



# Parameter-Free Stochastic Optimization

From Yair Carmon and Oliver Hinder

Presented by Yu-Hu Yan

2024.09.27

# Stochastic Optimization

## □ Problem Setup

$$\begin{aligned} \min \quad & \textcolor{red}{f}(x) \\ \text{s.t.} \quad & x \in \mathbb{R}^d \end{aligned}$$

In this talk, we consider: (i) *unbounded* domain  $\mathbb{R}^d$ ;  
(ii) *convex G-Lipschitz* function  $f(\cdot)$ .

Stochastic noise formulation:

$$\boxed{\begin{aligned} \mathbb{E}[\mathbf{g}(\mathbf{x})] &= \nabla f(\mathbf{x}), \\ \mathbb{E} [\|\nabla f(\mathbf{x}) - \mathbf{g}(\mathbf{x})\|^2 \mid \mathbf{x}] &\leq \sigma^2 \end{aligned}}$$

# Parameter-Freeness

---

*Achieving **parameter-freeness** is important in both **online** and **offline** optimization.*

Common parameters in *stochastic* optimization:

- $G$ -Lipschitzness:  $f(x) - f(y) \leq G\|x - y\|$
- $L$ -smoothness:  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$
- initial optimality:  $d_0 \triangleq \|x_0 - x^*\|$
- stochastic noise:  $\|g(x) - \nabla f(x)\| \leq \sigma$

# Benchmark

---

## □ Benchmark method

**SGD:**  $x_{t+1} = x_t - \eta_t g_t$  as simple as **OGD**

Although SGD is very simple, its analysis in stochastic optimization is far from complete.

## □ Benchmark rate for *convex Lipschitz* case

Tuned SGD achieves the *optimal* rate of

$$f(\bar{x}) - f(x^*) \leq \mathcal{O}\left(\frac{d_0(G + \sigma)}{\sqrt{T}}\right) \quad (\bar{x} \text{ being some statistic of } \{x_t\}_{t=1}^T)$$

- $G$ -Lipschitzness:  $f(x) - f(y) \leq G\|x - y\|$
- $L$ -smoothness:  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$
- initial optimality:  $d_0 \triangleq \|x_0 - x^*\|$
- stochastic noise:  $\|g(x) - \nabla f(x)\| \leq \sigma$

Can we obtain the optimal rate **without these parameters?**

# Making SGD Parameter-Free



Oliver Hinder  
University of Pittsburgh



Yair Carmon  
Tel Aviv University

# Ideal Step Size

- With parameters, what is the *ideal* step size?

To see this, we start with a preliminary analysis of SGD.

If we run SGD for  $T$  rounds and submit  $\bar{x} \triangleq \frac{1}{T} \sum_{t=1}^T x_t$ .

$$\begin{aligned} f(\bar{x}) - f(x_\star) &\leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x_\star) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x_\star \rangle \\ &= \boxed{\frac{1}{T} \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle} + \boxed{\frac{1}{T} \sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle} \quad (\Delta_t \triangleq \nabla f(x_t) - g_t) \end{aligned}$$

*algorithm-related*      *algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

**Next step:** analyze  $\sum_{t=1}^T \langle g_t, x_t - x_\star \rangle$  for ideal step size.

# Ideal Step Size

- With parameters, what is the *ideal* step size?

To see this, we start with a preliminary analysis of SGD.

If we run SGD for  $T$  rounds and submit  $\bar{x} \triangleq \frac{1}{T} \sum_{t=1}^T x_t$ .

$$\begin{aligned} f(\bar{x}) - f(x_\star) &\leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x_\star) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x_\star \rangle \\ &= \boxed{\frac{1}{T} \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle} + \boxed{\frac{1}{T} \sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle} \quad (\Delta_t \triangleq \nabla f(x_t) - g_t) \end{aligned}$$

*algorithm-related*      *algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

**Simplification:** consider *noiseless* case of  $g_t = \nabla f(x_t)$ .

# Ideal Step Size

□ With parameters, what is the *ideal* step size?

Next step: analyze  $\sum_{t=1}^T \langle g_t, x_t - x_\star \rangle$ .  $\Rightarrow$  similar analysis to OGD

For simplicity, we denote by  $d_t \triangleq \|x_t - x_\star\|$  and  $G_t \triangleq \sum_{i \leq t} \|g_i\|^2$ .

$$d_{t+1}^2 = \|x_{t+1} - x_\star\|^2 = \|x_t - \eta g_t - x_\star\|^2 = d_t^2 - 2\eta \langle g_t, x_t - x_\star \rangle + \eta^2 \|g_t\|^2 \quad (\text{consider constant step size})$$

$$\Rightarrow \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \approx \frac{d_0^2 - d_T^2}{\eta} + \eta G_T \leq \frac{d_0^2}{\eta} + \eta G_T \quad \Rightarrow \quad \text{ideal step size: } \eta_\star \approx \frac{d_0}{\sqrt{G_T}}$$

**Issue:** to achieve the optimal rate  $\mathcal{O}(d_0 G / \sqrt{T})$ , the step size  $\eta$  must use initial optimality  $d_0 \triangleq \|x_0 - x_\star\|$ , which is *impossible*.

