

Making SGD Parameter-Free

Presented by Qilin Ye

May 9, 2023

§0 Stochastic Gradient Descent

... the same old gradient descent:

$$x_{t+1} := x_t - \eta \nabla F(x_t)$$

where F is convex & differentiable.

Non-differentiable? Use unbiased **subgradients**: $x_{t+1} := x_t - \eta g_t$.¹

¹A subgradient of f satisfies $f(z) \geq f(x) + g^T(z - x)$ for all z .

§0 Stochastic Gradient Descent

... the same old gradient descent:

$$x_{t+1} := x_t - \eta \nabla F(x_t)$$

where F is convex & differentiable.

Non-differentiable? Use unbiased **subgradients**: $x_{t+1} := x_t - \eta g_t$.¹

¹A subgradient of f satisfies $f(z) \geq f(x) + g^T(z - x)$ for all z .

§0 Stochastic Gradient Descent... Problems?

Apparently, choosing the correct learning rate is not a trivial job.

- (1) Too large? Possible oscillation. Too small? Slow!
- (2) Distance between starting point and optimum matters.
- (3) The rate of convergence is affected by scaling.
- (4) ...

§0 Stochastic Gradient Descent... Problems?

Apparently, choosing the correct learning rate is not a trivial job.

- (1) Too large? Possible oscillation. Too small? Slow!
- (2) Distance between starting point and optimum matters.
- (3) The rate of convergence is affected by scaling.
- (4) ...

§0 Stochastic Gradient Descent... Problems?

Apparently, choosing the correct learning rate is not a trivial job.

- (1) Too large? Possible oscillation. Too small? Slow!
- (2) Distance between starting point and optimum matters.
- (3) The rate of convergence is affected by scaling.
- (4) ...

§0 Parameter-Free Optimizations and Regrets

We aim to design **parameter-free** algorithms that “automatically” tune the learning rate.

And we aim to obtain “good” **regret** guarantees.

§0 Parameter-Free Optimizations and Regrets

We aim to design **parameter-free** algorithms that “automatically” tune the learning rate.

And we aim to obtain “good” **regret** guarantees.

§0 Notations and Problem Setup

- (1) Let $\mathcal{X} \subset \mathbb{R}^d$ be convex closed and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be convex.
- (2) Let x^* a minimum of f , assuming existence.
- (3) Let \mathcal{O} be an oracle that is a subgradient of f in expectation:
 $\mathbb{E}[\mathcal{O}(x) \mid x] \in \partial f(x)$.
- (4) Denote the iterates by $x_0, x_1(\eta), x_2(\eta), \dots$ and (sub)gradients $g_0, g_1(\eta), g_2(\eta), \dots$. Define $\bar{x}(\eta) := T^{-1} \sum_{i < T} x_i(\eta)$.
- (5) Distance to optimum and running maximum distance:

$$d_t(\eta) := \|x_t(\eta) - x^*\| \quad \bar{d}_t(\eta) := \max_{i \leq t} d_i(\eta).$$

- (6) Distance to x_0 and running max distance: $r_t(\eta), \bar{r}_t(\eta)$.
- (7) Oracle error & running squared norms of oracles:

$$\Delta_i := g_i - \nabla f(x_i(\eta)) \quad G_t(\eta) := \sum_{i < t} \|g_i(\eta)\|^2.$$

§0 Notations and Problem Setup

- (1) Let $\mathcal{X} \subset \mathbb{R}^d$ be convex closed and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be convex.
- (2) Let x^* a minimum of f , assuming existence.
- (3) Let \mathcal{O} be an oracle that is a subgradient of f in expectation:
 $\mathbb{E}[\mathcal{O}(x) \mid x] \in \partial f(x)$.
- (4) Denote the iterates by $x_0, x_1(\eta), x_2(\eta), \dots$ and (sub)gradients $g_0, g_1(\eta), g_2(\eta), \dots$. Define $\bar{x}(\eta) := T^{-1} \sum_{i < T} x_i(\eta)$.
- (5) Distance to optimum and running maximum distance:

$$d_t(\eta) := \|x_t(\eta) - x^*\| \quad \bar{d}_t(\eta) := \max_{i \leq t} d_i(\eta).$$

- (6) Distance to x_0 and running max distance: $r_t(\eta), \bar{r}_t(\eta)$.
- (7) Oracle error & running squared norms of oracles:

$$\Delta_i := g_i - \nabla f(x_i(\eta)) \quad G_t(\eta) := \sum_{i < t} \|g_i(\eta)\|^2.$$

§0 Notations and Problem Setup

- (1) Let $\mathcal{X} \subset \mathbb{R}^d$ be convex closed and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be convex.
- (2) Let x^* a minimum of f , assuming existence.
- (3) Let \mathcal{O} be an oracle that is a subgradient of f in expectation:
 $\mathbb{E}[\mathcal{O}(x) \mid x] \in \partial f(x)$.
- (4) Denote the iterates by $x_0, x_1(\eta), x_2(\eta), \dots$ and (sub)gradients $g_0, g_1(\eta), g_2(\eta), \dots$. Define $\bar{x}(\eta) := T^{-1} \sum_{i < T} x_i(\eta)$.
- (5) Distance to optimum and running maximum distance:

$$d_t(\eta) := \|x_t(\eta) - x^*\| \quad \bar{d}_t(\eta) := \max_{i \leq t} d_i(\eta).$$

- (6) Distance to x_0 and running max distance: $r_t(\eta), \bar{r}_t(\eta)$.
- (7) Oracle error & running squared norms of oracles:

$$\Delta_i := g_i - \nabla f(x_i(\eta)) \quad G_t(\eta) := \sum_{i < t} \|g_i(\eta)\|^2.$$

