
Breaking the Heavy-Tailed Noise Barrier in Stochastic Optimization Problems

Nikita Puchkin
HSE University,
IITP RAS

Eduard Gorbunov
MBZUAI

Nikolay Kutuzov
MIPT

Alexander Gasnikov
University Innopolis,
MIPT, ISP RAS

Abstract

We consider stochastic optimization problems with heavy-tailed noise with structured density. For such problems, we show that it is possible to get faster rates of convergence than $\mathcal{O}(K^{-2(\alpha-1)/\alpha})$, when the stochastic gradients have finite moments of order $\alpha \in (1, 2]$. In particular, our analysis allows the noise norm to have an unbounded expectation. To achieve these results, we stabilize stochastic gradients, using smoothed medians of means. We prove that the resulting estimates have negligible bias and controllable variance. This allows us to carefully incorporate them into clipped-SGD and clipped-SSTM and derive new high-probability complexity bounds in the considered setup.

1 INTRODUCTION

Stochastic optimization problems with heavy-tailed noise have been gaining a lot of attention in the machine learning community. This phenomenon can be partially explained due to the growing popularity of large language models (Brown et al., 2020; OpenAI, 2023) where stochastic gradients are often far from being well-concentrated (Zhang et al., 2020). In theoretical studies, such behaviour is reflected in the so-called bounded (central) α -th moment assumption with $\alpha \in (1, 2]$ (Nemirovskij and Yudin, 1983; Zhang et al., 2020), written as

$$\mathbb{E}[\|g(x) - \nabla f(x)\|^\alpha] \leq \sigma^\alpha, \quad (1)$$

where $\nabla f(x)$ is the gradient of the objective function $f(x)$, $g(x)$ is the stochastic gradient, and $\sigma \geq 0$. While

Proceedings of the 27th International Conference on Artificial Intelligence and Statistics (AISTATS) 2024, Valencia, Spain. PMLR: Volume 238. Copyright 2024 by the author(s).

in the classical literature on stochastic optimization the authors usually require the noise to have bounded variance (see, for instance, (Nemirovski et al., 2009; Ghadimi and Lan, 2013)), many recent results on high-probability/in-expectation rates of convergence were obtained under a strictly weaker condition $1 < \alpha < 2$.

In deep learning and machine learning communities, one of the most popular techniques to deal with heavy-tailed noise is gradient clipping. In (Pascanu et al., 2013), the authors showed that such a simple trick helps to stabilize neural network training via stochastic gradient descent. Their algorithm called clipped-SGD and clipping in general were then studied in a series of papers including (Abadi et al., 2016; Zhang et al., 2019; Chen et al., 2020; Zhang et al., 2020; Mai and Johansson, 2021; Karimireddy et al., 2021). In particular, for strongly convex functions Zhang et al. (2020) showed that the expected error of clipped-SGD (Pascanu et al., 2013) decreases as $\mathcal{O}(K^{-2(\alpha-1)/\alpha})$ when the number of iterations K grows. In (Sadiev et al., 2023), the authors extended this result and proved that a similar bound (up to logarithmic factors) holds with high probability. According to (Zhang et al., 2020, Theorem 5), the rate of convergence $\mathcal{O}(K^{-2(\alpha-1)/\alpha})$ is tight and cannot be improved if no assumptions, except for (1), are made. However, this rate deteriorates when α is close to 1, and, if $\alpha = 1$, the convergence is not even guaranteed. Fortunately, authors usually have to construct quite specific families of discrete distributions to attain the lower bound $\Omega(K^{-2(\alpha-1)/\alpha})$ (see, e. g., (Nemirovskij and Yudin, 1983; Devroye et al., 2016; Zhang et al., 2020; Cherapanamjeri et al., 2022; Vural et al., 2022)). Such an extreme situation is unlikely to hold in practice and we can hope for more optimistic error guarantees. This brings us to a natural question: *is it possible to achieve better rates of convergence in stochastic optimization problems with heavy-tailed noise under refined assumptions on its structure?* In this paper, we give an affirmative answer to this question.