# Proxy for Ideal Step Size

□ What kind of step size is implementable?

Back to  $\sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \approx \frac{d_0^2 - d_T^2}{\eta} + \eta G_T$

**Key insight:** the negative term of  $-d_T^2$  can be useful.

For simplicity, we denote by  $r_t \triangleq \|x_t - x_0\|$  and  $\bar{r}_t \triangleq \max_{i \leq t} \|x_i - x_0\|$ .

$$\frac{d_0^2 - d_T^2}{\eta} + \eta G_T = \frac{(d_0 + d_T)(d_0 - d_T)}{\eta} + \eta G_T$$

Consider two cases:

- Case  $d_T \leq d_0$ : then  $d_0 + d_T \leq 2d_0$  and  $d_0 - d_T = \|x_0 - x_\star\| - \|x_T - x_\star\| \leq \|x_T - x_0\| = r_T \leq \bar{r}_T$ .
- Case  $d_T > d_0$ : then  $d_0^2 - d_T^2 < 0 \leq 2d_0 \bar{r}_T$ .

→  $d_0^2 - d_T^2 \lesssim d_0 \bar{r}_T$ .    *What is different?*

# Proxy for Ideal Step Size

□ What kind of step size is implementable?

$$d_0^2 - d_T^2 \lesssim d_0 \bar{r}_T \quad \Rightarrow \quad \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \approx \frac{d_0^2 - d_T^2}{\eta} + \eta G_T \leq \frac{d_0 \bar{r}_T}{\eta} + \eta G_T$$

---

⇒ We can choose a *proxy* for the ideal step size:  $\eta_\dagger \approx \frac{\bar{r}_T}{\sqrt{\alpha G_T}}$  ( $\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\|$ )  
( $\alpha$  only for scalability)

Since  $x_0$  is known, the quantity  $\{\bar{r}_t\}_{t \in [T]}$  is *trackable*.

$$\Rightarrow \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \lesssim \frac{d_0 \bar{r}_T}{\eta} + \eta G_T \lesssim (d_0 + \bar{r}_T) \sqrt{G_T} \quad \begin{matrix} \text{very close to the optimal rate} \\ \mathcal{O}(d_0 G / \sqrt{T}). \end{matrix}$$

**Next step:** connect  $\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\|$  to  $d_0 \triangleq \|x_0 - x_\star\|$  with  $\eta = \eta_\dagger$ .

# Proxy for Ideal Step Size

## □ Proof

**Next step:** relate  $\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\|$  to  $d_0 \triangleq \|x_0 - x_\star\|$  with  $\eta = \eta_\dagger$ .

---

From standard analysis:  $0 \leq \sum_{t=1}^T f(x_t) - \sum_{t=1}^T f(x_\star) \leq \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \leq \frac{d_0^2 - d_T^2}{\eta} + \eta G_T$

$$\Rightarrow \textcolor{blue}{d_T^2} \leq d_0^2 + \eta^2 G_T \quad \Rightarrow \text{Taking } \eta = \eta_\dagger = \frac{\bar{r}_T}{\sqrt{\alpha G_T}} \text{ leads to } \textcolor{blue}{d_T^2} \leq d_0^2 + \frac{1}{\alpha} \bar{r}_T^2.$$

$$\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\| \leq \max_{t \leq T} \|x_t - x_\star\| + \|x_0 - x_\star\| = \bar{d}_T + d_0 \quad (\bar{d}_T \triangleq \max_{t \leq T} \|x_t - x_\star\|)$$

$$\Rightarrow \textcolor{blue}{\bar{d}_T^2} \leq d_0^2 + \frac{1}{\alpha} (\bar{d}_T + d_0)^2 \quad \Rightarrow \textcolor{red}{\bar{d}_T} \leq \frac{\alpha + 1}{\alpha - 1} \cdot d_0 \quad \& \quad \textcolor{red}{\bar{r}_T} \leq \frac{2\alpha}{\alpha - 1} \cdot d_0$$

(taking max on both sides)

# One More Thing

---

□ Is this step size implementable?  $\eta_t \approx \frac{\bar{r}_T}{\sqrt{\alpha G_T}}$

**No!** Because  $\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\|$  and  $G_T \triangleq \sum_{t \leq T} \|g_t\|^2$  are only available *after* the algorithm ends.

