

# Optimizing $(L_0, L_1)$ -Smooth Functions by Gradient Methods

Anton Rodomanov (CISPA)

(joint work with D. Vankov, A. Nedich, L. Sankar, and S. Stich)

5 November 2024

Research Seminar at Université Grenoble Alpes  
Grenoble, France

# Outline

- 1 Motivation
- 2  $(L_0, L_1)$ -Smooth Functions
- 3 Gradient Method
- 4 Other Algorithms
- 5 Experiments
- 6 Conclusions

# Motivation

# Classical Theory for Gradient Method

**Optimization problem:**  $f^* := \min_{x \in \mathbb{R}^d} f(x)$ , where  $f$  is **smooth**.

**Gradient Method (GM):**

$$x_{k+1} = x_k - \eta \nabla f(x_k), \quad k \geq 0.$$

The standard assumption for analyzing GM is that  $f$  is **Lipschitz-smooth**:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d,$$

which is equivalent to the boundedness of the second derivative:

$$\boxed{\|\nabla^2 f(x)\| \leq L}, \quad \forall x \in \mathbb{R}^d.$$

Under this assumption, the theory suggests choosing the stepsize

$$\eta = \frac{1}{L}$$

which ensures the good convergence rate of the method.

# Are All Smooth Functions Lipschitz-Smooth?

No, many smooth functions arising in applications are not Lipschitz-smooth. . .

For example,  $f(x) = |x|^p$  for  $p > 2$  or  $f(x) = e^x$ .

How do we solve optimization problems involving such functions?

## Relative Smoothness [Bauschke et al. 2017; Lu et al. 2018]

Instead of Lipschitz-smoothness, we can consider **relative smoothness**:

$$\boxed{\nabla^2 f(x) \preceq L \nabla^2 \rho(x)}, \quad x \in \mathbb{R}^d,$$

where  $\rho$  is a certain convex “reference function”.

Then, we can apply the **Bregman GM / Mirror Descent**:

$$x_{k+1} = \operatorname{argmin}_{x \in \mathbb{R}^d} \{f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + L \beta_\rho(x_k, x)\},$$

where  $\beta_\rho(x, y) := \rho(y) - \rho(x) - \langle \nabla \rho(x), y - x \rangle$  is the **Bregman distance** generated by  $\rho$ .

**Example:**  $f(x) = \frac{1}{4} \|Ax - b\|^4 + \frac{1}{2} \|Cx - d\|^2$  is smooth relative to  $\rho(x) = \frac{1}{4} \|x\|^4 + \frac{1}{2} \|x\|^2$ .

This is a very powerful technique but requires **fixing the reference function  $\rho$  in advance**.

## $(L_0, L_1)$ -Smooth Functions [Zhang et al. 2020]

In this work, we concentrate instead on another interesting smoothness assumption referred to as  $(L_0, L_1)$ -smoothness:

$$\boxed{\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|,} \quad \forall x \in \mathbb{R}^d.$$

**Original motivation:** Empirical study of loss functions in Neural Networks for Natural Language Processing (NLP) problems.

**NB:**  $f$  is  $L$ -smooth  $\iff f$  is  $(L, 0)$ -smooth.

**Basic example:** Any polynomial  $f(x) = \sum_{i=0}^d a_i x^i$  ( $a_i \in \mathbb{R}$ ) of degree  $d \geq 3$  is  $(L_0, L_1)$ -smooth but not Lipschitz-smooth.

Indeed,  $f'(x) = \sum_{i=1}^d i a_i x^{i-1}$ ,  $f''(x) = \sum_{i=2}^d i(i-1) a_i x^{i-2}$ . Therefore  $\frac{|f''(x)|}{|f'(x)|} \rightarrow 0$  as  $|x| \rightarrow \infty$ , while  $|f''(x)|$  is bounded on any compact interval.

# Clipped Gradient Method

A popular algorithm that provably works for  $(L_0, L_1)$ -smooth functions is the **Clipped GM**:

$$x_{k+1} = x_k - \eta_k \nabla f(x_k), \quad \eta_k = \min \left\{ \eta, \frac{\gamma}{\|\nabla f(x_k)\|} \right\},$$

where  $\eta = \Theta(\frac{1}{L_0})$  and  $\gamma = \Theta(\frac{1}{L_1})$ .