§0 Notations and Problem Setup

- (1) Let $\mathcal{X} \subset \mathbb{R}^d$ be convex closed and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be convex.
- (2) Let x^* a minimum of f , assuming existence.
- (3) Let \mathcal{O} be an oracle that is a subgradient of f in expectation:
 $\mathbb{E}[\mathcal{O}(x) \mid x] \in \partial f(x)$.
- (4) Denote the iterates by $x_0, x_1(\eta), x_2(\eta), \dots$ and (sub)gradients $g_0, g_1(\eta), g_2(\eta), \dots$. Define $\bar{x}(\eta) := T^{-1} \sum_{i < T} x_i(\eta)$.
- (5) Distance to optimum and running maximum distance:

$$d_t(\eta) := \|x_t(\eta) - x^*\| \quad \bar{d}_t(\eta) := \max_{i \leq t} d_i(\eta).$$

- (6) Distance to x_0 and running max distance: $r_t(\eta), \bar{r}_t(\eta)$.
- (7) Oracle error & running squared norms of oracles:

$$\Delta_i := g_i - \nabla f(x_i(\eta)) \quad G_t(\eta) := \sum_{i < t} \|g_i(\eta)\|^2.$$

§0 Notations and Problem Setup

- (1) Let $\mathcal{X} \subset \mathbb{R}^d$ be convex closed and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be convex.
- (2) Let x^* a minimum of f , assuming existence.
- (3) Let \mathcal{O} be an oracle that is a subgradient of f in expectation:
 $\mathbb{E}[\mathcal{O}(x) \mid x] \in \partial f(x)$.
- (4) Denote the iterates by $x_0, x_1(\eta), x_2(\eta), \dots$ and (sub)gradients $g_0, g_1(\eta), g_2(\eta), \dots$. Define $\bar{x}(\eta) := T^{-1} \sum_{i < T} x_i(\eta)$.
- (5) Distance to optimum and running maximum distance:

$$d_t(\eta) := \|x_t(\eta) - x^*\| \quad \bar{d}_t(\eta) := \max_{i \leq t} d_i(\eta).$$

- (6) Distance to x_0 and running max distance: $r_t(\eta), \bar{r}_t(\eta)$.
- (7) Oracle error & running squared norms of oracles:

$$\Delta_i := g_i - \nabla f(x_i(\eta)) \quad G_t(\eta) := \sum_{i < t} \|g_i(\eta)\|^2.$$

Fact 1

If all $\|g_i\|$'s are uniformly bounded by $L > 0$, then setting η to be the fixed point of

$$\eta \mapsto \frac{\|x_0 - x^*\|}{(\sum_{i < T} \|g_i(\eta)\|^2)^{1/2}} = \frac{d_0}{\sqrt{G_T(\eta)}}$$

satisfies the optimal error bound for the average iterate after T iterations:

$$f(\bar{x}) - f(x^*) \leq \frac{d_0 \sqrt{G_T(\eta)}}{T} = O(d_0 L T^{-1/2}).$$

Fact 2: SoTA w/out Knowing $d_0 = \|x_0 - x^*\|$ a priori

... gains an additional logarithmic factor:

$$O\left(d_0 \sqrt{\log(1 + T d_0^2 \epsilon^{-2})} / T + \epsilon / T\right).$$

§0 What Did This Paper Do?

- (1) For any prescribed $\epsilon > 0$ and $\delta \in (0, 1)$, this paper provides a $1 - \delta$ probability optimality gap with an additional log factor:

$$O\left((d_0 T^{-1/2} + \epsilon T^{-1}) \cdot \log^2(\delta^{-1} \log(d_0 T \epsilon^{-1}))\right)$$

- (2) Strong localization guarantee: the average iterate (as well as other intermediate outputs) \bar{x} satisfies $\|\bar{x} - x^*\| = O(\|x_0 - x^*\|)$.
- (3) Good adaptivity to gradient norms and other scenarios.

§0 What Did This Paper Do?

- (1) For any prescribed $\epsilon > 0$ and $\delta \in (0, 1)$, this paper provides a $1 - \delta$ probability optimality gap with an additional log factor:

$$O\left((d_0 T^{-1/2} + \epsilon T^{-1}) \cdot \log^2(\delta^{-1} \log(d_0 T \epsilon^{-1}))\right)$$

- (2) Strong localization guarantee: the average iterate (as well as other intermediate outputs) \bar{x} satisfies $\|\bar{x} - x^*\| = O(\|x_0 - x^*\|)$.
- (3) Good adaptivity to gradient norms and other scenarios.

§0 What Did This Paper Do?

- (1) For any prescribed $\epsilon > 0$ and $\delta \in (0, 1)$, this paper provides a $1 - \delta$ probability optimality gap with an additional log factor:

$$O\left((d_0 T^{-1/2} + \epsilon T^{-1}) \cdot \log^2(\delta^{-1} \log(d_0 T \epsilon^{-1}))\right)$$

- (2) Strong localization guarantee: the average iterate (as well as other intermediate outputs) \bar{x} satisfies $\|\bar{x} - x^*\| = O(\|x_0 - x^*\|)$.
- (3) Good adaptivity to gradient norms and other scenarios.

This page is intentionally left blank.

§1 High-Level Idea: Using Proxy for d_0

In SGD, the output iterates $x_t(\eta)$ should ideally converge to x^* (recall the **optimal bound** with knowledge of x^*)

$$\Rightarrow \frac{r_t(\eta)}{\sqrt{G_T(\eta)}} \text{ converges to } \frac{d_0}{\sqrt{G_T(\eta)}}.$$

Instead of computing the uncomputable fixed point, we resort to approximating the fixed point of

$$\eta \mapsto \frac{\bar{r}_T(\eta)}{\sqrt{\alpha G_T(\eta) + \beta}}. \quad (\text{FP1})$$

(Why \bar{r}_T instead of r_T ?)