□ An observation

Consider *noiseless* case (i.e.,  $g_t = \nabla f(x_t)$ ), given step size  $\eta$ , the decision  $x_t$  is a function of  $\eta$ , marked as  $x_t(\eta)$ . As a result,  $\bar{r}_T(\eta)$  and  $G_T(\eta)$  are both functions of  $\eta$ .

**fixed-point** problem:  $\eta = \frac{\bar{r}_T(\eta)}{\sqrt{\alpha G_T(\eta)}}$

⇒ **Solution:** *bisection* method

# Bisection

## □ Common bisection

Given initial bound  $[a^{(1)}, b^{(1)}]$ , w.l.o.g., choosing  $a^{(2)} = \frac{a^{(1)} + b^{(1)}}{2}$  and  $b^{(2)} = b^{(1)}$ ,

$$b^{(2)} - a^{(2)} = b^{(1)} - \frac{a^{(1)} + b^{(1)}}{2} = \frac{1}{2} \cdot (b^{(1)} - a^{(1)}) \quad \Rightarrow \text{Each bisection halves } (b^{(i)} - a^{(i)}).$$

With  $(b^{(1)} - a^{(1)})$ , it runs  $\log T$  times to make it  $\frac{1}{T}$  times smaller.

## □ This work's bisection

Given initial bound  $[a^{(1)}, b^{(1)}]$ , w.l.o.g., choosing  $a^{(2)} = \sqrt{a^{(1)}b^{(1)}}$  and  $b^{(2)} = b^{(1)}$ ,

$$\log \frac{b^{(2)}}{a^{(2)}} = \log \frac{b^{(1)}}{\sqrt{a^{(1)}b^{(1)}}} = \log \sqrt{\frac{b^{(1)}}{a^{(1)}}} = \frac{1}{2} \cdot \log \frac{b^{(1)}}{a^{(1)}} \quad \Rightarrow \text{Each bisection halves } \log \frac{b^{(i)}}{a^{(i)}}.$$

With  $\frac{b^{(1)}}{a^{(1)}}$ , it runs  $\log \log T$  times to make it  $\frac{1}{T}$  times smaller.

# Bisection

□ This work's bisection

$$\text{fixed-point problem: } \eta = h(\eta) \triangleq \frac{\bar{r}_T(\eta)}{\sqrt{\alpha G_T(\eta)}}$$

If we can ensure that there exists an interval  $[\eta_-, \eta_+]$  such that  $\eta_- \leq h(\eta_-)$  and  $\eta_+ \geq h(\eta_+)$ , the bisection method can find  $\eta_\dagger$  with  $\log \log \frac{\eta_+}{\eta_-}$  steps.

**Next step:** how to find a valid interval  $[\eta_-, \eta_+]$  with  $\eta_- \leq h(\eta_-)$  and  $\eta_+ \geq h(\eta_+)$ .

**Solution:** start with  $\eta_- = \eta_\epsilon$  with some small enough user-defined  $\eta_\epsilon$  and  $\eta_+ = 2^{2^k} \eta_\epsilon$  as upper bound.

**Lemma 1.** If  $\eta \geq \eta_{\max} \triangleq \frac{2\alpha}{\alpha-1} \cdot \frac{d_0}{\sqrt{\alpha \|g_0\|^2}}$  then  $\eta > h(\eta)$ . Then  $\mathcal{O}(\log \log \frac{\eta_{\max}}{\eta_\epsilon})$  bisection steps are enough.

The proof is by *contradiction*, simple, and thus omitted here.

# Useful in Online Ensemble?

- Can we use a similar step-size pool here?

Set the step size pool as

$$\left\{ \eta_1 = \frac{1}{\sqrt{T}}, \eta_i = 2^{2^i} \eta_1, \eta_N = \mathcal{O}(1) \right\} \quad \Rightarrow N \approx \mathcal{O}(\log \log T)$$

*More efficient* than the classic  $\mathcal{O}(\log T)$  size.

---

**How will this new step size pool affect the regret bound?**

Dynamic regret of the  $i$ -th base learner can be bounded by  $\mathcal{O}\left(\frac{1+P_T}{\eta_i} + \eta_i T\right)$ .

Suppose there exist  $\eta_{i^*}$  such that  $\eta_{i^*} \leq \eta_* \leq \eta_{i^*+1}$ .  $\Rightarrow \eta_{i^*} T \leq \eta_* T$

*(this term remains to be tuned)*

**Next step:** find a *lower* bound for  $\eta_{i^*}$  to handle the term  $\frac{1+P_T}{\eta_{i^*}}$ .

# Useful in Online Ensemble?

- Can we use a similar step-size pool here?