- [Zhang et al. 2020] showed that, to find an  $\epsilon$ -stationary point ( $\|\nabla f(\bar{x})\| \leq \epsilon$ ), Clipped GM needs at most  $O(\frac{L_0 F_0}{\epsilon^2} + \frac{L_1^2 F_0}{L_0})$  gradient computations, where  $F_0 := f(x_0) - f^*$ .
- [Koloskova et al. 2023] further improved it up to  $O(\frac{L_0 F_0}{\epsilon^2} + \frac{L_1 F_0}{\epsilon})$ .

**NB:** Standard GM for  $L$ -smooth functions has complexity of  $O(\frac{L F_0}{\epsilon^2})$ .



# Motivation for This Work

- Further study of  $(L_0, L_1)$ -class: **main inequalities and properties**.
- Why does Clipped GM work for this class? How “natural” is this method and is there any good **interpretation** for it?
- What is the efficiency of gradient methods when our problem is additionally **convex**?

## $(L_0, L_1)$ -Smooth Functions

# Basic Examples

Recall the definition:  $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$ .

## Examples:

- ① (exponent)  $f(x) = e^x$  is  $(L_0, L_1)$ -smooth with  $L_0 = 0$  and  $L_1 = 1$ .
- ② (logistic function)  $f(x) = \ln(1 + e^x)$  is  $(L_0, L_1)$ -smooth with arbitrary  $L_1 \in [0, 1]$  and  $L_0 = \frac{1}{4}(1 - L_1)^2$ .
- ③ (power of Euclidean norm)  $f(x) = \frac{1}{p}\|x\|^p$ , where  $p > 2$ , is  $(L_0, L_1)$ -smooth with arbitrary  $L_1 > 0$  and  $L_0 = (\frac{p-2}{L_1})^{p-2}$ .

**NB:** For the same function, the choice of  $(L_0, L_1)$  may not be unique.

# Calculus of $(L_0, L_1)$ -Smooth Functions

In general, the class is not closed under summation or affine substitution of the arguments. Nevertheless, the class is still closed under some operations.

- 1 If  $f_i$  is  $(L_{0,i}, L_{1,i})$ -smooth for each  $1 \leq i \leq n$ , then  $f(x) = \sum_{i=1}^n f_i(x_i)$ , where  $x \equiv (x_1, \dots, x_n)$ , is  $(L_0, L_1)$ -smooth with  $L_0 = \max_{1 \leq i \leq n} L_{0,i}$  and  $L_1 = \max_{1 \leq i \leq n} L_{1,i}$ .
- 2 If  $f$  is  $(L_0, L_1)$ -smooth and  $g$  is  $L$ -smooth and  $M$ -Lipschitz, then  $f + g$  is  $(L'_0, L'_1)$ -smooth with  $L'_0 = L_0 + ML_1 + L$  and  $L'_1 = L_1$ .
- 3 If  $h(x) = f(\langle a, x \rangle + b)$  and  $f$  is  $(L_0, L_1)$ -smooth, then  $h$  is  $(L'_0, L'_1)$ -smooth with  $L'_0 = \|a\|^2 L_0$  and  $L'_1 = \|a\| L_1$ .

# Main Inequalities

**Theorem.** Function  $f$  is  $(L_0, L_1)$ -smooth iff any of the following inequalities holds for any  $x, y \in \mathbb{R}^d$ :

$$\|\nabla f(y) - \nabla f(x)\| \leq (L_0 + L_1 \|\nabla f(x)\|) \frac{e^{L_1 \|y-x\|} - 1}{L_1},$$

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq (L_0 + L_1 \|\nabla f(x)\|) \frac{\phi(L_1 \|y - x\|)}{L_1^2},$$

where  $\phi(t) := e^t - t - 1$ .