§1 Proposition 1

Assuming we have found the η satisfying (FP1), and with probability 1 our oracle $\mathcal{O}(x) = \nabla f(x)$ (i.e. *true gradient*):

Proposition 1

If $\alpha > 1, \beta = 0$, then the average iterate $\bar{x} := T^{-1} \sum_{i < T} x_i(\eta)$ satisfies

$$\|\bar{x} - x^*\| \leq \frac{2\alpha}{\alpha - 1} \|x_0 - x^*\| = \frac{2\alpha}{\alpha - 1} d_0$$

and

$$f(\bar{x}) - f(x^*) \leq \frac{\alpha^{3/2}}{\alpha - 1} \cdot \frac{d_0 \sqrt{G_T(\eta)}}{T} \sim \frac{d_0 \sqrt{G_T(\eta)}}{T}.$$

(This is the optimal regret bound!)

§1 Limitations of Proposition 1?

The mapping $\varphi : \eta \mapsto \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$ may not be continuous...
 \Rightarrow a fixed point may not exist!

Workaround: find a small interval where $\eta \mapsto \varphi(\eta) - \eta$ changes sign...
 \Rightarrow Bisection!

§1 Limitations of Proposition 1?

The mapping $\varphi : \eta \mapsto \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$ may not be continuous...
 \Rightarrow a fixed point may not exist!

Workaround: find a small interval where $\eta \mapsto \varphi(\eta) - \eta$ changes sign...
 \Rightarrow Bisection!

§1 Limitations of Proposition 1?

The mapping $\varphi : \eta \mapsto \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$ may not be continuous...
 \Rightarrow a fixed point may not exist!

Workaround: find a small interval where $\eta \mapsto \varphi(\eta) - \eta$ changes sign...
 \Rightarrow Bisection!

§1 Limitations of Proposition 1 (continued)?

Also, what if the exact gradient assumption is removed?

This page is intentionally left blank.

§2 The Algorithm (exact gradient version)

Algorithm 1: Parameter-free SGD step size tuning (exact gradient version)²

- 1 **Inputs:** initial learning rate $\eta_\epsilon > 0$, total gradient budget $B \in \mathbb{N}$, damping parameters $\{\alpha^{(k)}, \beta^{(k)}\}$.
 - 2 **for** $k = 2, 4, 8, 16, \dots$ **do**
 - 3 **if** $k > B/4$ **then return** x_0 (??) \triangleright edge case, bad B ; make it larger!
 - 4 $T_k \leftarrow \lfloor B/(2k) \rfloor$ \triangleright dynamically adjust SGD complexity based on k
 - 5 $\eta_0 \leftarrow \text{RootFindingBisection}(\eta_\epsilon, 2^{2^k} \eta_\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$
 - 6 **if** Bisection says k is OK (??) **then return** $T_k^{-1} \sum_{i < T_k} x_i(\eta_0)$
-

²(??) to be explained later.

§2 The Algorithm (exact gradient version)

Algorithm 2: Parameter-free SGD step size tuning (exact gradient version)²

- 1 **Inputs:** initial learning rate $\eta_\epsilon > 0$, total gradient budget $B \in \mathbb{N}$, damping parameters $\{\alpha^{(k)}, \beta^{(k)}\}$.
 - 2 **for** $k = 2, 4, 8, 16, \dots$ **do**
 - 3 **if** $k > B/4$ **then return** x_0 (??) ▷ edge case, bad B ; make it larger!
 - 4 $T_k \leftarrow \lfloor B/(2k) \rfloor$ ▷ dynamically adjust SGD complexity based on k
 - 5 $\eta_0 \leftarrow \text{RootFindingBisection}(\eta_\epsilon, 2^{2^k} \eta_\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$
 - 6 **if** Bisection says k is OK (??) **then return** $T_k^{-1} \sum_{i < T_k} x_i(\eta_0)$
-

²(??) to be explained later.

§2 The Algorithm (exact gradient version)

Algorithm 3: Parameter-free SGD step size tuning (exact gradient version)²

- 1 **Inputs:** initial learning rate $\eta_\epsilon > 0$, total gradient budget $B \in \mathbb{N}$, damping parameters $\{\alpha^{(k)}, \beta^{(k)}\}$.
 - 2 **for** $k = 2, 4, 8, 16, \dots$ **do**
 - 3 **if** $k > B/4$ **then return** x_0 (??) ▷ edge case, bad B ; make it larger!
 - 4 $T_k \leftarrow \lfloor B/(2k) \rfloor$ ▷ dynamically adjust SGD complexity based on k
 - 5 $\eta_0 \leftarrow \text{RootFindingBisection}(\eta_\epsilon, 2^{2^k} \eta_\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$
 - 6 **if** Bisection says k is OK (??) **then return** $T_k^{-1} \sum_{i < T_k} x_i(\eta_0)$
-

²(??) to be explained later.

§2 The Algorithm (exact gradient version)

Algorithm 4: Parameter-free SGD step size tuning (exact gradient version)²

- 1 **Inputs:** initial learning rate $\eta_\epsilon > 0$, total gradient budget $B \in \mathbb{N}$, damping parameters $\{\alpha^{(k)}, \beta^{(k)}\}$.
 - 2 **for** $k = 2, 4, 8, 16, \dots$ **do**
 - 3 **if** $k > B/4$ **then return** x_0 (??) ▷ edge case, bad B ; make it larger!
 - 4 $T_k \leftarrow \lfloor B/(2k) \rfloor$ ▷ dynamically adjust SGD complexity based on k
 - 5 $\eta_0 \leftarrow \text{RootFindingBisection}(\eta_\epsilon, 2^{2^k} \eta_\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$
 - 6 **if** Bisection says k is OK (??) **then return** $T_k^{-1} \sum_{i < T_k} x_i(\eta_0)$
-

²(??) to be explained later.

§2 Two Lemmas & Termination Guarantees

Lemma 1

With appropriate parameters, under exact gradient setting,

$$\eta \leq \varphi(\eta) \Rightarrow \bar{d}_T(\eta) \leq \frac{\alpha + 1}{\alpha - 1} \cdot d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{2\alpha}{\alpha - 1} d_0.$$

Proof. First notice that $d_{i+1}^2 = \|x_i - \eta g_i - x^*\|^2 = d_i^2 - 2\eta \langle g_i, x_i - x^* \rangle + \eta^2 \|g_i\|^2$. Then, by (sub)gradient and convexity of f , we also have $\langle g_i, x_i - x^* \rangle \geq f(x_i) - f(x^*) \geq 0$. Summation over all $i < t$ gives $d_t^2 \leq d_0^2 + \eta^2 G_t$. The remainder of the proof are pure algebraic manipulations and are likely not ideal for presentation.