Set the step size pool as

$$\left\{ \eta_1 = \frac{1}{\sqrt{T}}, \eta_i = 2^{2^i} \eta_1, \eta_N = \mathcal{O}(1) \right\}$$

We take **full-information dynamic regret minimization** as an example.

$$\Rightarrow N \approx \mathcal{O}(\log \log T)$$

*More efficient* than the classic  $\mathcal{O}(\log T)$  size.

---

## How will this new step size pool affect the regret bound?

Suppose there exist  $\eta_{i^*}$  such that  $\eta_{i^*} \leq \eta_* \leq \eta_{i^*+1}$ .

**Next step:** find a *lower* bound for  $\eta_{i^*}$  to handle the term  $\frac{1+P_T}{\eta_{i^*}}$ .

$$\eta_{i+1} = 2^{2^{i+1}} \eta_1 \text{ and } \eta_i = 2^{2^i} \eta_1 \quad \Rightarrow \quad \frac{\eta^{i+1}}{\eta^i} = \frac{2^{2^{i+1}}}{2^{2^i}} = 2^{2^{i+1}-2^i} = 2^{2^i} \quad \Rightarrow \quad \eta_{i^*} \geq \frac{\eta_*}{2^{2^{i^*}}}$$

$$\Rightarrow \frac{1+P_T}{\eta_{i^*}} \leq \frac{1+P_T}{\eta_*} \cdot 2^{2^{i^*}}$$

**Last question:** How large is  $2^{2^{i^*}}$ ?

# Useful in Online Ensemble?

- Can we use a similar step-size pool here?

*We take **full-information dynamic regret minimization** as an example.*

Set the step size pool as

$$\left\{ \eta_1 = \frac{1}{\sqrt{T}}, \eta_i = 2^{2^i} \eta_1, \eta_N = \mathcal{O}(1) \right\}$$

$$\Rightarrow N \approx \mathcal{O}(\log \log T)$$

*More efficient* than the classic  $\mathcal{O}(\log T)$  size.

---

**How will this new step size pool affect the regret bound?**

Suppose there exist  $\eta_{i^*}$  such that  $\eta_{i^*} \leq \eta_* \leq \eta_{i^*+1}$ . **Last question:** How large is  $2^{2^{i^*}}$ ?

$$(\eta_{i^*} \triangleq 2^{2^{i^*}} \eta_1) \leq \eta_* \quad \Rightarrow \quad 2^{2^{i^*}} \leq \frac{\eta_*}{\eta_1} \approx \eta_* \cdot \sqrt{T}$$

$$\frac{1 + P_T}{\eta_{i^*}} + \eta_{i^*} T \leq \frac{1 + P_T}{\eta_*} \cdot 2^{2^{i^*}} + \eta_* T \leq \frac{1 + P_T}{\eta_*} \cdot \eta_* \cdot \sqrt{T} + \eta_{i^*} T \lesssim \sqrt{T}(1 + P_T)$$

*Useless.*  $\mathcal{O}(\log T)$  seems necessary.

# Stochasticity Analysis

## □ Recall the starting analysis

If we run SGD for  $T$  rounds and submit  $\bar{x} \triangleq \frac{1}{T} \sum_{t=1}^T x_t$ .

$$\begin{aligned} f(\bar{x}) - f(x_\star) &\leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x_\star) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x_\star \rangle \\ &= \frac{1}{T} \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle + \frac{1}{T} \sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle \quad (\Delta_t \triangleq \nabla f(x_t) - g_t) \end{aligned}$$

*algorithm-related*      *algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

**Next step:** analyze the stochastic gap  $\sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle$ .  $\Rightarrow$  equals zero for *expectation* bounds.

For *high-probability* result  $\Rightarrow$  Solution: *concentration*

# Stochasticity Analysis

- Concentration  $\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$  **Stochastic part:** gap between the true gradient  $\nabla f(x_t)$  and its stochastic sample  $g_t$ .
- ➡ **concentration** inequality

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference* sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

$$\text{where } \theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}.$$

How to construct a bounded sequence?  $\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

First try to import boundedness:  $(d_t \triangleq \|x_t - x_\star\|) \leq (\bar{d}_t \triangleq \max_{i \leq t} \|x_i - x_\star\|)$

$$\Rightarrow \sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle = \bar{d}_T \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{\bar{d}_T} \right\rangle$$

Is the concentration applicable?

No!  $X_t = \left\langle \Delta_t, \frac{x_t - x_\star}{\bar{d}_T} \right\rangle$  is not adapted to  $\mathcal{F}_t$ .