Contribution. We consider a novel stochastic convex optimization setup with smooth (quasi-

strongly/strongly) convex objective and structured noise (see Assumption 2.1 below), going beyond the standard bounded α -th moment condition with $\alpha \in (1, 2]$. We provide new high-probability upper bounds on the error of versions of clipped-SGD and its accelerated variant called clipped-SSTM in smooth (quasi-strongly/strongly) problems, properly tailored to our setting. In particular, we do not assume boundedness of α -th moments and get $\tilde{\mathcal{O}}(K^{-1/2})$ bound, which outperforms $\mathcal{O}(K^{-2(\alpha-1)/\alpha})$ for $\alpha < 4/3$. Moreover, for symmetric noise distributions, we obtain rates of convergence, which match (up to logarithmic factors) the state-of-the-art ones derived under the bounded variance assumption (Nazin et al., 2019; Davis et al., 2021; Gorbunov et al., 2020). In particular, for smooth strongly convex problems, the dominating term in our upper bound decreases as $\tilde{\mathcal{O}}(K^{-1})$. Our approach relies on new non-asymptotic results on the performance of smoothed median of means.

Paper structure. The rest of the paper is organized as follows. In Section 2, we introduce our notation and formulate problem setup. Section 3 is devoted to an overview of related work. In Sections 4 and 5, we present our main results and illustrate the performance of suggested algorithms in Section 6. Many technical details are deferred to Appendix.

2 SETUP AND NOTATION

Before we formulate our main contributions, we need to introduce the notation and formalize the problem and setup we focus on.

Notation. Throughout the paper, we denote the standard Euclidean norm in \mathbb{R}^d as $\|\cdot\|$. To simplify the bounds in the main text, we use $\tilde{\mathcal{O}}(\cdot)$ notation that hides constant and polylogarithmic factors. For any $x_1, \dots, x_n \in \mathbb{R}^d$, we denote $\text{Mean}(x_1, \dots, x_n) = (x_1 + \dots + x_n)/n$. For any random vectors ξ_1, \dots, ξ_{2m+1} , $\text{Med}(\xi_1, \dots, \xi_{2m+1})$ stands for the $(m+1)$ -th order statistic (also called the median), taken in the component-wise fashion. For any non-zero $x \in \mathbb{R}^d$ and $\lambda > 0$, $\text{clip}(x, \lambda) = \min\{1, \lambda/\|x\|\}x$ denotes the clipping operator. We also define $\text{clip}(0, \lambda) = 0$ for all $\lambda > 0$. For any $\theta > 0$,

$$\Phi_\theta(t) = \frac{1}{\sqrt{2\pi\theta}} \int_1^t e^{-u^2/(2\theta^2)} du$$

stands for the CDF of a Gaussian random variable with zero mean and variance θ^2 . Sometimes, we use the notation $a \wedge b$ and $a \vee b$, instead of $\min\{a, b\}$ and $\max\{a, b\}$, respectively. Along with the standard $\mathcal{O}(\cdot)$ notation, we use the relations $g \lesssim h$ and $h \gtrsim g$, which

are equivalent to $g = \mathcal{O}(h)$. Finally, for any functions $g : \mathbb{R}^d \rightarrow \mathbb{R}$ and $h : \mathbb{R}^d \rightarrow \mathbb{R}$, their convolution is denoted as $g * h(x) = \int_{\mathbb{R}^d} g(x-y)h(y) dy$. We also adopt the notation $g^{*k}(x) = \underbrace{g * \dots * g}_{k \text{ times}}(x)$.

Setup. We consider an unconstrained smooth convex optimization problem

$$\min_{x \in \mathbb{R}^d} f(x), \quad (2)$$

where the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is accessible through the stochastic first-order oracle $\mathcal{G} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that for given point $x \in \mathbb{R}^d$ returns some estimate of $\nabla f_\xi(x)$ of the true gradient $\nabla f(x)$. We make the following assumption about the distribution of $\nu = \nabla f_\xi(x) - \nabla f(x)$.

Assumption 2.1. For any $x \in \mathbb{R}^d$ and each $j \in \{1, \dots, d\}$, the marginal density \mathfrak{p}_j of the j -th component of the noise $\nu = \nabla f_\xi(x) - \nabla f(x)$ satisfies the following conditions:

- there exists $M_j > 0$, such that $\mathfrak{r}_j(u) = (\mathfrak{p}_j(u) - \mathfrak{p}_j(-u))/2$ fulfils

$$\int_1^{+7} u \mathfrak{r}_j(u) du = 0 \quad \text{and} \quad \int_1^{+7} u^2 |\mathfrak{r}_j(u)| du \leq M_j.$$

- there are $B_j > 0$ and $\beta_j \geq 1$, such that, for any $k \in \mathbb{N}$,

$$s_j^k(u) \leq \frac{B_j k}{k^{(\beta_j+1)/\beta_j} + |u|^{1+\beta_j}},$$

where $s_j(x) = (\mathfrak{p}_j(x) + \mathfrak{p}_j(-x))/2$.

