# Exact convergence rates of the last iterate in subgradient methods

François Glineur and Moslem Zamani



Information and Communication Technologies, Electronics and Applied Mathematics Institute, and Center for Operations Research and Econometrics UCLouvain

Results from preprint https://arxiv.org/abs/2307.11134

ALGOPT2024

August 30, 2024, LLN

Last-iterate convergence

Performance estimation

Convergence rates

Extensions

Last-iterate optimal subgradient method

Normalized step sizes

Conclusions

*Objective*: minimize a function  $f : \mathbb{R}^d \to \mathbb{R}$  that is

convex

 $\partial f(x) = \{g \text{ such that } f(y) \ge f(x) + g^T(y - x) \text{ for all } y\} \neq \emptyset$ 

B-Lipschitz continous

$$g \in \partial f(x) \Rightarrow \|g\| \leq B$$

▶ with minimizer x\*

*Method*: subgradient method with fixed step sizes  $\{h_k\}$  $x_{k+1} = x_k - h_k g_k$  for some  $g_k \in \partial f(x_k)$ 

starting from  $x_0$ 

## Performance criteria

#### *Target*: convergence rate after *N* iterations, either

Initial iterate assumption:

$$\|x_0-x_*\|\leq R$$

Homogeneity: rates in function values must be proportional to BR

Lower bound: no method can achieve better rate than

 $\frac{BR}{\sqrt{N+1}}$ 

#### Lower bound proof (variation of [Drori, Teboulle 2016])

Consider following function with d = N + 1, B = 1 and  $x_* = 0$  $f(x) = \max\{0, x_1, x_2, \dots, x_{N+1}\} = \begin{bmatrix}\max_{1 \le k \le N+1} x_k\end{bmatrix}_+$ 

Choose starting point  $x_0 = (1, 1, ..., 1)$  with  $R = \sqrt{N+1}$ 

- As long as f(x<sub>k</sub>) > 0, subgradient g<sub>k</sub> ∈ ∂f(x<sub>k</sub>) can be chosen as a basis vector e<sub>i</sub> for some 1 ≤ i ≤ N + 1 (and ||g<sub>k</sub>|| = B)
- ► Induction hypothesis (*H<sub>k</sub>*) (easy to check for *k* = 0) *x<sub>k</sub>* contains at least *N* + 1 − *k* components equal to 1
- ► Assume (*H<sub>k</sub>*) for *k* ≤ *N*. Then *f*(*x<sub>k</sub>*) ≥ 1. So subgradient *g<sub>k</sub>* can be chosen as some basis vector *e<sub>i</sub>*, and *x<sub>k+1</sub>* can differ only by at most one component from *x<sub>k</sub>*, implying (*H<sub>k+1</sub>*) holds

Conclusion: 
$$f(x_k) \ge 1 = \frac{BR}{\sqrt{N+1}}$$
 for all  $0 \le k \le N$ 

(also for other criteria / for steps with several past subgradients)

Only two ingredients:

(1) subgradient inequality and (2) square distance telescoping

Ingredient (1)

$$\begin{split} \|x^{k+1} - x^{\star}\|^{2} &= \|x^{k} - h_{k}g^{k} - x^{\star}\|^{2} \\ &= \|x^{k} - x^{\star}\|^{2} + h_{k}^{2}\|g^{k}\|^{2} - 2h_{k}\langle g^{k}, x^{k} - x^{\star}\rangle \\ &\leq \|x^{k} - x^{\star}\|^{2} + h_{k}^{2}\|g^{k}\|^{2} - 2h_{k}\left(f(x^{k}) - f(x^{\star})\right). \end{split}$$

(where we have only used *subgradient inequality*  $f(x^*) - f(x^k) \ge \langle g^k, x^* - x^k \rangle$  between  $x^*$  and  $x^k$ )