# Lower Bound for Convex Functions

**Theorem.** Let  $f$  be a convex  $(L_0, L_1)$ -smooth function. Then, for any  $x, y \in \mathbb{R}^d$ , we have

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L_0 + L_1 \|\nabla f(y)\|}{L_1^2} \phi_* \left( \frac{L_1 \|\nabla f(y) - \nabla f(x)\|}{L_0 + L_1 \|\nabla f(y)\|} \right),$$

where  $\phi_*(\gamma) = (1 + \gamma) \ln(1 + \gamma) - \gamma$  ( $\geq \frac{\gamma^2}{2+\gamma}$ ) is conjugate to  $\phi$ .

**Corollary:**

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\|\nabla f(y) - \nabla f(x)\|^2}{2(L_0 + L_1 \|\nabla f(y)\|) + L_1 \|\nabla f(y) - \nabla f(x)\|}.$$

# Gradient Method

# Minimizing Upper Bound

**Natural idea:** Minimize the upper bound on the objective:

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + (L_0 + L_1 \|\nabla f(x)\|) \frac{\phi(L_1 \|y - x\|)}{L_1^2},$$

where  $\phi(t) = e^t - t - 1$ .

The optimal point  $y^* = T(x)$  is the result of the [gradient step](#):

$$T(x) = x - r^* \frac{\nabla f(x)}{\|\nabla f(x)\|}, \quad r^* = \frac{1}{L_1} \ln \left( 1 + \frac{L_1 \|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|} \right),$$

resulting in the following bound on improving the function value:

$$\begin{aligned} f(x) - f(T(x)) &\geq \max_{r \geq 0} \left\{ \|\nabla f(x)\| r - \frac{L_0 + L_1 \|\nabla f(x)\|}{L_1^2} \phi(L_1 r) \right\} \\ &= \frac{L_0 + L_1 \|\nabla f(x)\|}{L_1^2} \phi_* \left( \frac{L_1 \|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|} \right). \end{aligned}$$



## Optimal Stepsize

Thus, the point  $y^*$  minimizing the upper bound on the objective is the result of the gradient step

$$T(x) = x - \eta^* \nabla f(x),$$

where the **optimal stepsize** is given by

$$\eta^* = \frac{1}{L_1 \|\nabla f(x)\|} \ln \left( 1 + \frac{L_1 \|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|} \right).$$

The corresponding progress in decreasing the objective is

$$f(x) - f(T(x)) \geq \frac{L_0 + L_1 \|\nabla f(x)\|}{L_1^2} \phi_* \left( \frac{L_1 \|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|} \right) =: \Delta(x).$$

## Simplified Stepsize

The function  $\phi_*$  satisfies  $\frac{\gamma^2}{2+\gamma} \leq \phi_*(\gamma) \leq \frac{\gamma^2}{2}$ .

From this estimate, it follows that  $\Delta(x) \sim \frac{\|\nabla f(x)\|^2}{L_0 + L_1 \|\nabla f(x)\|}$ . More precisely:

$$\frac{\|\nabla f(x)\|^2}{2L_0 + 3L_1 \|\nabla f(x)\|} \leq \Delta(x) \leq \frac{\|\nabla f(x)\|^2}{2(L_0 + L_1 \|\nabla f(x)\|)}.$$

Thus, the guarantee for the optimal stepsize can be simplified:

$$f(x) - f(T(x)) \geq \frac{\|\nabla f(x)\|^2}{2L_0 + 3L_1 \|\nabla f(x)\|}. \quad (*)$$

We can obtain the same guarantee by using the [simplified stepsize](#)

$$\eta = \frac{1}{L_0 + \frac{3}{2}L_1 \|\nabla f(x)\|}.$$

With this stepsize, we still have the same guarantee (\*).

## Clipping Step Size

Note that our simplified stepsize is essentially the clipping stepsize:

$$\eta \sim \frac{1}{L_0 + L_1 \|\nabla f(x)\|} \sim \frac{1}{\max\{L_0, L_1 \|\nabla f(x)\|\}} = \min\left\{\frac{1}{L_0}, \frac{1}{L_1 \|\nabla f(x)\|}\right\}.$$

For the **clipping stepsize**

$$\eta_{\text{cl}} = \min\left\{\frac{1}{2L_0}, \frac{1}{3L_1 \|\nabla f(x)\|}\right\},$$

we can show a similar bound on the function progress as before:

$$f(x) - f(T(x)) \geq \frac{\|\nabla f(x)\|^2}{2(2L_0 + 3L_1 \|\nabla f(x)\|)}.$$

# Various Stepsize Choices: Summary

We have shown that the gradient step

$$T(x) = x - \eta(x) \nabla f(x)$$

is a natural operation minimizing the upper bound on the objective.

