbrace{|f(x) - f(y)| \leq G\|x - y\|_2}_{\Leftrightarrow} \quad \forall x, y$$

this is equivalent to $\|g\|_2 \leq G$ for all x and $g \in \partial f(x)$ (see next page)

$$\underbrace{\sup\{\|y\| \mid y \in \partial f(x), \forall x\}} \leq G.$$

Proof.

- assume $\|g\|_2 \leq G$ for all subgradients; choose $g_y \in \partial f(y)$, $g_x \in \partial f(x)$:

$$g_x^T(x - y) \geq f(x) - f(y) \geq g_y^T(x - y)$$

by the Cauchy-Schwarz inequality

$$G\|x - y\|_2 \geq f(x) - f(y) \geq -G\|x - y\|_2$$

- assume $\|g\|_2 > G$ for some $g \in \partial f(x)$; take $y = x + g/\|g\|_2$: $\Rightarrow \|y - x\| = 1$

$$\begin{aligned} |f(y) - f(x)| &\leq G\|y - x\| & f(y) &\geq \underline{f(x) + g^T(y - x)} \\ & & &= f(x) + \underline{\|g\|_2} \\ & & &> \underline{f(x) + G} \end{aligned}$$

$$f(y) - f(x) > G\|y - x\|$$

$$2 \sum_{i=1}^k t_i (f_{\text{best}}^{(k)} - f^*) \leq 2 \sum_{i=1}^k t_i (f(x_i) - f^*)$$

Analysis

- the subgradient method is not a descent method

- the key quantity in the analysis is the distance to the optimal set

with $x^+ = x^{(i)}$, $x = x^{(i-1)}$, $g = g^{(i-1)}$, $t = t_i$:

$$\|x^+ - x^*\|_2^2 = \|x - tg - x^*\|_2^2$$

$$\underbrace{f_{\text{best}}^k}_{\Delta} = \min \{ f_{\text{best}}^{k-1}, f^k \} \leq \|x - x^*\|_2^2 - 2tg^T(x - x^*) + t^2\|g\|_2^2$$

$$\leq \|x - x^*\|_2^2 - 2t(f(x) - f^*) + t^2\|g\|_2^2$$

combine inequalities for $i = 1, \dots, k$, and define $f_{\text{best}}^{(k)} = \min_{0 \leq i < k} f(x^{(i)})$:

$$\underbrace{\Delta}_{\Delta} 2 \left(\sum_{i=1}^k t_i \right) (f_{\text{best}}^{(k)} - f^*) \leq \|x^{(0)} - x^*\|_2^2 - \underbrace{\|x^{(k)} - x^*\|_2^2}_{\Delta} + \sum_{i=1}^k \underbrace{t_i^2 \|g^{(i-1)}\|_2^2}_{\Delta}$$

$$\leq \|x^{(0)} - x^*\|_2^2 + \sum_{i=1}^k t_i^2 \|g^{(i-1)}\|_2^2$$

$$\Rightarrow \underbrace{f_{\text{best}}^{(k)} - f^*}_{\Delta} \leq \frac{\|x^{(0)} - x^*\|_2^2 + \sum_{i=1}^k t_i^2 G^2}{2 \sum_{i=1}^k t_i}$$

$$\rightarrow 0 ?$$

Fixed step size: $t_i = t$

$$f_{\text{best}}^{(k)} - f^* \leq \frac{\|x^{(0)} - x^*\|_2^2}{\underbrace{2kt}_{\downarrow 0}} + \underbrace{\frac{G^2 t}{2}}_{\rightarrow 0}$$

- does not guarantee convergence of $f_{\text{best}}^{(k)}$

- for large k , $f_{\text{best}}^{(k)}$ is approximately $G^2 t / 2$ -suboptimal $O(G^2)$. $\left. \begin{array}{c} \rightarrow \\ \downarrow \end{array} \right\} f^* + \frac{G^2 t}{2}$

Fixed step length: $t_i = s / \|g^{(i-1)}\|_2$ $t_i \|g^{(i-1)}\|_2 = s$

$$f_{\text{best}}^{(k)} - f^* \leq \frac{G \|x^{(0)} - x^*\|_2^2}{2ks} + \frac{Gs}{2}$$