§2 Two Lemmas & Termination Guarantees

Lemma 1

With appropriate parameters, under exact gradient setting,

$$\eta \leq \varphi(\eta) \Rightarrow \bar{d}_T(\eta) \leq \frac{\alpha + 1}{\alpha - 1} \cdot d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{2\alpha}{\alpha - 1} d_0.$$

Proof. First notice that $d_{i+1}^2 = \|x_i - \eta g_i - x^*\|^2 = d_i^2 - 2\eta \langle g_i, x_i - x^* \rangle + \eta^2 \|g_i\|^2$. Then, by (sub)gradient and convexity of f , we also have $\langle g_i, x_i - x^* \rangle \geq f(x_i) - f(x^*) \geq 0$. Summation over all $i < t$ gives $d_t^2 \leq d_0^2 + \eta^2 G_t$. The remainder of the proof are pure algebraic manipulations and are likely not ideal for presentation.

Lemma 2

With appropriate parameters, under exact gradient setting, if the following holds, then $\eta > \varphi(\eta)$:

$$\eta > \eta_{\max} := \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}}.$$

Proof. If $\eta > \eta_{\max}$ but $\eta \leq \varphi(\eta)$, using $\|g_0\|^2 \leq \sum \|g_i\|^2 = G_T(\eta)$ we obtain

$$\frac{\bar{r}_T(\eta)}{\sqrt{\alpha\|g_0\|^2 + \beta}} \geq \frac{\bar{r}_T(\eta)}{\sqrt{\alpha G_T(\eta) + \beta}} = \varphi(\eta) \geq \eta > \eta_{\max} = \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}},$$

contradicting the previous lemma.

§2 Termination Guarantees

Upshot: if k is such that $2^{2^k} \eta_\epsilon > \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0}{\sqrt{\alpha \|g_0\|^2 + \beta}}$, then

$$\eta \mapsto \varphi(\eta) - \eta$$

changes sign on $[\eta_\epsilon, 2^{2^k} \eta_\epsilon]$. Time for bisection!

§2 The Algorithm – RootFindingBisection

1 **Function** RootFindingBisection($\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta$):

```
2   define  $\varphi$  by  $\varphi(\eta) = \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta)} + \beta$            ▷ bisection target
3   if  $\eta_{\text{high}} \leq \varphi(\eta_{\text{high}})$  then return  $\infty$            ▷  $\eta_{\text{high}}$  too low, need to increase
4   if  $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$  then return  $\eta_{\text{low}}$            ▷  $\eta_{\text{low}}$  is sufficient (if small)
5   while  $\eta_{\text{high}} > 2\eta_{\text{low}}$  do
6     ▷ loop invariant:  $\eta_{\text{low}} < \eta_{\text{high}}, \eta_{\text{low}} \leq \varphi(\eta_{\text{low}}), \eta_{\text{high}} > \varphi(\eta_{\text{high}})$ 
7      $\eta_{\text{mid}} \leftarrow \sqrt{\eta_{\text{low}}\eta_{\text{high}}}$ 
8     if  $\eta_{\text{mid}} \leq \varphi(\eta_{\text{mid}})$  then  $\eta_{\text{low}} \leftarrow \eta_{\text{mid}}$  else  $\eta_{\text{high}} \leftarrow \eta_{\text{mid}}$ 
9   if  $\bar{r}_T(\eta_{\text{high}}) \leq \bar{r}_T(\eta_{\text{low}}) \cdot \varphi(\eta_{\text{high}}) / \eta_{\text{high}}$  then return  $\eta_{\text{high}}$ 
   else return  $\eta_{\text{low}}$ 
```

§2 The Algorithm – RootFindingBisection

```
1 Function RootFindingBisection( $\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta$ ):
2   define  $\varphi$  by  $\varphi(\eta) = \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$             $\triangleright$  bisection target
3   if  $\eta_{\text{high}} \leq \varphi(\eta_{\text{high}})$  then return  $\infty$             $\triangleright \eta_{\text{high}}$  too low, need to increase
4   if  $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$  then return  $\eta_{\text{low}}$             $\triangleright \eta_{\text{low}}$  is sufficient (if small)
5   while  $\eta_{\text{high}} > 2\eta_{\text{low}}$  do
6      $\triangleright$  loop invariant:  $\eta_{\text{low}} < \eta_{\text{high}}, \eta_{\text{low}} \leq \varphi(\eta_{\text{low}}), \eta_{\text{high}} > \varphi(\eta_{\text{high}})$ 
7      $\eta_{\text{mid}} \leftarrow \sqrt{\eta_{\text{low}}\eta_{\text{high}}}$ 
8     if  $\eta_{\text{mid}} \leq \varphi(\eta_{\text{mid}})$  then  $\eta_{\text{low}} \leftarrow \eta_{\text{mid}}$  else  $\eta_{\text{high}} \leftarrow \eta_{\text{mid}}$ 
9   if  $\bar{r}_T(\eta_{\text{high}}) \leq \bar{r}_T(\eta_{\text{low}}) \cdot \varphi(\eta_{\text{high}}) / \eta_{\text{high}}$  then return  $\eta_{\text{high}}$ 
   else return  $\eta_{\text{low}}$ 
```