The following three stepsizes are equivalent (up to absolute constants) in terms of the objective progress:

- ① (Optimal stepsize)  $\eta^*(x) = \frac{1}{L_1 \|\nabla f(x)\|} \ln(1 + \frac{L_1 \|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|})$ .
- ② (Simplified stepsize)  $\eta(x) = \frac{1}{L_0 + \frac{3}{2} L_1 \|\nabla f(x)\|}$ .
- ③ (Clipping stepsize)  $\eta_{\text{cl}}(x) = \min\{\frac{1}{2L_0}, \frac{1}{3L_1 \|\nabla f(x)\|}\}$ .

They all ensure that

$$f(x) - f(T(x)) \geq \frac{\|\nabla f(x)\|^2}{c(2L_0 + 3L_1 \|\nabla f(x)\|)},$$

where  $c = 1$  for the first two choices and  $c = 2$  for the third one.

# GM: Convergence to Stationary Point

Consider now the gradient method:

$$x_{k+1} = x_k - \eta(x_k) \nabla f(x_k), \quad k \geq 0,$$

where  $\eta(\cdot)$  is one of the stepsize formulas considered before.

**Theorem.** For any given  $\epsilon > 0$ , to reach  $\min_{0 \leq i \leq k-1} \|\nabla f(x_i)\| \leq \epsilon$ , it suffices to make the following number of iterations:

$$k \geq \frac{(2c)L_0F_0}{\epsilon^2} + \frac{(3c)L_1F_0}{\epsilon},$$

where  $F_0 = f(x_0) - f^*$ ,  $c = 1$  for the optimal and simplified stepsizes, and  $c = 2$  for the clipping stepsize.

# Convergence to Stationary Point: Proof

According to the main inequality, we have

$$f_k - f_{k+1} \geq \psi(g_k), \quad \psi(g) := \frac{g^2}{c(2L_0 + 3L_1g)},$$

where  $f_k = f(x_k) - f^*$  and  $g_k = \|\nabla f(x_k)\|$ . Note that  $\psi$  is increasing.

Summing up, we get

$$F_0 \geq f_0 - f_k \geq \sum_{i=0}^{k-1} \psi(g_i) \geq k\psi(g_k^*),$$

where  $g_k^* := \min_{0 \leq i \leq k-1} g_i$ .

Hence,

$$g_k^* \leq \psi^{-1}\left(\frac{F_0}{k}\right) \leq \epsilon$$

whenever

$$k \geq \frac{F_0}{\psi(\epsilon)} \equiv F_0 \frac{c(2L_0 + 3L_1\epsilon)}{\epsilon^2} \equiv \frac{(2c)L_0 F_0}{\epsilon^2} + \frac{(3c)L_1 F_0}{\epsilon}.$$

# Efficiency on Convex Functions

Consider the same method but now additionally assume that  $f$  is **convex**.

**Theorem.** For any given  $\epsilon > 0$ , we have  $f(x_k) - f^* \leq \epsilon$  whenever

$$k \geq O\left(\frac{L_0 R^2}{\epsilon} + L_1^2 R^2\right),$$

where  $R := \|x_0 - x^*\|$  is the distance from the initial point to the solution  $x^*$  of our problem. Furthermore, the distance  $\|x_k - x^*\|$  decreases monotonically.

# Efficiency on Convex Functions: Overview of Proof

We consider the method with the simplified stepsize:

$$x_{k+1} = x_k - \eta_k \nabla f(x_k), \quad \eta_k = \frac{1}{L_0 + \frac{3}{2}L_1 g_k},$$

where  $g_k := \|\nabla f(x_k)\|$ . The proof for the other two stepsizes is similar. Denote  $r_k := \|x_k - x^*\|$ . Then,

$$r_{k+1}^2 = r_k^2 - 2\eta_k \beta_k + \eta_k^2 g_k^2,$$

where  $\beta_k := \langle \nabla f(x_k), x_k - x^* \rangle (\geq f(x_k) - f^*)$ .