- does not guarantee convergence of $f_{\text{best}}^{(k)}$

- for large k , $f_{\text{best}}^{(k)}$ is approximately $Gs/2$ -suboptimal $O(G)$

$$x^{k+1} = x^k - s \frac{g^k}{\|g^k\|_2}$$

Diminishing step size: $t_i \rightarrow 0, \sum_{i=1}^{\infty} t_i = \infty$ $O\left(\frac{1}{k}\right)$ $O\left(\frac{1}{\sqrt{k}}\right)$

$$f_{\text{best}}^{(k)} - f^* \leq \frac{\underbrace{\|x^{(0)} - x^*\|_2^2}_{= R^2} + G^2 \sum_{i=1}^k t_i^2}{2 \sum_{i=1}^k t_i} \rightarrow 0$$

$$t_k \sim O\left(\frac{1}{k}\right)$$

$$t_k \Rightarrow \frac{3}{k}$$

$$\downarrow \frac{3000}{k}$$

$$\underline{\underline{\Delta}}$$

can show that $\left(\sum_{i=1}^k t_i^2\right) / \left(\sum_{i=1}^k t_i\right) \rightarrow 0$; hence, $f_{\text{best}}^{(k)}$ converges to f^*

Proof: $\forall \epsilon > 0, \exists N_1$ s.t. $t_i \leq \frac{\epsilon}{G^2}, \forall i > N_1$

$\exists N_2$ s.t. $\sum_{i=1}^{N_2} t_i \geq \frac{1}{\epsilon} \left(R^2 + G^2 \sum_{i=1}^{N_2} t_i^2 \right)$

Let $N = \max\{N_1, N_2\}$, For $k > N$.

$$\frac{R^2 + G^2 \sum_{i=1}^k t_i^2}{2 \sum_{i=1}^k t_i} = \frac{R^2 + G^2 \sum_{i=1}^{N_1} t_i^2}{2 \sum_{i=1}^k t_i} + \frac{G^2 \sum_{i=N_1+1}^k t_i^2}{2 \sum_{i=1}^k t_i} \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon$$