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

How to construct a bounded sequence?  $\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

Second try to import boundedness:  $(d_t \triangleq \|x_t - x_\star\|) \leq (\bar{d}_t \triangleq \max_{i \leq t} \|x_i - x_\star\|)$

$$\sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle = \sum_{t=1}^T \textcolor{blue}{d_t} \cdot \left\langle \Delta_t, \frac{x_t - x_\star}{d_t} \right\rangle \leq \textcolor{red}{\bar{d}_T} \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{d_t} \right\rangle \quad X_t = \langle \Delta_t, (x_t - x_\star)/d_t \rangle \text{ is now adapted to } \mathcal{F}_t.$$

Is this correct? No!  $\langle \Delta_t, (x_t - x_\star)/d_t \rangle$  can be negative.

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

How to construct a bounded sequence?  $\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

$$\sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle = s_{k_T} \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle$$

The same issue that  $X_t = \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle$  is not adapted to  $\mathcal{F}_t$ .

**Solution:**  $s_{k_t} \triangleq 2^{k_t} d_0$ , where  $k_t \triangleq \left\lceil \log \frac{d_t}{d_0} \right\rceil$   $\Rightarrow$  Intuitively,  $d_t \leq s_{k_t} \leq 2d_t$

**Key idea:** use *discrete*  $s_{k_t}$  to replace the *continuous*  $d_t$ .

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a martingale difference sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

$$\text{where } \theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}.$$

How to construct a bounded sequence?

$\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

$$\sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle = s_{k_T} \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle$$

**Solution:**  $s_{k_t} \triangleq 2^k d_0$ , where  $k_t \triangleq \left\lceil \log \frac{d_t}{d_0} \right\rceil$   
 ↳ Intuitively,  $d_t \leq s_{k_t} \leq 2d_t$

**Key idea:** use *discrete*  $s_{k_t}$  to replace the *continuous*  $d_t$ .

$$\Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle \geq X \right] \leq \sum_k \Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_k} \right\rangle \geq X \right]$$

(union bound)

$s_k$  is now *unrelated* to the filtration.

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

$$\text{where } \theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}.$$

How to construct a bounded sequence?

$\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

$$\sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle = s_{k_T} \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle$$

**Solution:**  $s_{k_t} \triangleq 2^k d_0$ , where  $k_t \triangleq \left\lceil \log \frac{d_t}{d_0} \right\rceil$   
 ↳ Intuitively,  $d_t \leq s_{k_t} \leq 2d_t$

**Key idea:** use *discrete*  $s_{k_t}$  to replace the *continuous*  $d_t$ .

$$\Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle \geq X \right] \leq \sum_k \Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_k} \right\rangle \geq X \right]$$

**Next step:** how many possible  $s_k$  in the  $T$ -th round.

# Stochasticity Analysis

□ How to handle *boundedness*?

$$\sum_{t=1}^T \langle \nabla f(x_t) - g_t, x_t - x_\star \rangle$$

$x_t - x_\star$  is *unbounded!*

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

$$\text{where } \theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}.$$

How to construct a bounded sequence?

$\nabla f(x_t) - g_t$  is bounded by assuming *bounded noise*.

**Solution:**  $s_{k_t} \triangleq 2^k d_0$ , where  $k_t \triangleq \left\lceil \log \frac{d_t}{d_0} \right\rceil$

→ Intuitively,  $d_t \leq s_{k_t} \leq 2d_t$

**Key idea:** use *discrete*  $s_{k_t}$  to replace the *continuous*  $d_t$ .

$$\Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_{k_T}} \right\rangle \geq X \right] \leq \sum_k \Pr \left[ \sum_{t=1}^T \left\langle \Delta_t, \frac{x_t - x_\star}{s_k} \right\rangle \geq X \right]$$

**Next step:** how many possible  $s_k$  in the  $T$ -th round.

$$d_{t+1} = \|x_{t+1} - x_\star\| = \|x_t - \eta g_t - x_\star\| \leq d_t + \eta(G + \sigma)$$

$$\rightarrow d_t = \mathcal{O}(t)$$

$$\rightarrow k_t = \mathcal{O}(\log t)$$

# Overall

## □ Final result overview

**Theorem 7** In the noiseless setting Algorithm 1, with parameters  $\alpha^{(k)} = 3$ ,  $\beta^{(k)} = 0$ ,  $\eta_\varepsilon > 0$ ,  $B \in \mathbb{N}$ , and  $x_0 \in \mathbb{R}^d$ , performs at most  $B$  subgradient queries and returns  $\bar{x} = \frac{1}{T} \sum_{i < T} x_i(\eta) \in \mathbb{R}^d$  for some  $\eta \geq \eta_\varepsilon$  and integer  $T$  satisfying

$$T \geq \max \left\{ \frac{B}{12 \log \log_+ \frac{\|x_0 - x_\star\|}{\eta_\varepsilon \|g_0\|}}, 1 \right\} \quad (10)$$

such that either

$$\|\bar{x} - x_\star\| \leq 4\|x_0 - x_\star\| \text{ and } f(\bar{x}) - f(x_\star) \leq \sqrt{27} \frac{\|x_0 - x_\star\| \sqrt{G_T(\eta')}}{T} \quad (11)$$

for some  $\eta' \in [\eta, 2\eta]$ , or  $\eta = \eta_\varepsilon$  and

$$\|\bar{x} - x_\star\| \leq \eta_\varepsilon \sqrt{3G_T(\eta_\varepsilon)} \text{ and } f(\bar{x}) - f(x_\star) \leq 2 \frac{\eta_\varepsilon G_T(\eta_\varepsilon)}{T}. \quad (12)$$