This gives an upper bound on the accuracy  $f(x^k) - f(x^*)$ 

$$h_k(f(x^k) - f(x^*)) \le \frac{1}{2} ||x^k - x^*||^2 - \frac{1}{2} ||x^{k+1} - x^*||^2 + \frac{1}{2} h_k^2 B^2$$

using bound on subgradient norm  $||g_k|| \leq B$ 

Ingredient (2) From

$$h_k(f(x^k) - f(x^*)) \le \frac{1}{2} ||x^k - x^*||^2 - \frac{1}{2} ||x^{k+1} - x^*||^2 + \frac{1}{2} h_k^2 B^2$$

*telescoping* (summing from k = 0 to k = N) gives

$$\sum_{k=0}^{N} h_k (f(x^k) - f(x^*)) \le \frac{1}{2} \left\| x^0 - x^* \right\|^2 - \frac{1}{2} \| x^{N+1} - x^* \|^2 + \frac{1}{2} B^2 \sum_{k=0}^{N} h_k^2$$

hence

$$\min_{0 \le k \le N} f(x^k) - f(x^*) \le \frac{\frac{1}{2} \|x^0 - x^*\|^2 + \frac{1}{2} B^2 \sum_{k=0}^N h_k^2}{\sum_{k=0}^N h_k}$$

#### Standard convergence analysis (end.)

$$\min_{0 \le k \le N} f(x^k) - f(x^*) \le \frac{\frac{1}{2} \|x^0 - x^*\|^2 + \frac{1}{2} B^2 \sum_{k=0}^N h_k^2}{\sum_{k=0}^N h_k}$$

Right-hand side is convex and symmetric in stepsizes h<sub>k</sub>, hence optimal values are constant h<sub>k</sub> = h for all k

$$\min_{0 \le k \le N} f(x^k) - f(x^*) \le \frac{\frac{1}{2} \|x^0 - x^*\|^2 + \frac{1}{2} B^2 (N+1) h^2}{(N+1)h}$$

► Optimal *h* is then  $h_k = \frac{R}{B} \frac{1}{\sqrt{N+1}}$  leading to an *optimal* rate  $\min_{0 \le k \le N} f(x^k) - f(x^*) \le \frac{BR}{\sqrt{N+1}}$ (same rate holds for average iterate since  $f(\frac{1}{N+1} \sum_{k=0}^{N} x_k) \le \frac{1}{N+1} \sum_{k=0}^{N} f(x_k)$ )

#### End of story?

#### What about last-iterate convergence?

$$\min_{0 \le k \le N} f(x^k) - f(x^*) \le \frac{BR}{\sqrt{N+1}}$$

- Says nothing about convergence of last iterate x<sub>N</sub>
- O. Shamir, Open problem: Is averaging needed for strongly convex stochastic gradient descent? JMLR (2012)
- Practitioners often use the last iterate
- Storing best iterate might not be feasible (storage requirements, objective computation)
- Algorithm may correspond to a real-word dynamical system (see for example work by Nesterov and Shikhman)

*Goal of this talk*: study last-iterate convergence with and without performance estimation

#### Short history of our results

- 2012-2013: Drori and Teboulle introduce performance estimation problems (PEP) main idea: compute worst-case convergence rates
- 2013-2017: with Taylor and Hendrickx we further develop SDP-based PEP approach
- 2017: Yurii asks us "With your tool, can you tell the convergence rate of the last iterate in subgradient method?" We find a purely numerical rate (see next page), and no proof
- 2023: with Zamani we get back to the question and obtain a full PEP proof and a bit later a classic proof

Puzzle: can you guess the convergence rate?

For constant stepsize h = 1 one can compute using either PESTO (Matlab) or PEPIT (Python) toolboxes

$$f(x_N) - f(x_*) \leq BR\left[1 - N + \frac{1}{2}\left(s_N - s_N^{-1}\right)^2\right]$$

where the rate involves a mysterious sequence  $\{s_k\}$ :

$$s_0 = 1, \ s_1 = 2, \ s_2 = 2.5, \ s_3 = 2.9,$$

#### **Puzzle time**

Puzzle: can you guess the convergence rate?