$G^2 \sum_{i=N_1+1}^k t_i^2 \leq t_i \left(\frac{\epsilon}{G^2}\right)$

Example: 1-norm minimization

$$\text{minimize } \|Ax - b\|_1$$

$$f(x) = h(Ax - b) \rightarrow A(x) = \|x\|_1$$

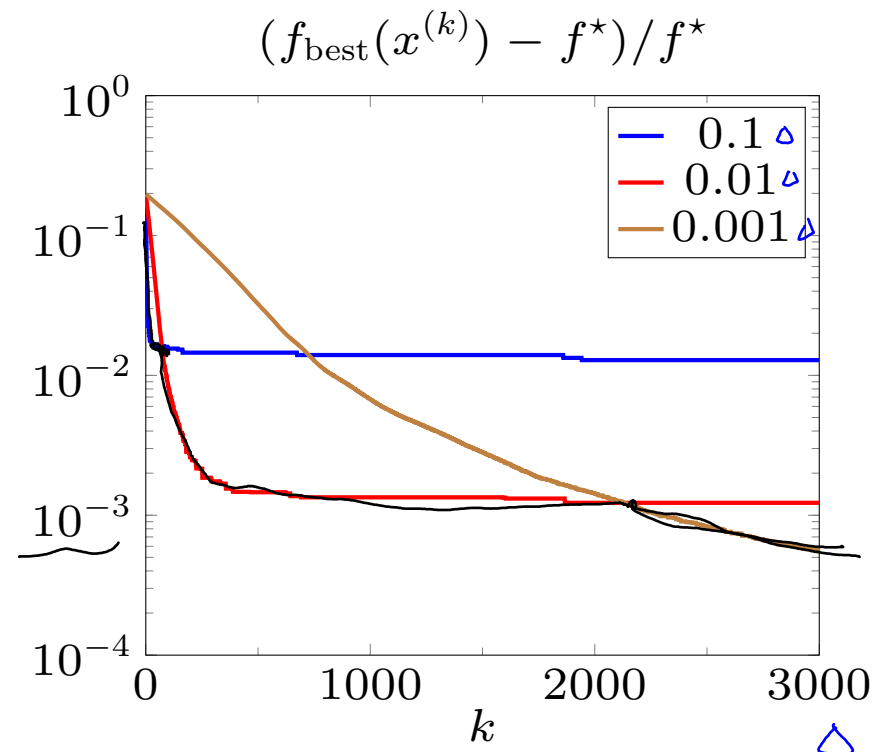
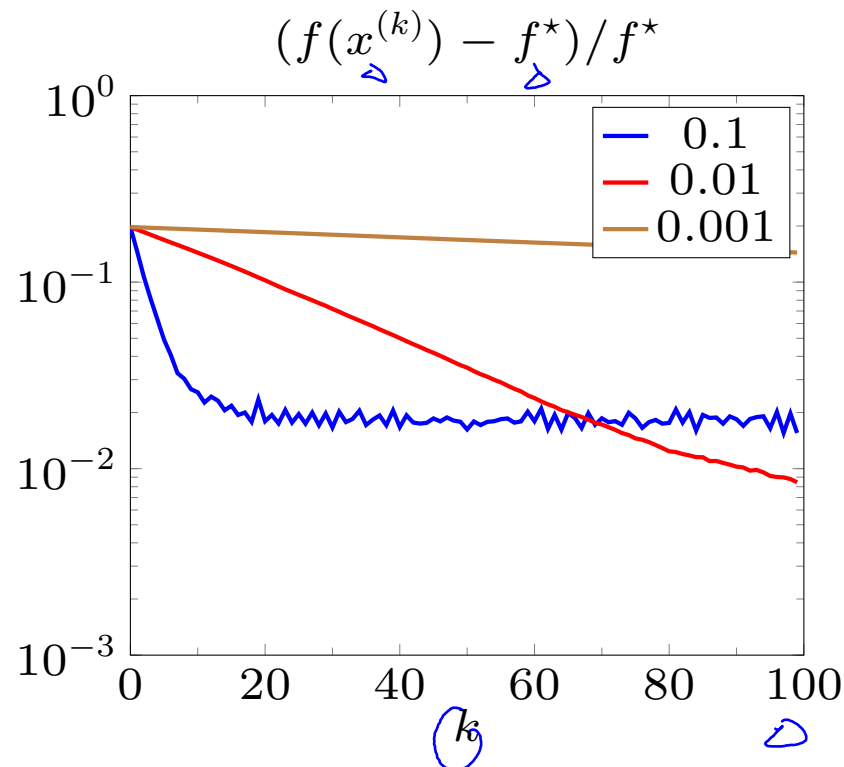
$$\Rightarrow \partial f(x) = A^T \partial h(Ax - b)$$

