



rijksuniversiteit  
 groningen



Université  
Paris Cité

# SGD without Variance Assumption

New Tight Bounds via a Computer-Aided Lyapunov Analysis

---

**Daniel Cortild, Lucas Ketels, J. Peypouquet, G. Garrigos**

Journées Franco-Chilliennes pour l'Optimization

INSA Rouen, France, July 11th, 2025

# Stochastic Gradient Descent

Consider the problem

$$\min \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) : x \in \mathbb{R}^d \right\},$$

where all  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  are convex and  $L$ -smooth, and  $f$  has minimizers.

**Stochastic Gradient Descent** (SGD) iterates

$$x_0 \in \mathbb{R}^d, \quad x_{t+1} = x_t - \gamma \nabla f_{i_k}(x_t) \quad \text{for } t = 0, 1, \dots,$$

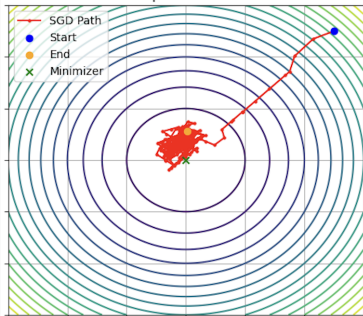
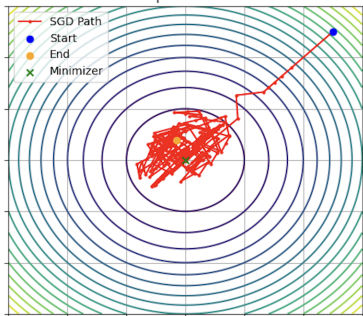
where  $i_k$  is chosen i.i.d. from the uniform distribution on  $\{1, \dots, n\}$ .

# Type of Results for SGD

Convergence results for SGD are usually presented as

$$\text{Performance}(t) \leq \text{Bias}(t) + \text{Variance}(t),$$

where  $\text{Bias}(t) \rightarrow 0$  as  $t \rightarrow \infty$ , and (ideally)  $\text{Variance}(t)$  remains bounded.



**Goal:** Minimize the bias term first, and the variance term second.

# Variance Assumption

**Typical Assumption:** Uniformly bounded gradient variance;

$$\sup_{x \in \mathbb{R}^d} \mathbb{E}[\|\nabla f_{i_k}(x) - \nabla f(x)\|^2] < +\infty.$$

However, this is unrealistic in practice.<sup>1</sup>

**Alternative Assumptions:** Weak growth, Strong growth, Maximal strong growth, Relaxed growth, etc.

**Our setting:** We define

$$\sigma_*^2 := \mathbb{E}[\|\nabla f_{i_k}(x_*)\|^2] \quad \text{for some } x_* \in \operatorname{argmin} f.$$

Note this is automatically finite in our setting.

---

<sup>1</sup>Bottou, Curtis, and Nocedal, "Optimization Methods for Large-Scale Machine Learning", 2018.

## Our Results in the Convex Setting

We obtain a result on the Cesàro average  $\bar{x}_T = \frac{x_0 + \dots + x_{T-1}}{T}$  of the form

$$\mathbb{E}[f(\bar{x}_T) - \min f] \leq \text{Bias}(T) \cdot \|x_0 - x_*\|^2 + \text{Variance}(T) \cdot \sigma_*^2,$$

where

	$\gamma L \in (0, 1)$	$\gamma L = 1$	$\gamma L \in (1, 2)$
Bias(T)	$\frac{1}{2\gamma T}$	$\frac{1}{(2 - \varepsilon)\gamma T}$	$\frac{1}{2\gamma(2 - \gamma L)T}$
Variance(T)	$\frac{\gamma}{2(1 - \gamma L)}$	$\frac{\gamma(2 + \varepsilon)}{\varepsilon(2 - \varepsilon)}$	$\frac{\exp(T)}{2 - \gamma L}$

**Observation 1:** Singularity at  $\gamma L = 1$  for optimal step-size.

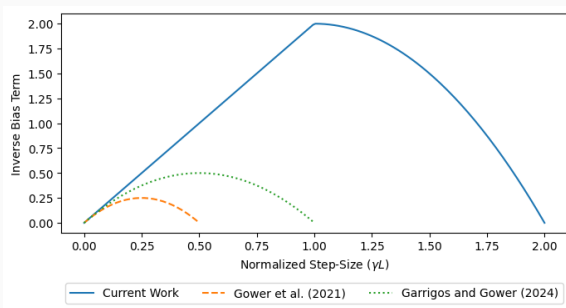
**Observation 2:** No uniform bound in  $T$  for  $\gamma L > 1$ .