For constant stepsize h = 1 one can compute using either PESTO (Matlab) or PEPIT (Python) toolboxes

$$f(x_N) - f(x_*) \le BR \Big[ 1 - N + \frac{1}{2} (s_N - s_N^{-1})^2 \Big]$$

where the rate involves a mysterious sequence  $\{s_k\}$ :

$$s_0 = 1, \ s_1 = 2, \ s_2 = 2.5, \ s_3 = 2.9, \ s_4 = 3.24482758621, \ \ldots$$

or

$$s_0 = 1, \ s_1 = 2, \ s_2 = \frac{5}{2}, \ s_3 = \frac{29}{10},$$

#### **Puzzle time**

Puzzle: can you guess the convergence rate?

For constant stepsize h = 1 one can compute using either PESTO (Matlab) or PEPIT (Python) toolboxes

$$f(x_N) - f(x_*) \le BR \left[ 1 - N + \frac{1}{2} \left( s_N - s_N^{-1} \right)^2 \right]$$

where the rate involves a mysterious sequence  $\{s_k\}$ :

$$s_0 = 1, \ s_1 = 2, \ s_2 = 2.5, \ s_3 = 2.9, \ s_4 = 3.24482758621, \ \dots$$

or

$$s_0 = 1, \ s_1 = 2, \ s_2 = \frac{5}{2}, \ s_3 = \frac{29}{10}, \ s_4 = \frac{941}{290}, \ \dots$$

#### **Puzzle time**

Puzzle: can you guess the convergence rate?

For constant stepsize h = 1 one can compute using either PESTO (Matlab) or PEPIT (Python) toolboxes

$$f(x_N) - f(x_*) \le BR \left[ 1 - N + \frac{1}{2} \left( s_N - s_N^{-1} \right)^2 \right]$$

where the rate involves a mysterious sequence  $\{s_k\}$ :

$$s_0 = 1, \ s_1 = 2, \ s_2 = 2.5, \ s_3 = 2.9, \ s_4 = 3.24482758621, \ \dots$$

or

$$s_0 = 1, \ s_1 = 2, \ s_2 = \frac{5}{2}, \ s_3 = \frac{29}{10}, \ s_4 = \frac{941}{290}, \ \dots$$

Answer:  $s_{k+1} = s_k + \frac{1}{s_k}$ 

#### Take-home messages:

- Performance estimation applied to subgradient methods
- Exact convergence rates can be obtained for the last iterate: suboptimal by a factor  $O(\sqrt{\log(N)})$
- New last-iterate optimal method can be designed with linearly decreasing step sizes
- Extensions to constrained case, to normalized steps
- Inspiration for results provided by performance estimation but ultimately all proofs converted to classical style using a new key lemma

Last-iterate convergence

Performance estimation

Convergence rates

Extensions

Last-iterate optimal subgradient method

Normalized step sizes

Conclusions

Last-iterate convergence

For a given PEP (Performance Estimation Problem) we can

- compute the exact value of the performance criteria's worst-case = optimal value of PEP problem
- identify an explicit function (and starting point) achieving this worst-case value = primal solution of PEP problem + interpolation
- obtain an independently-checkable proof that this worst-case value is a valid (upper) bound on the performance criteria = dual multiplier of PEP problem
- all three steps can be done either numerically or analytically

For a large class of first-order methods, including fixed-step subgradient methods, these can be computed *exactly* using a semidefinite programming (SDP) problem.

### Interpolation conditions for nonsmooth convex functions

To perform PEP for subgradient methods on a class of functions

we need the corresponding interpolation conditions explicitly

given a list of values  $(x_i, f_i, g_i)_{i \in I}$ , does there exist a convex f with B-bounded subgradients such that  $f(x_i) = f_i$  and  $g_i \in \partial f(x_i)$  for all  $i \in I = \{*, 0, 1, \dots, N\}$ 

Necessary and sufficient conditions:

$$f(x_i) = f_i \text{ and } g_i \in \partial f(x_i) \text{ for every } i \in I$$
  
 $\Leftrightarrow$   
 $f_j \ge f_i + g_i^T(x_j - x_i) \text{ for every } i, j \in I$   
 $\|g_i\| \le B \text{ for every } i \in I$ 