• subgradient is given by $A^T \text{sign}(Ax - b)$

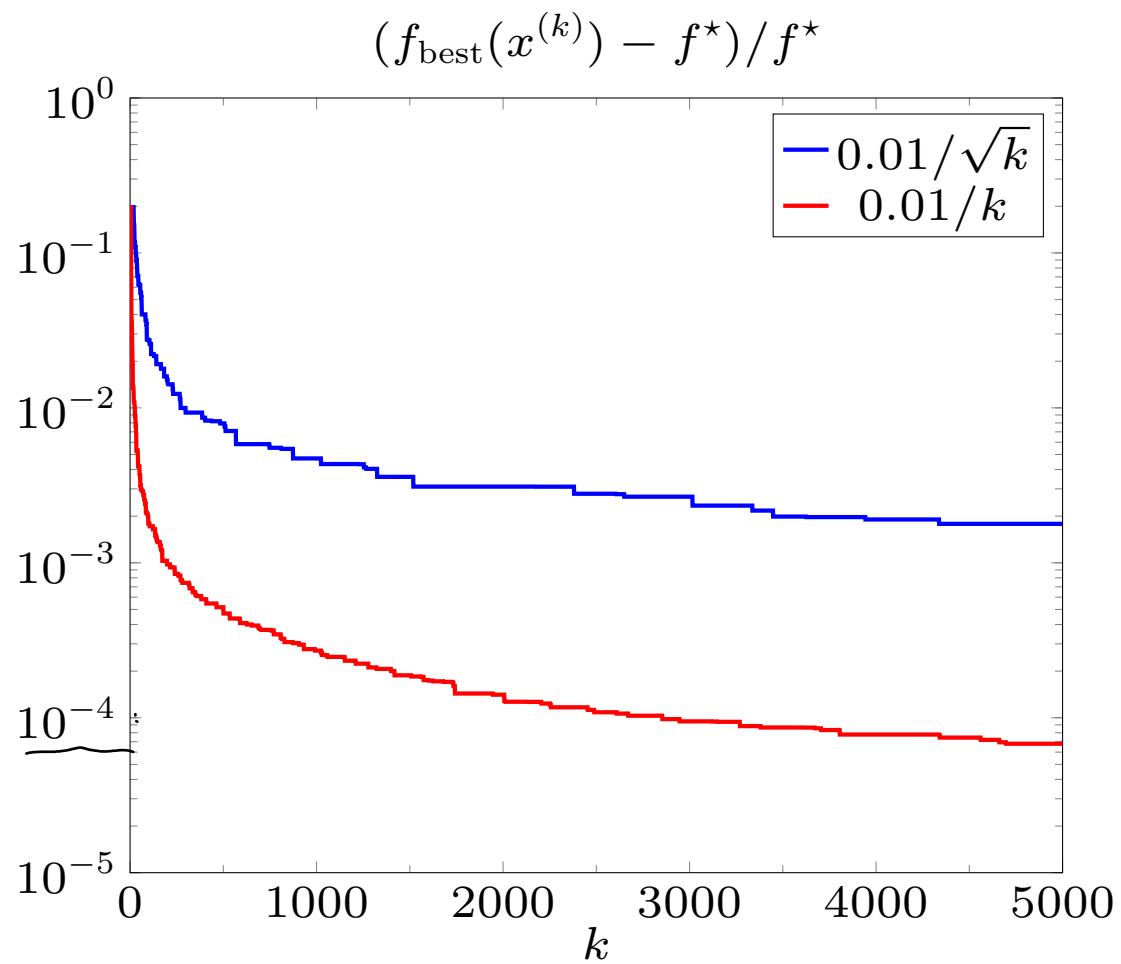
• example with $A \in \mathbf{R}^{500 \times 100}$, $b \in \mathbf{R}^{500}$

Fixed steplength $t_k = s / \|g^{(k-1)}\|_2$ for $s = 0.1, 0.01, 0.001$

$$\frac{GS}{2}$$



Diminishing step size: $t_k = 0.01/\sqrt{k}$ and $t_k = 0.01/k$



Optimal step size for fixed number of iterations

from page 5-5: if $s_i = t_i \|g^{(i-1)}\|_2$ and $\|x^{(0)} - x^*\|_2 \leq R$:

$$f_{\text{best}}^{(k)} - f^* \leq \frac{R^2 + \sum_{i=1}^k s_i^2}{2 \sum_{i=1}^k s_i / G}$$

$$\min_{s>0} \frac{R^2}{2ks/G} + \frac{s}{2/G}$$

$$\Leftrightarrow \frac{R^2}{2ks} = \frac{s}{2}$$

$$\Rightarrow s = \frac{R}{\sqrt{k}}$$

- for given k , bound is minimized by fixed step length $s_i = s = R/\sqrt{k}$
- resulting bound after k steps is

$$G.D.: \underline{\underline{O\left(\frac{1}{k}\right)}}$$

$$\boxed{f_{\text{best}}^{(k)} - f^* \leq \frac{GR}{\sqrt{k}}}$$

- guarantees accuracy $f_{\text{best}}^{(k)} - f^* \leq \epsilon$ in $k = \underline{\underline{O(1/\epsilon^2)}}$ iterations