### Remarks:

- (i) Need to consider the corner case where  $\eta_\varepsilon$  cannot serve as the lower bound for bisection.
- (ii) Not fully parameter-free since the learner requires  $T$ .
- (iii) Only with a  $\log \log T$  overhead compared with the tuned optimal bound.

# Overall

---

## □ Three key ideas overview

- Formulate the problem as bisection
- An aggressive bisection method
- Construct bounded sequence for concentration

# One Question

- Still fixed-point problem with stochastic gradients?

$$\text{fixed-point problem: } \eta = \frac{\bar{r}_T(\eta)}{\sqrt{\alpha G_T(\eta)}}$$

For *noiseless* (deterministic) case, i.e.,  $g_t = \nabla f(x_t)$ , given  $\eta$ ,  $G_T \triangleq \sum_{t \leq T} \|g_t\|^2$  can be *uniquely* decided, and we can consider  $G_T$  as  $G_T(\eta)$ .

For *stochastic* case, i.e.,  $\mathbb{E}[g_t] = \nabla f(x_t)$ , given  $\eta$ ,  $G_T \triangleq \sum_{t \leq T} \|g_t\|^2$  *cannot* be uniquely decided due to the additional *noise in  $\{g_t\}_{t=1}^T$* . Can  $G_T$  still be represented as  $G_T(\eta)$ ?



# DOG is SGD's Best Friend: A Parameter-Free Dynamic Step Size Schedule



Maor Ivgi  
Tel Aviv University



Oliver Hinder  
University of Pittsburgh



Yair Carmon  
Tel Aviv University

# Difference

## □ What's different from COLT'22?

COLT 2022

*time-invariant* step size

$$\eta_{\dagger} \approx \frac{\bar{r}_{\textcolor{blue}{T}}}{\sqrt{\alpha G_{\textcolor{blue}{T}}}} \quad (\bar{r}_T \triangleq \max_{t \leq T} \|x_t - x_0\|) \\ (G_T \triangleq \sum_{t \leq T} \|g_t\|^2)$$

Bisection is needed

Analysis is *complicated*

This work

*time-varying* step size  
**DOG (Distance over Gradients)**

$$\eta_t = \frac{\bar{r}_{\textcolor{red}{t}}}{\sqrt{\alpha G_{\textcolor{red}{t}}}} \quad (\bar{r}_t \triangleq \max_{i \leq t} \|x_i - x_0\|) \\ (G_t \triangleq \sum_{i \leq t} \|g_i\|^2)$$

Bisection is *not* needed

Analysis is *much simpler*

# Analysis

## □ Starting from the beginning

If we run SGD for  $T$  rounds and submit  $\bar{x} \triangleq \frac{1}{T} \sum_{t=1}^T x_t$ .

$$\begin{aligned} f(\bar{x}) - f(x_\star) &\leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x_\star) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x_\star \rangle \\ &= \frac{1}{T} \sum_{t=1}^T \langle g_t, x_t - x_\star \rangle + \frac{1}{T} \sum_{t=1}^T \langle \Delta_t, x_t - x_\star \rangle \quad (\Delta_t \triangleq \nabla f(x_t) - g_t) \end{aligned}$$

*algorithm-related*      *algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

**Next step:** analyze  $\sum_{t=1}^T \langle g_t, x_t - x_\star \rangle$ .

# Algorithm-Related Analysis

## □ What's different when using **DOG** step size?

By standard analysis as before:  $\sum_{t=1}^T \langle g_t, x_t - x_\star \rangle \lesssim \sum_{t=1}^T \frac{d_t^2 - d_{t+1}^2}{\eta_t} + \sum_{t=1}^T \eta_t \|g_t\|^2$

$$\sum_{t=1}^T \frac{d_t^2 - d_{t+1}^2}{\eta_t} \quad \rightarrow \quad \boxed{\eta_t = \frac{\bar{r}_{\textcolor{red}{t}}}{\sqrt{G_{\textcolor{red}{t}}}} \quad (\bar{r}_t \triangleq \max_{i \leq t} \|x_i - x_0\|) \quad (G_t \triangleq \sum_{i \leq t} \|g_i\|^2)} \quad \rightarrow \quad \sum_{t=1}^T \frac{d_t^2 - d_{t+1}^2}{\bar{r}_t} \cdot \sqrt{G_t}$$