§2 The Algorithm – RootFindingBisection

```
1 Function RootFindingBisection( $\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta$ ):
2   define  $\varphi$  by  $\varphi(\eta) = \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta)} + \beta$             $\triangleright$  bisection target
3   if  $\eta_{\text{high}} \leq \varphi(\eta_{\text{high}})$  then return  $\infty$             $\triangleright \eta_{\text{high}}$  too low, need to increase
4   if  $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$  then return  $\eta_{\text{low}}$             $\triangleright \eta_{\text{low}}$  is sufficient (if small)
5   while  $\eta_{\text{high}} > 2\eta_{\text{low}}$  do
6      $\eta_{\text{mid}} \leftarrow \sqrt{\eta_{\text{low}}\eta_{\text{high}}}$ 
7     if  $\eta_{\text{mid}} \leq \varphi(\eta_{\text{mid}})$  then  $\eta_{\text{low}} \leftarrow \eta_{\text{mid}}$  else  $\eta_{\text{high}} \leftarrow \eta_{\text{mid}}$ 
8   if  $\bar{r}_T(\eta_{\text{high}}) \leq \bar{r}_T(\eta_{\text{low}}) \cdot \varphi(\eta_{\text{high}}) / \eta_{\text{high}}$  then return  $\eta_{\text{high}}$ 
9   else return  $\eta_{\text{low}}$ 
```

§2 The Algorithm – RootFindingBisection

```
1 Function RootFindingBisection( $\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta$ ):
2   define  $\varphi$  by  $\varphi(\eta) = \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$             $\triangleright$  bisection target
3   if  $\eta_{\text{high}} \leq \varphi(\eta_{\text{high}})$  then return  $\infty$             $\triangleright \eta_{\text{high}}$  too low, need to increase
4   if  $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$  then return  $\eta_{\text{low}}$             $\triangleright \eta_{\text{low}}$  is sufficient (if small)
5   while  $\eta_{\text{high}} > 2\eta_{\text{low}}$  do
6      $\triangleright$  loop invariant:  $\eta_{\text{low}} < \eta_{\text{high}}, \eta_{\text{low}} \leq \varphi(\eta_{\text{low}}), \eta_{\text{high}} > \varphi(\eta_{\text{high}})$ 
7      $\eta_{\text{mid}} \leftarrow \sqrt{\eta_{\text{low}}\eta_{\text{high}}}$ 
8     if  $\eta_{\text{mid}} \leq \varphi(\eta_{\text{mid}})$  then  $\eta_{\text{low}} \leftarrow \eta_{\text{mid}}$  else  $\eta_{\text{high}} \leftarrow \eta_{\text{mid}}$ 
9   if  $\bar{r}_T(\eta_{\text{high}}) \leq \bar{r}_T(\eta_{\text{low}}) \cdot \varphi(\eta_{\text{high}}) / \eta_{\text{high}}$  then return  $\eta_{\text{high}}$ 
   else return  $\eta_{\text{low}}$ 
```

§2 The Algorithm – RootFindingBisection

```
1 Function RootFindingBisection( $\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta$ ):
2   define  $\varphi$  by  $\varphi(\eta) = \bar{r}_T(\eta) / \sqrt{\alpha G_T(\eta) + \beta}$             $\triangleright$  bisection target
3   if  $\eta_{\text{high}} \leq \varphi(\eta_{\text{high}})$  then return  $\infty$             $\triangleright \eta_{\text{high}}$  too low, need to increase
4   if  $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$  then return  $\eta_{\text{low}}$             $\triangleright \eta_{\text{low}}$  is sufficient (if small)
5   while  $\eta_{\text{high}} > 2\eta_{\text{low}}$  do
6      $\triangleright$  loop invariant:  $\eta_{\text{low}} < \eta_{\text{high}}, \eta_{\text{low}} \leq \varphi(\eta_{\text{low}}), \eta_{\text{high}} > \varphi(\eta_{\text{high}})$ 
7      $\eta_{\text{mid}} \leftarrow \sqrt{\eta_{\text{low}}\eta_{\text{high}}}$ 
8     if  $\eta_{\text{mid}} \leq \varphi(\eta_{\text{mid}})$  then  $\eta_{\text{low}} \leftarrow \eta_{\text{mid}}$  else  $\eta_{\text{high}} \leftarrow \eta_{\text{mid}}$ 
9   if  $\bar{r}_T(\eta_{\text{high}}) \leq \bar{r}_T(\eta_{\text{low}}) \cdot \varphi(\eta_{\text{high}}) / \eta_{\text{high}}$  then return  $\eta_{\text{high}}$ 
10  else return  $\eta_{\text{low}}$ 
```

§2 Properties of RootFindingBisection

- (1) Each iteration halves $\log(\eta_{\text{high}}/\eta_{\text{low}})$
 \Rightarrow number of iterations is $\log \log(\eta_{\text{high}}/\eta_{\text{low}})$.
Consequently, by **Lemma 2**, when our algorithm terminates, $k \leq 2 \log \log^+(\eta_{\text{max}}/\eta_{\epsilon})$.
- (2) Approximates the (possibly non-existent) root $\eta = \varphi(\eta)$ up to a factor of 2, even when root is non-existent. If the returned interval is $[\eta_{\text{low}}^*, \eta_{\text{high}}^*]$ then

$$\frac{\bar{r}_T(\eta_0)}{2\sqrt{\alpha G_T(\eta_{\text{high}}^*) + \beta}} \leq \eta_0 \leq \frac{\bar{r}_T(\eta_{\text{low}}^*)}{\sqrt{\alpha G_T(\eta_0) + \beta}}$$

§2 Properties of RootFindingBisection

- (1) Each iteration halves $\log(\eta_{\text{high}}/\eta_{\text{low}})$
 \Rightarrow number of iterations is $\log \log(\eta_{\text{high}}/\eta_{\text{low}})$.
Consequently, by **Lemma 2**, when our algorithm terminates,
 $k \leq 2 \log \log^+(\eta_{\text{max}}/\eta_{\epsilon})$.
- (2) Approximates the (possibly non-existent) root $\eta = \varphi(\eta)$ up to a factor of 2, even when root is non-existent. If the returned interval is $[\eta_{\text{low}}^*, \eta_{\text{high}}^*]$ then

$$\frac{\bar{r}_T(\eta_0)}{2\sqrt{\alpha G_T(\eta_{\text{high}}^*) + \beta}} \leq \eta_0 \leq \frac{\bar{r}_T(\eta_{\text{low}}^*)}{\sqrt{\alpha G_T(\eta_0) + \beta}}$$