Leads to a convex, tractable formulation as a SDP

#### Results: average iterate with constant stepsize

Worst-case for constant stepsize subgradient method

$$x_{i+1} = x_i - h(\frac{R}{B})g_i$$

applied to convex function with B-bounded subgradients

► For *average* value of iterates  $\hat{f}_N = \frac{f(x_0) + f(x_1) + ... + f(x_N)}{N+1}$ , tight worst-case is

$$\hat{f}_N - f(x_*) \le egin{cases} BRig(rac{1}{2}h + rac{1}{2(N+1)}rac{1}{h}ig) & ext{when } h \ge rac{1}{N+1} \ BRig(1 - rac{N}{2}hig) & ext{when } h \le rac{1}{N+1} \end{cases}$$

(recovers result shown earlier for large h)

▶ Optimal constant step-size is then h<sup>\*</sup> = 1/√N+1 (belongs to "large step" case) leading to tight worst-case

$$\widehat{f}_N - f(x_*) \leq rac{BR}{\sqrt{N+1}}$$

15

#### Results: last iterate with constant stepsize

• Define sequence 
$$\{s_i\}_{i\geq 0} = \{1, 2, \frac{5}{2}, \frac{29}{10}, ...\}$$
 with  
 $s_0 = 1, s_{i+1} = s_i + \frac{1}{s_i}$  for all  $i \geq 0$ 

- No closed form but s<sub>N</sub><sup>2</sup> grows like 2(N + 1) + ½ log(N), also appears in [Nesterov 2009] (again!) for primal-dual subgradient
- ► For value of *last* iterate  $f(x_N)$ , tight worst-case is  $f(x_N) - f(x_*) \le \begin{cases} BR\left[(\frac{1}{2}s_N^2 - N)h + \frac{1}{2s_N^2}\frac{1}{h}\right] & \text{when } h \ge \frac{1}{s_N^2} \\ BR(1 - Nh) & \text{when } h \le \frac{1}{s_N^2} \end{cases}$
- No previous result with correct asymptotic rate for last iterate
   [Harvey,Liaw,Plan,Randhawa 2019] prove a log N/32√N lower bound when B = 1 with stepsize h<sub>i</sub> = 1/√i, and prove a high probability O(log N/√N) upper bound in stochastic case

16

#### Results: optimal stepsize and variants

• To perform N subgradient iterations, optimal stepsize is then  $h^* = \frac{1}{s_N \sqrt{s_N^2 - 2N}}$ 

and corresponding exact worst-case convergence rate becomes

$$f(x_N) - f(x_*) \leq BR \sqrt{1 - rac{2N}{s_N^2}} \lesssim BR \cdot \sqrt{rac{1 + rac{1}{4}\log(N)}{N + 1}}$$

► Using  $h = \frac{1}{\sqrt{N+1}}$  (now known to be suboptimal for last iterate) leads to slightly worse

$$f(x_N) - f(x_*) \leq BR \cdot \left(\frac{\frac{5}{4} + \frac{1}{4}\log(N)}{\sqrt{N+1}}\right)$$

$$\max f^{N+1} - f^*$$
s.t.  $f^i - f^j - \left\langle \frac{B}{Rh} (x^j - x^{j+1}), x^i - x^j \right\rangle \ge 0 \quad i \in \{1, \dots, N+1, \star\}, j \in \{1, \dots, N\}$ 
 $f^i - f^{N+1} - \left\langle g^{N+1}, x^i - x^{N+1} \right\rangle \ge 0 \quad i \in \{1, \dots, N+1, \star\}$ 
 $f^i - f^* \ge 0 \quad i \in \{1, \dots, N+1\}$ 
 $R^2 h^2 - \left\| x^k - x^{k+1} \right\|^2 \ge 0 \quad k \in \{1, \dots, N\}$ 
 $B^2 - \left\| g^{N+1} \right\|^2 \ge 0$ 
 $R^2 - \left\| x^1 - x^* \right\|^2 \ge 0.$ 