**Observation 3:** If  $\sigma_*^2 = 0$ , these are not problems.

# Comparison to State-of-the-Art

Comparison to

- Gower et al. (2021)<sup>2</sup>,
- Garrigos and Gower (2024)<sup>3</sup>.



<sup>2</sup>Gower, Sebbouh, and Loizou, "SGD for Structured Nonconvex Functions: Learning Rates, Minibatching and Interpolation", 2021.

<sup>3</sup>Garrigos and Gower, *Handbook of Convergence Theorems for (Stochastic) Gradient Methods*, 2024.

## Proof Strategy 1/2

Our proofs are based on a Lyapunov analysis with an energy of the form

$$E_t := a_t \|x_t - x_*\|^2 + \rho \sum_{s=0}^{t-1} [f(x_s) - \min f] - \sum_{s=0}^{t-1} e_s \sigma_*^2,$$

where  $(a_t), (e_t), \rho \geq 0$ .

If we can prove a decrease in energy, namely  $\mathbb{E}[E_{t+1}] \leq \mathbb{E}[E_t]$ , then;

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(x_t) - \min f] \leq \frac{a_0}{\rho T} \cdot \|x_0 - x_*\|^2 + \frac{1}{\rho T} \sum_{t=0}^{T-1} e_t \sigma_*^2.$$

We aim at solving

$$\text{Bias}_{\text{opt}}(T) = \inf \left\{ \frac{a_0}{\rho T} : (a_t), (e_t), \rho \text{ are Lyapunov parameters} \right\}.$$

$$\inf_{(a_t), (e_t), \rho} \left\{ \frac{a_0}{\rho T} : \mathbb{E}[E_{t+1}] \leq \mathbb{E}[E_t], \text{ for all convex smooth functions} \right\}$$

$$\inf_{(a_t), (e_t), \rho} \left\{ \frac{a_0}{\rho T} : \mathbb{E}[E_{t+1}] \leq \mathbb{E}[E_t], \text{ for all convex smooth functions} \right\}$$

- Using standard tools from the *Performance Estimation Problem* methodology,<sup>456</sup> we transform the problem into a finite-dimensional optimization problem.
- This problem may be solved numerically.
- The dual problem of the equivalent SDP provides dual variables that help us inspire the proof.

---

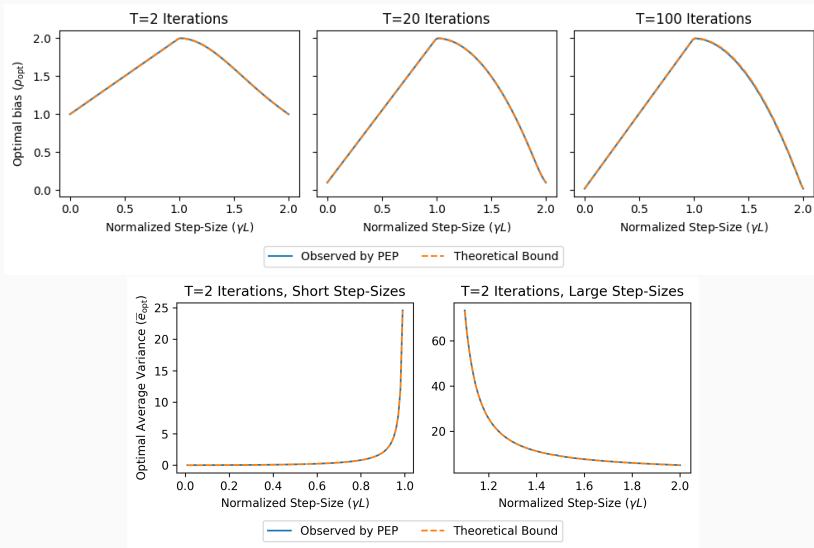
<sup>4</sup>Drori and Teboulle, "Performance of first-order methods for smooth convex minimization: a novel approach", 2014.

<sup>5</sup>Taylor, Hendrickx, and Glineur, "Smooth Strongly Convex Interpolation and Exact Worst-case Performance of First-order Methods", 2017.

<sup>6</sup>Taylor and Bach, "Stochastic first-order methods: non-asymptotic and computer-aided analyses via potential functions", 2019.