Optimal step size when f^* is known

- right-hand side in first inequality of page 5-5 is minimized by

$$t_i = \frac{f(x^{(i-1)}) - f^*}{\|g^{(i-1)}\|_2^2}$$

Quadratic w.r.t. t_i

$$\|x^i - x^*\|_2^2 \leq \|x^{i-1} - x^*\|_2^2 - 2t_i(f^{i-1} - f^*) + t_i^2 \|g^i\|^2$$

min

- optimized bound is

$$\frac{(f(x^{(i-1)}) - f^*)^2}{\|g^{(i-1)}\|_2^2} \leq \|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2$$

- applying recursively (with $\|x^{(0)} - x^*\|_2 \leq R$ and $\|g^{(i)}\|_2 \leq G$) gives

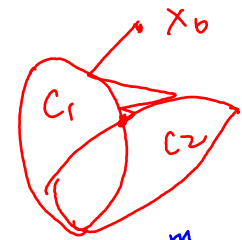
$$f_{\text{best}}^{(k)} - f^* \leq \frac{GR}{\sqrt{k}}$$

$O\left(\frac{1}{\sqrt{k}}\right) \rightarrow \text{best?}$

Exercise: find point in intersection of convex sets

Alternating projection.

find a point in the intersection of m closed convex sets C_1, \dots, C_m :



$$\text{minimize } \underline{f(x) = \max \{f_1(x), \dots, f_m(x)\}}$$

Find $x \in \bigcap_{j=1}^m C_j$

where $f_j(x) = \inf_{y \in C_j} \|x - y\|_2$ is Euclidean distance of x to C_j
 ≥ 0

- $\underline{f^* = 0}$ if the intersection is nonempty
- (from p. 4-14): $g \in \partial f(\hat{x})$ if $g \in \partial f_j(\hat{x})$ and C_j is farthest set from \hat{x}
- (from p. 4-20) subgradient $g \in \partial f_j(\hat{x})$ follows from projection $P_j(\hat{x})$ on C_j :

$$f_j(\hat{x}) > 0$$

$$g = 0 \quad (\text{if } \hat{x} \in C_j), \quad g = \frac{1}{\|\hat{x} - P_j(\hat{x})\|_2} (\hat{x} - P_j(\hat{x})) \quad (\text{if } \hat{x} \notin C_j)$$

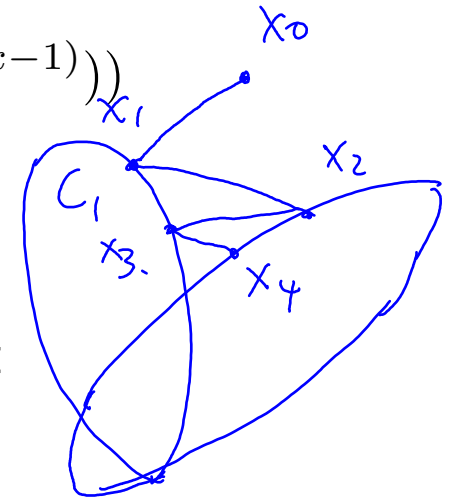
note that $\|g\|_2 = 1$ if $\hat{x} \notin C_j$

Subgradient method

$$t_i = \frac{f^{i-1} - f^*}{\|g^{i-1}\|_2} = 0$$

- optimal step size (page 5-11) for $f^* = 0$ and $\|g^{(i-1)}\|_2 = 1$ is $t_i = f(x^{(i-1)})$
- at iteration k , find farthest set C_j (with $f(x^{(k-1)}) = f_j(x^{(k-1)})$), and take

$$\begin{aligned} \underline{x^{(k)}} &= x^{(k-1)} - \frac{f(x^{(k-1)})}{f_j(x^{(k-1)})} (x^{(k-1)} - P_j(x^{(k-1)})) \\ &= \underline{P_j(x^{(k-1)})} \end{aligned}$$



at each step, we project the current point onto the farthest set

- a version of the alternating projections algorithm $O(\frac{1}{\sqrt{k}})$
- for $m = 2$, projections alternate onto one set, then the other $\bigcap_{i=1}^m \{x \mid \underbrace{f_i(x)}_{A_i x + b_i} \leq 0\}$
- later, we will see faster versions of this that are almost as simple

$$\underline{\{f(x) \leq 0\} \text{ feasible.}}$$

Optimality of the subgradient method

can the $f_{\text{best}}^{(k)} - f^* \leq GR/\sqrt{k}$ bound on page 5-10 be improved?