# PEP-based proof is ... straightforward?

Define 
$$f^i=f(x^i)$$
 and  $\sigma_i=rac{1}{s_{i+1}},~~i\in\{0,1,...,N\}$  and observe that

$$f^{N+1} - BR\left(\left(\frac{1}{2}s_{N+1}^2 - N\right)h + \frac{1}{2s_{N+1}^2h}\right) + \sum_{i=1}^{N} \frac{B\sigma_{N-i}^2}{2Rh} \left(R^2h^2 - \left\|x^i - x^{i+1}\right\|^2\right)$$

$$+\sum_{i=1}^{N}\sum_{j=i+1}^{N}\sigma_{N-j}(\sigma_{N-i}-\sigma_{N+1-i})(f^{i}-f^{j}-\langle \frac{B}{Rh}(x^{j}-x^{j+1}),x^{i}-x^{j}\rangle)$$

$$+\sum_{i=1}^{N} \left(\sigma_{N-i} - \sigma_{N+1-i}\right) \left(f^{i} - f^{N+1} - \left\langle g^{N+1}, x^{i} - x^{N+1} \right\rangle\right) + \frac{B\sigma_{N}^{2}}{2Rh} \left(R^{2} - \left\|x^{1}\right\|^{2}\right)$$

$$+ \sigma_{N} \sum_{i=1}^{N} \sigma_{N-i} \left( -f^{i} - \left\langle \frac{B}{Rh} (x^{i} - x^{i+1}), -x^{i} \right\rangle \right) + \frac{Rh}{2B} \left( B^{2} - \left\| g^{N+1} \right\|^{2} \right) \\ + \sigma_{N} \left( -f^{N+1} + \left\langle g^{N+1}, x^{i} \right\rangle \right) \\ = \frac{Rh}{2B} \left\| z^{N+1} - B z^{N+1} + B \sum_{i=1}^{N} \left( z^{N-1} - z^{N-1} \right) z^{N-1} \right\|^{2} \leq 0$$

$$= \frac{-Rh}{2B} \left\| g^{N+1} - \frac{B}{Rh} x^{N+1} + \frac{B}{Rh} \sum_{i=1} \left( \sigma_{N-i} - \sigma_{N+1-i} \right) x^i \right\| \leq 0.$$
<sup>19</sup>

- After staring at the PEP proof, we noticed similarities between inequality multipliers
- ► Grouping similar terms, we obtain Jensen-like inequalities (insight: applying Jensen ↔ some sum of interpolation inequalities)
- Simplifying further we obtain a classic-style proof, that is no longer looking computer generated
- ▶ We encapsulate the main part of the proof in a key Lemma
- Key Lemma fully reverse-engineered from PEP but can be easily checked by hand

#### Key Lemma for subgradient methods

Lemma ([Zamani,G 2023])

Consider the subgradient method with fixed step sizes  $\{h_k\}$ 

 $x_{k+1} = x_k - h_k g_k$  for some  $g_k \in \partial f(x_k)$  for  $k = 0, 1, \dots, N-1$ 

Choose  $h_N > 0$  and introduce N + 2 weights  $v_k$  that satisfy

$$1 = v_{-1} \leq v_0 \leq v_1 \leq \cdots \leq v_{N-1} \leq v_N$$

Then iterates after N iterations of the subgradient method satisfy

$$\sum_{k=0}^{N} \left( h_{k} v_{k}^{2} - (v_{k} - v_{k-1}) \sum_{i=k}^{N} h_{i} v_{i} \right) \left( f(x^{k}) - f(x^{*}) \right)$$

$$\leq \frac{1}{2} \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}^{2} + \frac{1}{2} \sum_{k=0}^{N} h_{k}^{2} v_{k}^{2} \underbrace{\left\| g^{k} \right\|}_{B}^{2}$$

Key Lemma inequality:

$$\sum_{k=0}^{N} \left( h_{k} v_{k}^{2} - (v_{k} - v_{k-1}) \sum_{i=k}^{N} h_{i} v_{i} \right) \left( f(x^{k}) - f(x^{*}) \right) \leq \frac{1}{2} v_{0}^{2} R^{2} + \frac{1}{2} B^{2} \sum_{k=0}^{N} h_{k}^{2} v_{k}^{2}$$

for any weights  $1 = v_{-1} \le v_0 \le v_1 \le \cdots \le v_{N-1} \le v_N$ 

- constant  $v_k = 1$  recovers usual (average) rate
- ▶ but a suitable choice of {v<sub>k</sub>} allows us to modify coefficients in front of f(x<sup>k</sup>) - f(x<sup>\*</sup>)
- ► in particular one can cancel all coefficients except last one in front of f(x<sup>N</sup>) f(x\*)