According to the lower bound (presented before),

$$\beta_k \geq \frac{g_k^2}{2L_0 + 3L_1 g_k} + \frac{g_k^2}{2L_0 + L_1 g_k} \geq \frac{g_k^2}{L_0 + \frac{3}{2}L_1 g_k} \equiv \xi(g_k) = \eta_k g_k^2.$$

Note that  $\xi$  is increasing.

Hence,

$$r_k^2 - r_{k+1}^2 \equiv \eta_k (2\beta_k - \eta_k g_k^2) \geq \eta_k \beta_k = \frac{\beta_k \xi(g_k)}{g_k^2} \geq \frac{\beta_k^2}{[\xi^{-1}(\beta_k)]^2}.$$



## Efficiency on Convex Functions: Overview of Proof – II

Summing up, we get

$$R^2 \geq r_0^2 - r_k^2 \geq \sum_{i=0}^{k-1} \frac{\beta_i^2}{[\xi^{-1}(\beta_i)]^2} \geq k \frac{(\beta_k^*)^2}{[\xi^{-1}(\beta_k^*)]^2},$$

where  $\beta_k^* := \min_{0 \leq i \leq k-1} \beta_i$ .

Hence,

$$\xi^{-1}(\beta_k^*) \geq \frac{\sqrt{k}}{R} \beta_k^*.$$

Applying  $\xi$  on both sides, we get

$$\beta_k^* \geq \xi\left(\frac{\sqrt{k}}{R} \beta_k^*\right) \equiv \frac{(\frac{\sqrt{k}}{R} \beta_k^*)^2}{L_0 + \frac{3}{2} L_1 \frac{\sqrt{k}}{R} \beta_k^*} \equiv \frac{(\beta_k^*)^2}{\frac{L_0 R^2}{k} + \frac{3}{2} \frac{L_1 R}{\sqrt{k}} \beta_k^*}.$$

Thus,

$$\beta_k^* \leq \frac{L_0 R^2}{k(1 - \frac{3}{2} \frac{L_1 R}{\sqrt{k}})} \leq \epsilon$$

whenever  $\frac{3L_1 R}{\sqrt{k}} \leq 1$  and  $\frac{2L_0 R^2}{k} \leq \epsilon$ . Thus,  $k \geq \max\{\frac{2L_0 R^2}{\epsilon}, 9L_1^2 R^2\}$ .

## Other Algorithms

## Normalized Gradient Method

We can also consider the **Normalized Gradient Method (NGM)**:

$$x_{k+1} = x_k - \frac{\beta_k}{\|\nabla f(x_k)\|} \nabla f(x_k), \quad k \geq 0.$$

**Theorem.** Consider NGM run for  $K$  iterations with constant coefficients:

$$\beta_k = \frac{\hat{R}}{\sqrt{K}}, \quad 0 \leq k \leq K-1.$$

Then, for any given  $\epsilon > 0$ , we have  $\min_{0 \leq k \leq K} f(x_k) - f^* \leq \epsilon$  whenever

$$K + 1 \geq \max \left\{ \frac{L_0 \bar{R}^2}{\epsilon}, \frac{4}{9} L_1^2 \bar{R}^2 \right\},$$

where  $\bar{R} := \frac{R^2}{\hat{R}} + \hat{R}$  and  $R := \|x_0 - x^*\|$ .

**NB:** We can also use time-varying coefficients  $\beta_k = \frac{\hat{R}}{\sqrt{k+1}}$ . The complexity is the same up to an extra logarithmic factor.

# Gradient Method with Polyak Stepsize

Another interesting method is **GM with Polyak Stepsize**:

$$x_{k+1} = x_k - \frac{f(x_k) - f^*}{\|\nabla f(x_k)\|^2} \nabla f(x_k), \quad k \geq 0.$$

It also achieves the same complexity (up to absolute constants).

# Acceleration

We also propose an **acceleration procedure** with the complexity of

$$O\left(m\sqrt{\frac{L_0 R^2}{\epsilon}} + L_1^2 R^2\right),$$

where  $m$  is the complexity of “line search”.