*This term is hard to analyze due to the **DOG** step size.*

**Solution:** use a weighted combination  $\bar{x} \triangleq \frac{1}{\sum_{t=1}^T \bar{r}_t} \sum_{t=1}^T \bar{r}_t x_t$

# Algorithm-Related Analysis

□ What's different when using DOG step size?

$$f(\bar{x}) - f(x_\star) \leq \frac{1}{\sum_{t=1}^T \bar{r}_t} \left( \sum_{t=1}^T \bar{r}_t \cdot \langle g_t, x_t - x_\star \rangle + \sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle \right)$$

*algorithm-related*

*algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

$$\implies \sum_{t=1}^T \bar{r}_t \cdot \langle g_t, x_t - x_\star \rangle \lesssim \sum_{t=1}^T \frac{\bar{r}_t}{\eta_t} (d_t^2 - d_{t+1}^2) + \sum_{t=1}^T \eta_t \bar{r}_t \|g_t\|^2$$

$$\implies \sum_{t=1}^T \frac{\bar{r}_t}{\eta_t} (d_t^2 - d_{t+1}^2) \leq \sum_{t=1}^T (d_t^2 - d_{t+1}^2) \cdot \sqrt{G_t} \leq \dots \lesssim \bar{r}_T \bar{d}_T \sqrt{G_T} \quad (\text{standard derivation})$$

$$\implies \sum_{t=1}^T \eta_t \bar{r}_t \|g_t\|^2 \leq \sum_{t=1}^T \bar{r}_t^2 \cdot \frac{\|g_t\|^2}{G_t} \lesssim \bar{r}_T^2 \sqrt{G_T} \quad (\text{self-confident tuning})$$

$$\eta_t = \frac{\bar{r}_t}{\sqrt{G_t}} \quad (\bar{r}_t \triangleq \max_{i \leq t} \|x_i - x_0\|) \quad (G_t \triangleq \sum_{i \leq t} \|g_i\|^2)$$

$$(\Delta_t \triangleq \nabla f(x_t) - g_t)$$

# Algorithm-Related Analysis

- What's different when using DOG step size?

$$f(\bar{x}) - f(x_\star) \leq \frac{1}{\sum_{t=1}^T \bar{r}_t} \left( \sum_{t=1}^T \bar{r}_t \cdot \langle g_t, x_t - x_\star \rangle + \sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle \right)$$

*algorithm-related*

*algorithm-irrelated* due to *unknown*  $\nabla f(x_t)$

Algorithm-related part:

$$\sum_{t=1}^T \bar{r}_t \cdot \langle g_t, x_t - x_\star \rangle \lesssim \sum_{t=1}^T \frac{\bar{r}_t}{\eta_t} (d_t^2 - d_{t+1}^2) + \sum_{t=1}^T \eta_t \bar{r}_t \|g_t\|^2 \lesssim \bar{r}_T (\bar{d}_T + \bar{r}_T) \sqrt{G_T}$$

**Next step:** analyze the stochastic gap  $\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle$ .

# Stochasticity Analysis

## □ Weighted Concentration

To handle the stochastic gap  $\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle$ , we need a *weighted* version of concentration inequality.

**Lemma 2 (Uniform).** Let  $c > 0$  and  $X_t$  be a *martingale difference sequence* adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

# Stochasticity Analysis

## □ Weighted Concentration

To handle the stochastic gap  $\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle$ , we need a *weighted* version of concentration inequality.

**Lemma 3 (Weighted).** Let  $c > 0$  and  $X_t$  be a martingale difference sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t Y_i X_i \right| \geq 8Y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\{Y_t\}$  is a non-negative and non-decreasing sequence and  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

# Stochasticity Analysis

## □ Weighted Concentration

The proof is simple with the *uniform concentration* and a *useful technical lemma*.

**Lemma 3.** Let  $c > 0$  and  $X_t$  be a *martingale difference* sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t Y_i X_i \right| \geq 8Y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\{Y_t\}$  is a *non-negative and non-decreasing* sequence and  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

**Lemma 4.** Let  $\{a_t\}_{t=1}^T$  be a *non-negative and non-decreasing* sequence and  $\{b_t\}_{t=1}^T$  be any sequence. Then it holds that

$$\left| \sum_{i=1}^t a_i b_i \right| \leq 2a_t \max_{i \leq t} \left| \sum_{j=1}^i b_j \right|.$$

# Stochasticity Analysis

## □ Weighted Concentration

The proof is simple with the *uniform concentration* and a *useful technical lemma*.