§2 Proposition 2 (RootFindingBisection)

Proposition 2

Let $\eta_0 = \text{RootFindingBisection}(\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta)$, where $\alpha > 1, \beta > 0, T \in \mathbb{N}$, and each $\eta > 0$. Assume $\eta_{\text{high}} > \varphi(\eta_{\text{high}})$. Let $\bar{x} = T^{-1} \sum_{i < T} x_i(\eta_0)$ be the average iterate. Under exact gradient setting:

(1) if $\eta_{\text{low}} \leq \varphi(\eta_{\text{low}})$ then for some $\eta' \in [\eta_0, 2\eta_0]$,

$$\|\bar{x} - x_0\| \leq \frac{2\alpha}{\alpha - 1} d_0 \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0 \sqrt{\alpha G_T(\eta') + \beta}}{T};$$

(2) if $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$, then $\eta_0 = \eta_{\text{low}}$, and

$$\|\bar{x} - x_0\| \leq \eta_0 \sqrt{\alpha G_T(\eta_0) + \beta} \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{d_0 \sqrt{\alpha G_T(\eta_0) + \beta} + \eta_0 G_T(\eta_0)}{T}.$$

§2 Proposition 2 (RootFindingBisection)

Proposition 2

Let $\eta_0 = \text{RootFindingBisection}(\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta)$, where $\alpha > 1, \beta > 0, T \in \mathbb{N}$, and each $\eta > 0$. Assume $\eta_{\text{high}} > \varphi(\eta_{\text{high}})$. Let $\bar{x} = T^{-1} \sum_{i < T} x_i(\eta_0)$ be the average iterate. Under exact gradient setting:

- (1) if $\eta_{\text{low}} \leq \varphi(\eta_{\text{low}})$ then for some $\eta' \in [\eta_0, 2\eta_0]$,

$$\|\bar{x} - x_0\| \leq \frac{2\alpha}{\alpha - 1} d_0 \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0 \sqrt{\alpha G_T(\eta') + \beta}}{T};$$

- (2) if $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$, then $\eta_0 = \eta_{\text{low}}$, and

$$\|\bar{x} - x_0\| \leq \eta_0 \sqrt{\alpha G_T(\eta_0) + \beta} \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{d_0 \sqrt{\alpha G_T(\eta_0) + \beta} + \eta_0 G_T(\eta_0)}{T}.$$

§2 Main Theorem (exact gradient)

Theorem: (exact gradient version)

Let $\alpha^{(k)} = 3, \beta^{(k)} = 0, n_\epsilon > 0, B \in \mathbb{N}$, and $x_0 \in \mathbb{R}^d$. With a total gradient budget of B , under exact gradient setting our algorithm will tune the learning rate to some $\eta \geq \eta_\epsilon$. Using the average \bar{x} of

$$T \geq \max\left(1, \frac{B}{12 \log \log^+(\|x_0 - x^*\|/(\eta_\epsilon \|g_0\|))}\right)$$

iterates, one of the following holds.

(1) $\eta > \eta_\epsilon$, and

$$\|\bar{x} - x^*\| \leq 4\|x_0 - x^*\|, \quad f(\bar{x}) - f(x^*) \leq \sqrt{27} \cdot \frac{\|x_0 - x^*\| \sqrt{G_T(\eta')}}{T},$$

(2) or $\eta = \eta_\epsilon$, and

$$\|\bar{x} - x^*\| \leq \eta_\epsilon \sqrt{3G_T(\eta_\epsilon)}, \quad f(\bar{x}) - f(x^*) \leq \frac{2\eta_\epsilon G_T(\eta_\epsilon)}{T}.$$

This page is intentionally left blank.

§3 Moving Forward — Defining “Good Events”

- (1) Key observation: when \mathcal{O} outputs exact gradients, $g_i(\eta) \equiv \nabla f(x_i(\eta))$.
- (2) This means that under exact gradient setting,

$$\sum_{i < T} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle = 0.$$

- (3) Generalize above into “approximately:” for $T \in \mathbb{N}$, and $\alpha, \beta, \eta > 0$, define the “**good events**” to be

$$\mathfrak{E}(\eta) = \mathfrak{E}(\eta; T, \alpha, \beta) := \bigcap_{t \leq T} \left\{ \sum_{i < t} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle \geq -\frac{1}{4} \max(\bar{d}_t(\eta), \eta\sqrt{\beta}) \sqrt{\alpha G_t(\eta) + \beta} \right\}.$$

§3 Moving Forward — Defining “Good Events”

- (1) Key observation: when \mathcal{O} outputs exact gradients, $g_i(\eta) \equiv \nabla f(x_i(\eta))$.
- (2) This means that under exact gradient setting,

$$\sum_{i < T} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle = 0.$$

- (3) Generalize above into “approximately:” for $T \in \mathbb{N}$, and $\alpha, \beta, \eta > 0$, define the “good events” to be

$$\mathfrak{E}(\eta) = \mathfrak{E}(\eta; T, \alpha, \beta) := \bigcap_{i < T} \left\{ \sum_{i < i} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle \geq -\frac{1}{4} \max(\bar{d}_i(\eta), \eta\sqrt{\beta}) \sqrt{\alpha G_i(\eta) + \beta} \right\}.$$

§3 Moving Forward — Defining “Good Events”

- (1) Key observation: when \mathcal{O} outputs exact gradients,
 $g_i(\eta) \equiv \nabla f(x_i(\eta))$.
- (2) This means that under exact gradient setting,

$$\sum_{i < T} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle = 0.$$

- (3) Generalize above into “approximately:” for $T \in \mathbb{N}$, and $\alpha, \beta, \eta > 0$, define the “**good events**” to be

$$\mathfrak{E}(\eta) = \mathfrak{E}(\eta; T, \alpha, \beta) := \bigcap_{t \leq T} \left\{ \sum_{i < t} \langle \Delta_i(\eta), x_i(\eta) - x^* \rangle \geq -\frac{1}{4} \max(\bar{d}_t(\eta), \eta\sqrt{\beta}) \sqrt{\alpha G_t(\eta) + \beta} \right\}.$$

What Was That Mess?