Problem class

- f is convex, with a minimizer x^* ✓
- we know a starting point $x^{(0)}$ with $\|x^{(0)} - x^*\|_2 \leq R$ ✓
- we know the Lipschitz constant G of f on $\{x \mid \|x - x^{(0)}\|_2 \leq R\}$
- f is defined by an oracle: given x , oracle returns $f(x)$ and a subgradient

Algorithm class: k iterations of any method that chooses $x^{(i)}$ in

$$\begin{aligned} & \underline{x^{(0)} + \text{span}\{g^{(0)}, g^{(1)}, \dots, g^{(i-1)}\}} \\ & x^{k+1} = x^k + t_k g^k \\ & \quad = x^{k-1} + t_{k-1} g^{k-1} + t_k g^k. \end{aligned}$$

Test problem and oracle

$$= \max_{1 \leq i \leq k} \{e_i^T x\} + \frac{1}{2} \|x\|_2^2$$

$$f(x) = \max_{i=1, \dots, k} x_i + \frac{1}{2} \|x\|_2^2, \quad x^{(0)} = 0$$

$$0 \in \partial f(x) = \text{conv}\{e_j \mid j \in I(x)\} + x$$

$$0 \in \partial f(x^*)$$

- solution: $x^* = -\frac{1}{k}(\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{n-k})$ and $f^* = -\frac{1}{2k}$

- $R = \|x^{(0)} - x^*\|_2 = 1/\sqrt{k}$ and $G = 1 + 1/\sqrt{k}$

- oracle returns subgradient $\underline{e}_{\hat{j}} + x$ where $\hat{j} = \min\{j \mid x_j = \max_{i=1, \dots, k} x_i\}$

Iteration: for $i = 0, \dots, k-1$, entries $\underline{x}_{i+1}^{(i)}, \dots, x_k^{(i)}$ are zero; therefore

$$f_{\text{best}}^{(k)} - f^* = \min_{i < k} f(x^{(i)}) - f^* \geq -f^* = \frac{GR}{2(1 + \sqrt{k})}$$

$f(x^{(i)}) \geq 0$

Conclusion: $O(1/\sqrt{k})$ bound cannot be improved

$$g^{(0)} = e_1 + x^{(0)} = e_1 \Rightarrow t_1 = \frac{f^0 - f^*}{\|g^{(0)}\|_2^2}$$

Subgradient method

$$\Rightarrow x^1 = x^{(0)} - t_1 e_1 \Rightarrow x^2 = x^{(0)} - t_2 (e_1 + x^{(1)})$$

Summary: subgradient method

$$\|x^k - x^{k-1}\| < \epsilon \quad (X)$$

- handles general nondifferentiable convex problem

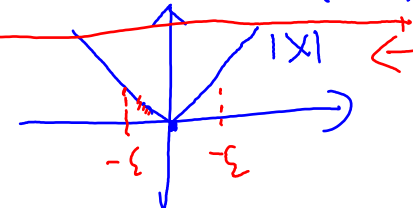
- often leads to very simple algorithms

- convergence can be very slow

→ Acceleration !!!

- no good stopping criterion

$$(\|\nabla f(x)\| \leq \epsilon)$$



if $x \neq 0$
 $\partial |x| = \text{sgn}(x)$
 $\Rightarrow \|\partial |x|\| = 1$

- theoretical complexity: $O(1/\epsilon^2)$ iterations to find ϵ -suboptimal point

practical
 stopping
 criterion !!!

- an 'optimal' 1st-order method: $O(1/\epsilon^2)$ bound cannot be improved

References

- S. Boyd, lecture notes and slides for EE364b, Convex Optimization II
- Yu. Nesterov, *Introductory Lectures on Convex Optimization. A Basic Course* (2004)

§3.2.1 with the example on page 5-15 of this lecture