**Lemma 4.** Let  $\{a_t\}_{t=1}^T$  be a *non-negative and non-decreasing* sequence and  $\{b_t\}_{t=1}^T$  be any sequence. Then it holds that

$$\left| \sum_{i=1}^t a_i b_i \right| \leq 2a_t \max_{i \leq t} \left| \sum_{j=1}^i b_j \right|.$$

**Lemma 3.** Let  $c > 0$  and  $X_t$  be a *martingale difference* sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t Y_i X_i \right| \geq 8Y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\{Y_t\}$  is a *non-negative and non-decreasing* sequence and  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

**Lemma 2.** Let  $c > 0$  and  $X_t$  be a *martingale difference* sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t X_i \right| \geq 4 \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .



⇒ With probability  $1 - \delta$ ,

$$\left| \sum_{i=1}^t Y_i X_i \right| \leq 2Y_t \max_{i \leq t} \left| \sum_{j=1}^i X_j \right| \leq 2Y_t \cdot \max_{i \leq t} 4 \sqrt{\theta_{i,\delta} \sum_{j=1}^i (X_j - \hat{X}_j)^2 + c^2 \theta_{i,\delta}^2} \leq 8Y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2}$$

# Stochasticity Analysis

## □ Weighted Concentration

**Lemma 3.** Let  $c > 0$  and  $X_t$  be a *martingale difference* sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq c$  with probability 1 for all  $t$ . Then for all  $\delta \in (0, 1)$ , and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq c$  with probability 1,

$$\Pr \left[ \exists t \leq T : \left| \sum_{i=1}^t Y_i X_i \right| \geq 8Y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^t (X_i - \hat{X}_i)^2 + c^2 \theta_{t,\delta}^2} \right] \leq \delta,$$

where  $\{Y_t\}$  is a *non-negative and non-decreasing* sequence and  $\theta_{t,\delta} \triangleq \log \frac{60 \log(6t)}{\delta}$ .

$$\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle = \sum_{t=1}^T \boxed{\bar{r}_t \bar{d}_t} \cdot \boxed{\left\langle \Delta_t, \frac{x_t - x_\star}{\bar{d}_t} \right\rangle} \quad (\bar{d}_t \triangleq \max_{i \leq t} \|x_i - x_\star\|)$$

$Y_t$                      $X_t$

Using **Lemma 3** accordingly, with probability  $1 - \delta$ ,

$$\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle \lesssim \bar{r}_T \bar{d}_T \sqrt{\theta_{T,\delta} G_T + \theta_{T,\delta}^2 G^2}$$

# Overall Analysis

## □ Combining two parts

**Algorithm-related part:**

$$\sum_{t=1}^T \bar{r}_t \cdot \langle g_t, x_t - x_\star \rangle \lesssim \bar{r}_T (\bar{d}_T + \bar{r}_T) \sqrt{G_T}$$

**Stochastic part:**

$$\sum_{t=1}^T \bar{r}_t \cdot \langle \Delta_t, x_t - x_\star \rangle \lesssim \bar{r}_T \bar{d}_T \sqrt{\theta_{T,\delta} G_T + \theta_{T,\delta}^2 G^2}$$

$$\Rightarrow f(\bar{x}) - f(x_\star) \leq \frac{1}{\sum_{t=1}^T \bar{r}_t} \left( \sum_{t=1}^T \bar{r}_t \langle g_t, x_t - x_\star \rangle + \sum_{t=1}^T \bar{r}_t \langle \Delta_t, x_t - x_\star \rangle \right)$$

$$\lesssim \frac{\bar{r}_T (\bar{d}_T + \bar{r}_T) \sqrt{\theta_{T,\delta} G_T + \theta_{T,\delta}^2 G^2}}{\sum_{t=1}^T \bar{r}_t} \lesssim \frac{(\textcolor{red}{d}_0 + \bar{r}_T) \sqrt{\theta_{T,\delta} G_T + \theta_{T,\delta}^2 G^2}}{\sum_{t=1}^T \bar{r}_t / \textcolor{red}{r}_T}$$

(Later details omitted...)

# Overall

---

## □ Three key ideas overview

- **DOG** step size (time-varying version of COLT 2022)
- Weighted combination for final decision
- Weighted concentration

# Summary

---

□ **Problem:** Parameter-free stochastic optimization

□ **Two Works from Yair Carmon & Oliver Hinder**

➤ Time-invariant step size with bisection method [COLT 2022]

- Key technique: *bisection formulation*

➤ Time-varying step size (DOG) [ICML 2023]

- Key technique: *weighted concentration*

*Thanks!*

□ **Following Work:** para.-free  $1/T^2$  rate of convex & smooth function [COLT 2024]