# Tightness of Bias and Variance Terms



**Note:** We only claim tightness within our framework.

# Results in Strongly Convex Setting

We obtain a bound of the form

$$\mathbb{E}[\|x_T - x_*\|^2] \leq \text{Bias}(T) \cdot \|x_0 - x_*\|^2 + \text{Variance}(T) \cdot \sigma_*^2,$$

where

$$\text{Bias}(T) = \max\{1 - \gamma\mu, \gamma L - 1\}^{2T},$$

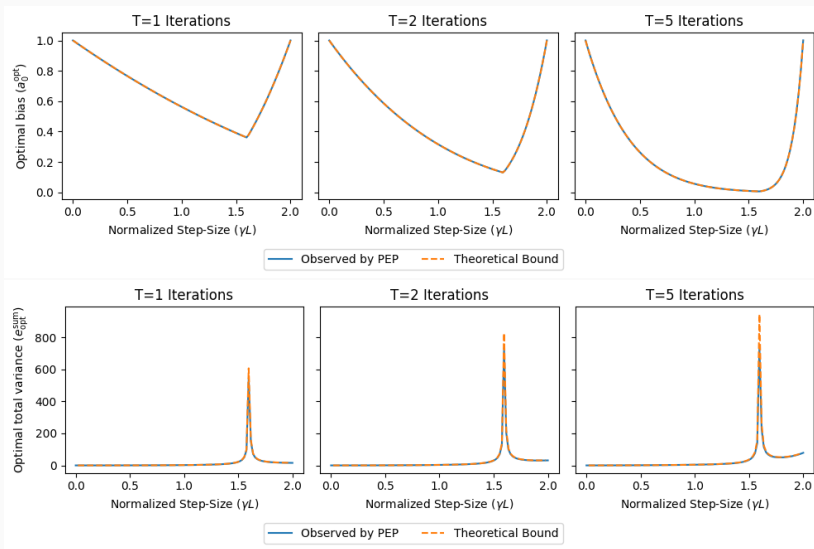
and

$$\text{Variance}(T) = \mathcal{O} \left( \gamma^2 + \frac{\gamma^4}{\left| \gamma - \frac{2}{L+\mu} \right|} \right).$$

**Observation 1:** If  $\gamma$  approaches  $\frac{2}{L+\mu}$ , the variance explodes.

**Observation 2:** To get finite variance at  $\gamma = \frac{2}{L+\mu}$ , we need to degrade the rate a little bit.

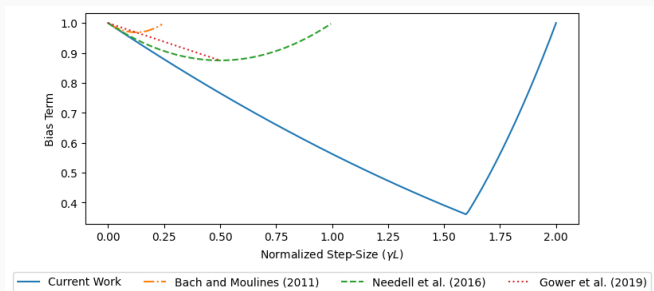
# Tightness in the Strongly Convex Setting



# Results in Strongly Convex Setting

- We obtain improved results for  $\gamma L \in (0, 2)$ ,
- The variance remains bounded for all step-sizes,
- There is a similar singularity at the optimal step-size  $\gamma_{\text{opt}} = \frac{2}{L+\mu}$ .

Comparison to state-of-the-art<sup>789</sup>:



<sup>7</sup>Bach and Moulines, "Non-Asymptotic Analysis of Stochastic Approximation Algorithms for Machine Learning", 2011.

<sup>8</sup>Needell, Srebro, and Ward, "Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm", 2016.

<sup>9</sup>Gower, Loizou, et al., "SGD: General Analysis and Improved Rates", 2019.

# Conclusion

- We provided the first study of SGD without variance assumptions for  $\gamma L \in (0, 2)$ , for convex and strongly convex functions, and improved the current results.
- There is a previously unobserved singularity at optimal step-sizes.
- Our proofs are computer-inspired and numerically shown to be tight within our Lyapunov framework.

**Based on:** Daniel Cortild, Lucas Ketels, Juan Peypouquet, and Guillaume Garrigos. **New Tight Bounds for SGD without Variance Assumption: A Computer-Aided Lyapunov Analysis.** arXiv preprint arXiv:2505.17965. May 2025

# Thank you!