Lemma 1 (exact gradient version)

With appropriate parameters, under exact gradient setting,

$$\eta \leq \varphi(\eta) \Rightarrow \bar{d}_T(\eta) \leq \frac{\alpha + 1}{\alpha - 1} \cdot d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{2\alpha}{\alpha - 1} d_0.$$

becomes ...

Lemma 1 (stochastic version)

With appropriate parameters, under $\mathfrak{E}(\eta; T, \alpha, \beta)$, i.e., the “good event” setting, if $\eta \leq \varphi(\eta)$, then

$$\bar{d}_T(\eta) \leq \frac{3\alpha + 2}{\alpha + 2} d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{4\alpha}{\alpha - 2} d_0.$$

What Was That Mess?

Lemma 1 (exact gradient version)

With appropriate parameters, under exact gradient setting,

$$\eta \leq \varphi(\eta) \Rightarrow \bar{d}_T(\eta) \leq \frac{\alpha + 1}{\alpha - 1} \cdot d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{2\alpha}{\alpha - 1} d_0.$$

becomes ...

Lemma 1 (stochastic version)

With appropriate parameters, under $\mathfrak{E}(\eta; T, \alpha, \beta)$, i.e., the “good event” setting, if $\eta \leq \varphi(\eta)$, then

$$\bar{d}_T(\eta) \leq \frac{3\alpha + 2}{\alpha + 2} d_0 \quad \text{and} \quad \bar{r}_T(\eta) \leq \frac{4\alpha}{\alpha - 2} d_0.$$

Proposition 2 (exact gradient version)

Let $\eta_0 = \text{RootFindingBisection}(\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta)$, where $\alpha > 1, \beta > 0, T \in \mathbb{N}$, and each $\eta > 0$. Assume $\eta_{\text{high}} > \varphi(\eta_{\text{high}})$. Let $\bar{x} = T^{-1} \sum_{i < T} x_i(\eta_0)$ be the average iterate. Under exact gradient setting:

(1) if $\eta_{\text{low}} \leq \varphi(\eta_{\text{low}})$ then for some $\eta' \in [\eta_0, 2\eta_0]$,

$$\|\bar{x} - x_0\| \leq \frac{2\alpha}{\alpha - 1} d_0 \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0 \sqrt{\alpha G_T(\eta') + \beta}}{T};$$

(2) if $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$, then $\eta_0 = \eta_{\text{low}}$, and

$$\|\bar{x} - x_0\| \leq \eta_0 \sqrt{\alpha G_T(\eta_0) + \beta} \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{d_0 \sqrt{\alpha G_T(\eta_0) + \beta} + \eta_0 G_T(\eta_0)}{T}.$$

Proposition 2 (stochastic version)

Let $\eta_0 = \text{RootFindingBisection}(\eta_{\text{low}}, \eta_{\text{high}}; T, \alpha, \beta)$, where $\alpha > 2, \beta > 0, T \in \mathbb{N}$, and $\eta_{\text{high}} = 2^{2^k} \eta_{\text{low}}$ for some k . Assume $\eta_{\text{high}} > \varphi(\eta_{\text{high}})$. Let $\bar{x} = T^{-1} \sum_{i < T} x_i(\eta_0)$ be the average iterate. Assume the “good events” $\bigcap_{j=0}^{2^k} \mathfrak{G}(2^j \eta_{\text{low}}; T, \alpha, \beta)$ all hold.

(1) If $\eta_{\text{low}} \leq \varphi(\eta_{\text{low}})$, then for some $\eta' \in [\eta_0, 2\eta_0]$,

$$\|\bar{x} - x_0\| \leq \frac{4\alpha}{\alpha - 2} d_0 \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{9\alpha - 2}{2(\alpha - 2)} \cdot \frac{d_0 \sqrt{\alpha G_T(\eta') + \beta}}{T};$$

(2) If $\eta_{\text{low}} > \varphi(\eta_{\text{low}})$ and in addition $\mathfrak{G}(\eta_{\text{low}}; T, \alpha, \beta)$ holds, then $\eta_0 = \eta_{\text{low}}$, and

$$\|\bar{x} - x_0\| \leq \eta_{\text{low}} \sqrt{\alpha G_T(\eta_{\text{low}}) + \beta} \quad \text{and} \quad f(\bar{x}) - f(x^*) \leq \frac{5}{4} \frac{d_0 \sqrt{\alpha G_T(\eta_{\text{low}}) + \beta} + \eta_{\text{low}} (\alpha G_T(\eta_{\text{low}} + \beta))}{T}.$$

What Was That Mess?

Lemma 2 (exact gradient version)

With appropriate parameters, under exact gradient setting, if the following holds, then $\eta > \varphi(\eta)$:

$$\eta > \eta_{\max} := \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}}.$$

Consequently, when our algorithm terminates, $k \leq 2 \log \log^+(\eta_{\max}/\eta\epsilon)$.

becomes...

Lemma 2 (stochastic version)

With appropriate parameters, if “good event” $\mathfrak{E}(\eta; T, \alpha, \beta)$ holds, then if the following implies $\eta > \varphi(\eta)$:

$$\eta > \eta_{\max} := \frac{4\alpha}{\alpha - 2} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}}.$$

Consequently, if $\bigcap_{k=2,4,8,\dots} \mathfrak{E}(2^{2^k} \eta\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$ holds, when our algorithm terminates, $k \leq 2 \log \log^+(\eta_{\max}/\eta\epsilon)$.