## Idea of the proof of the key Lemma

Inequality to prove:

$$\sum_{k=0}^{N} \left( h_{k} v_{k}^{2} - (v_{k} - v_{k-1}) \sum_{i=k}^{N} h_{i} v_{i} \right) \left( f(x^{k}) - f(x^{*}) \right) \leq \frac{1}{2} v_{0}^{2} R^{2} + \frac{1}{2} B^{2} \sum_{k=0}^{N} h_{k}^{2} v_{k}^{2}$$

Proof uses a *generalization* of the standard telescoping proof

1. From weights  $v_k$  define *auxiliary sequence*  $z^k$  recursively

$$z^0 = x^*$$
 and  $z^k = \left(1 - \frac{v_{k-1}}{v_k}\right)x^k + \left(\frac{v_{k-1}}{v_k}\right)z^{k-1}$ 

This implies

$$z^{k} = \left(\frac{v_{0}}{v_{k}}\right)x^{*} + \sum_{i=1}^{k} \left(\frac{v_{i} - v_{i-1}}{v_{k}}\right)x^{i}$$

(note  $z^k$  is a convex combination of  $x^*$  and iterates  $x^i$ )

# Idea of the proof of the key Lemma (cont.)

2. Subgradient inequality between  $x^k$  and  $z^k$  (instead of  $x^*$ ) gives  $h_k v_k^2 (f(x^k) - f(z^k)) \leq \frac{1}{2} v_{k-1}^2 \|x^k - z^k\|^2 - \frac{1}{2} v_k^2 \|x^{k+1} - z^{k+1}\|^2 + \frac{1}{2} B^2 h_k^2 v_k^2$ 

3. Telescoping (summing from k = 0 to k = N) gives that

$$\sum_{k=0}^{N} h_{k} v_{k}^{2} (f(x^{k}) - f(z^{k}))$$

$$\leq \frac{1}{2} v_{-1}^{2} \|x^{0} - z^{0}\|^{2} - \frac{1}{2} v_{N}^{2} \|x^{N+1} - z^{N+1}\|^{2} + \frac{1}{2} B^{2} \sum_{k=0}^{N} h_{k}^{2} v_{k}^{2}$$

implying

Λ /

$$\sum_{k=0}^{N} h_k v_k^2 (f(x^k) - f(z^k)) \le \frac{1}{2} \|x^0 - x^*\|^2 + \frac{1}{2} B^2 \sum_{k=0}^{N} h_k^2 v_k^2$$

4. Finally we need to find a lower bound on  $f(x^k) - f(z^k)$  terms:

$$\mathbf{z}^k = (rac{v_0}{v_k})\hat{x} + \sum_{i=1}^k (rac{v_i - v_{i-1}}{v_k})x^i$$

implies, by Jensen's inequality

$$f(z^k) \le (\frac{v_0}{v_k})f(\hat{x}) + \sum_{i=1}^k (\frac{v_i - v_{i-1}}{v_k})f(x^i)$$

hence

$$h_k v_k^2 (f(z^k) - f(x^*)) \ge h_k v_k \sum_{i=1}^k (v_i - v_{i-1}) (f(x^i) - f(x^*))$$

which combined with inequality from the previous step 3. gives

$$\sum_{k=0}^{N} \left( h_k v_k^2 - (v_k - v_{k-1}) \sum_{i=k}^{N} h_i v_i \right) \left( f(x^k) - f(x^*) \right)$$
  
$$\leq \sum_{k=0}^{N} h_k v_k^2 \left( f(x^k) - f(z^k) \right) \leq \frac{1}{2} \|x^0 - x^*\|^2 + \frac{1}{2} B^2 \sum_{k=0}^{N} h_k^2 v_k^2$$