## Procedure:

- 1 Run GM to find  $x_0$  such that  $f(x_0) - f^* \leq \frac{L_0}{5L_1^2}$ .
- 2 Run special monotone version of FGM from  $x_0$ .

**Main idea:** On the sublevel set  $x \in \mathcal{F}_0 := \{x : f(x) \leq f(x_0)\}$ , the function  $f$  is essentially standard  $2L_0$ -smooth:

$$\|\nabla f(x)\| \leq \frac{L_0}{L_1} \implies \|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\| \leq 2L_0.$$

# Experiments

# Experiments

We use the following test problem:

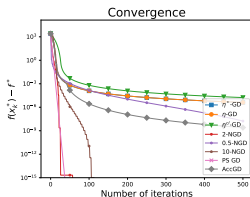
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{p} \|x\|^p \right\}.$$

The initial point  $x_0$  is chosen such that  $\|x_0\| = R$  with  $R = 10$ .

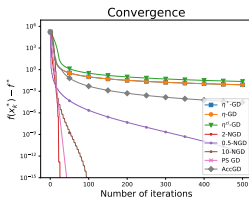
We choose

$$L_1 = 1, \quad L_0 = \left( \frac{p-2}{L_1} \right)^{p-2}.$$

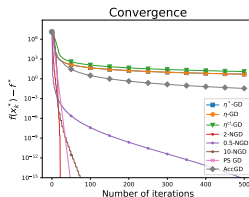
Comparison between different methods:



(a)  $p = 4$



(b)  $p = 6$

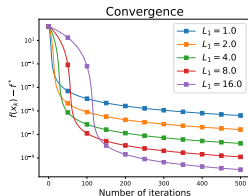


(c)  $p = 8$

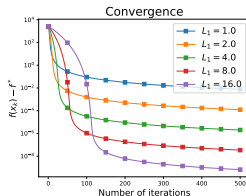
# Experiments – II

Recall that  $L_1 > 0$  can be arbitrary for the same problem.

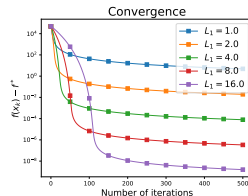
GM with optimal stepsize for different choices of  $L_1$ :



(a)  $p = 4$



(b)  $p = 6$



(c)  $p = 8$



## Conclusions

# Conclusions

- We have seen that GM is a natural method for  $(L_0, L_1)$ -smooth functions, obtained by minimizing the upper bound on the objective.
- The clipping stepsize is a simplification of the corresponding optimal stepsize ensuring the same bound on the function progress.
- In the convex case, we have obtained complexities of  $O(\frac{L_0 R^2}{\epsilon} + L_1^2 R^2)$  and  $O(m\sqrt{\frac{L_0 R^2}{\epsilon}} + L_1^2 R^2)$  for the basic and accelerated method, respectively.

## Open questions:

- Lower bounds?
- Alternative smoothness assumptions?

Paper ([arXiv:2410.10800](https://arxiv.org/abs/2410.10800))

## Optimizing $(L_0, L_1)$ -Smooth Functions by Gradient Methods

D. Vankov, A. Rodomanov, A. Nedich, L. Sankar, S. Stich



Thank you!

# References I



H. H. Bauschke, J. Bolte, and M. Teboulle. A Descent Lemma Beyond Lipschitz Gradient Continuity: First-Order Methods Revisited and Applications. **Mathematics of Operations Research**, 42(2):330–348, 2017.



A. Koloskova, H. Hendrikx, and S. U. Stich. Revisiting Gradient Clipping: Stochastic Bias and Tight Convergence Guarantees. In **International Conference on Machine Learning**, pages 17343–17363. PMLR, 2023.



H. Lu, R. M. Freund, and Y. Nesterov. Relatively Smooth Convex Optimization by First-Order Methods, and Applications. **SIAM Journal on Optimization**, 28(1):333–354, 2018.



J. Zhang, T. He, S. Sra, and A. Jadbabaie. Why Gradient Clipping Accelerates Training: A Theoretical Justification for Adaptivity. In **8th International Conference on Learning Representations**, 2020.