What Was That Mess?

Lemma 2 (exact gradient version)

With appropriate parameters, under exact gradient setting, if the following holds, then $\eta > \varphi(\eta)$:

$$\eta > \eta_{\max} := \frac{2\alpha}{\alpha - 1} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}}.$$

Consequently, when our algorithm terminates, $k \leq 2 \log \log^+(\eta_{\max}/\eta_\epsilon)$.

becomes...

Lemma 2 (stochastic version)

With appropriate parameters, if “good event” $\mathfrak{E}(\eta; T, \alpha, \beta)$ holds, then if the following implies $\eta > \varphi(\eta)$:

$$\eta > \eta_{\max} := \frac{4\alpha}{\alpha - 2} \cdot \frac{d_0}{\sqrt{\alpha\|g_0\|^2 + \beta}}.$$

Consequently, if $\bigcap_{k=2,4,8,\dots} \mathfrak{E}(2^{2^k} \eta_\epsilon; T_k, \alpha^{(k)}, \beta^{(k)})$ holds, when our algorithm terminates, $k \leq 2 \log \log^+(\eta_{\max}/\eta_\epsilon)$.

§3 “Good Events” Are Likely

For the remainder of the analysis, just like **Fact 1**, we assume the gradient oracle is uniformly bounded by $L > 0$.

Lemma 3: “good events” are likely

Let $T \in \mathbb{N}$, $\eta > 0$, $\delta \in (0, 1)$ be given. Define $C = \log(60\delta^{-1} \log^2(6T))$.

If $\alpha \geq 1024C$ and $\beta \geq 1024C^2L^2$ then $\mathbb{P}(\mathfrak{E}(\eta; T, \alpha, \beta)) \geq 1 - \delta$.

§3 “Good Events” Are Likely

For the remainder of the analysis, just like **Fact 1**, we assume the gradient oracle is uniformly bounded by $L > 0$.

Lemma 3: “good events” are likely

Let $T \in \mathbb{N}$, $\eta > 0$, $\delta \in (0, 1)$ be given. Define $C = \log(60\delta^{-1} \log^2(6T))$.

If $\alpha \geq 1024C$ and $\beta \geq 1024C^2L^2$ then $\mathbb{P}(\mathfrak{E}(\eta; T, \alpha, \beta)) \geq 1 - \delta$.

§3 “Good Events” Are Likely

Proposition 3

Let budget B , initial step size $\eta_\epsilon > 0$, and failure probability $\delta \in (0, 1)$ be given. Let

$\alpha^{(k)} = 1024C_k$ and $\beta^{(k)} = 1024C_k^2L^2$, where $C_k = 2k + \log(60\delta^{-1}\log^2(6B))$.

Then, $\mathbb{P}(\bigcap_{k=2,4,8,\dots} \bigcap_{j=0,1,\dots,2^k} \mathfrak{E}(2^j n_\epsilon; B, \alpha^{(k)}, \beta^{(k)})) \geq 1 - \delta$.

Proof. Notice that $C_k = \log(60 \log^2(6B)/(2^{-2k}\delta))$ so by the previous lemma, with $T = B$, $\alpha = \alpha^{(k)}$, $\beta = \beta^{(k)}$, and failure probability $2^{-2k}\delta$, for any η ,

$$1 - \mathbb{P}(\mathfrak{E}(\eta; B, \alpha^{(k)}, \beta^{(k)})) \leq 2^{-2k}\delta.$$

By union bound

$$1 - \mathbb{P}\left(\bigcap_{j=0}^{2^k} \mathfrak{E}(2^j \eta_\epsilon; B, \alpha^{(k)}, \beta^{(k)})\right) \leq (2^k + 1)2^{-2k}\delta \leq 2^{-(k-1)}\delta$$

and finally

$$1 - \mathbb{P}\left(\bigcap_{k=2,4,8,\dots} \bigcap_{j=0}^{2^k} \mathfrak{E}(2^j \eta_\epsilon; B, \alpha^{(k)}, \beta^{(k)})\right) \leq \sum_{k \geq 1} 2^{-k}\delta = \delta.$$

Theorem: (stochastic version)

For any failure probability $\delta \in (0, 1)$, budget $B \in \mathbb{N}$, starting point $x_0 \in \mathbb{R}^d$, and initial step size $\eta_\epsilon > 0$, with $\{\alpha^{(k)}, \beta^{(k)}\}$ specified as in the previous proposition, the algorithm (i) makes $\leq B$ gradient queries, (ii) fine-tunes the step size to $\eta \geq \eta_\epsilon$, and (iii) returns $\bar{x} = T^{-1} \sum_{i < T} x_i(\eta) \in \mathbb{R}^d$.

Define $C = -\log \delta + \log \log^+(B\|x^* - x_0\|/(\eta_\epsilon L))$. Then, for some $\eta' \in [\eta, 2\eta]$, the event $\{(1) \text{ and } ((2) \text{ or } (3))\}$ happens with probability $\geq 1 - \delta$.

$$T \geq \max \left(1, \frac{B}{8 \log \log^+(\|x_0 - x^*\|/(\eta_\epsilon L))} \right) \quad (1)$$

$$\|\bar{x} - x^*\| \leq 6\|x_0 - x^*\| \quad \text{and} \quad f(\bar{x}) - f(x^*) = O\left(\frac{\|x_0 - x^*\| \sqrt{CG_T(\eta') + C^2 L^2}}{T}\right) \quad (2)$$

$$\|\bar{x} - x^*\| = O\left(\eta_\epsilon \sqrt{CG_T(\eta_\epsilon) + C^2 L^2}\right) \quad \text{and} \quad f(\bar{x}) - f(x^*) = O\left(\frac{\eta_\epsilon (CG_T(\eta_\epsilon) + C^2 L^2)}{T}\right) \quad (3)$$



Yair Carmon and Oliver Hinder. “Making SGD Parameter-Free”. In: (2022). arXiv: 2205.02160 [math.OC].