So we have proved

#### Lemma

Iterates of the subgradient methods satisfy

$$\sum_{k=0}^{N} \left( h_k v_k^2 - (v_k - v_{k-1}) \sum_{i=k}^{N} h_i v_i \right) \left( f(x^k) - f(x^*) \right)$$
$$\leq \frac{1}{2} R^2 + \frac{1}{2} B^2 \sum_{k=0}^{N} h_k^2 v_k^2$$

*Proof* of last-iterate convergence rate:

Choose weights  $v_k$  that cancel all coefficients of  $f(x^k)$  except  $f(x^N)$ , which are

V

$$_{k}=\frac{1}{s_{N+1-k}}$$
 26

## All PEP rates are *exact* by design

(cannot be improved, even by a multiplicative/additive constant)

Follows from PEP solution, but can be made constructive by building an *explicit worst-case function* 

- Function of the type  $f(x) = [\max_k \{g_k^T x\}]_+$
- Recursive definition, coefficients  $g_k$  not straightforward
- Sugbradients for all iterates are  $g_k$ , have maximum norm B
- Subgradient inequality is satisfied between all pairs of iterates
- Matches exactly the announced convergence rate for the last iterate

# **Extensions**

## Last-iterate optimal subgradient method

Define the following new linearly decreasing stepsize schedule

$$x_{k+1} = x_k - \frac{R}{B} \frac{(N+1-k)}{(N+1)^{3/2}} g_k$$

Leads the optimal rate for the last iterate [Zamani,G 2023]

$$f(x_N) - f(x_*) \le \frac{BR}{\sqrt{N+1}}$$

- ► Improves <sup>15BD</sup>/<sub>√N+1</sub> [Jain,Nagaraj,Netrapalli 2021] for diameter D
   ► Same proof technique, key lemma with optimized weights v<sub>k</sub>
- Schedule dependence on N is forced for optimal method (already impossible to find fixed stepsizes h<sub>1</sub> and h<sub>2</sub> that are optimal for both N = 1 and N = 2)

#### Open question:

Existence of a last-iterate optimal method with stepsizes independent from N and with momentum terms?

# Subgradient method with normalized step sizes

Stepsizes so far feature a  $\frac{R}{B}$  factor, require knowledge of R and B

• constant stepsizes 
$$h_k = \frac{R}{B}h$$
 for some  $h$ 

• optimal stepsizes  $h_k = \frac{R}{B} \frac{(N+1-k)}{(N+1)^{3/2}}$ 

Need for *B* can be removed using normalized step sizes  $\{t_k\}$ 

$$x_{k+1} = x_k - t_k \frac{g_k}{\|g_k\|}$$
 for some  $g_k \in \partial f(x_k)$ 

- All previous results are also valid with exactly the same rates if we assume  $t_k = h_k B$
- constant stepsizes  $t_k = Rh$  for some h
- optimal stepsizes  $t_k = R \frac{(N+1-k)}{(N+1)^{3/2}}$
- Proof using key Lemma with adapted weights
- Removing dependence on R seems harder ( $\rightarrow$  parameter-free)

#### Solve convex constrained optimization

 $\min_{x\in X} f(x)$ 

with the projected subgradient method with fixed step sizes  $\{h_k\}$ 

 $x_{k+1} = \mathbb{P}[x_k - h_k g_k]$  for some  $g_k \in \partial f(x_k)$ 

( $\mathbb{P}$  is orthogonal projection on convex set X)

- All results are also valid, with exactly the same rates (both constant and optimal stepsizes, also normalized)
- Straightforward adaptation of the key Lemma using non-expansiveness of the projection operator

# Conclusions

#### Take-home messages:

- Performance estimation applied to subgradient methods
- Exact convergence rates can be obtained for the last iterate: suboptimal by a factor  $O(\sqrt{\log(N)})$
- New last-iterate optimal method can be designed with linearly decreasing step sizes
- Extensions to constrained case, to normalized steps
- Inspiration for results provided by performance estimation but ultimately all proofs converted to classical style using a new key lemma

For all your performance estimation needs:

https://github.com/PerformanceEstimation

Thank you Yurii!