In Assumption 2.1, we split marginal densities of noise components into a sum of symmetric and antisymmetric parts. For each $j \in \{1, \dots, d\}$, the antisymmetric remainder $\mathfrak{r}_j(u)$ is a signed density with a finite second moment. However, the symmetric term $s_j(u)$ may decay much slower, than $\mathfrak{r}_j(u)$. As a result, the density $\mathfrak{p}_j(u)$ has finite moments up to order $\alpha < (\beta_j \wedge 2)$. Note that if $\beta_j = 1$ for some $j \in \{1, \dots, d\}$, then the noise norm $\|\nu\|$ may have no expectation.

We proceed with several examples of $s_j(u)$, satisfying Assumption 2.1. Obviously, if the j -th component of $\nu(x)$ has a standard Cauchy distribution, that is, $s_j(u) = 1/\pi \cdot 1/(1+u^2)$, then, for any $k \in \mathbb{N}$, $s_j^k(u) = 1/\pi \cdot k/(k^2+u^2)$, and Assumption 2.1 is fulfilled with $B_j = 1/\pi$ and $\beta_j = 1$. This is a particular example of a symmetric α -stable distribution with parameter $\alpha = 1$. All symmetric α -stable distributions have a characteristic function of the form $\varphi(y) = e^{-|y|/\sigma_j^\alpha}$,

where $\epsilon > 0$ and $\beta \in (0; 2]$ (see, e.g., (Feller, 1971, Chapter XVII, Sections 5-6)). If $\beta \in [1; 2]$, they also satisfy Assumption 2.1 with $\beta_j = \beta$ and some $B_j > 0$. In general, if $\beta_j(u) = B_j u^{1+\beta_j}$, then it is known from probability theory that $\beta_j^k(u) = B_k u^{1+\beta_j}$ for any $k \geq 1$. Assumption 2.1 can be viewed as a non-asymptotic version of this property.

We also make several standard assumptions about function f itself. Similarly to (Gorbunov et al., 2021; Sadiev et al., 2023), it is sufficient for our analysis to make all the assumptions only on some compact subset (ball) of \mathbb{R}^d since we show that with high probability the considered methods do not leave this compact. We start with the standard smoothness assumption.

Assumption 2.2. There exists a set $Q \subset \mathbb{R}^d$ and constant $L > 0$ such that for all $x, y \in Q$

$$\| \nabla f(x) - \nabla f(y) \| \leq L \| x - y \|; \quad (3)$$

$$\| \nabla f(x) \|^2 \leq 2L (f(x) - f_*); \quad (4)$$

where $f_* = \inf_{x \in Q} f(x) > -\infty$.

When $Q = \mathbb{R}^d$, (4) follows from (3). However, when $Q \subset \mathbb{R}^d$, condition (4) can be derived from (3), if the latter is assumed on a slightly larger set, see (Sadiev et al., 2023, Appendix B) for additional discussion.

We assume convexity or strong convexity of f for the results with accelerated rates.

Assumption 2.3. There exists set $Q \subset \mathbb{R}^d$ and constant $\mu > 0$ such that f is μ -strongly convex, i.e., for all $x, y \in Q$, it holds that

$$f(y) > f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \| y - x \|^2; \quad (5)$$

Finally, for non-accelerated case, it is sufficient to assume a relaxed condition called quasi-strong convexity.

Assumption 2.4. There exists set $Q \subset \mathbb{R}^d$ and constant $\mu > 0$ such that f is μ -quasi-strongly convex, that is, for all $x \in Q$ and $x_* = \arg \min_{x \in \mathbb{R}^d} f(x)$

$$f(x) > f(x_*) + \langle \nabla f(x_*), x - x_* \rangle + \frac{\mu}{2} \| x - x_* \|^2; \quad (6)$$

The above assumption belongs to the class of conditions on structured non-convexity. When $\mu > 0$ it (together with smoothness) implies linear convergence for Gradient Descent (Necoara et al., 2019).

3 RELATED WORK

High-probability complexity bounds. Under sub-Gaussian noise assumption, optimal (up to logarithmic

factors) high-probability complexity bounds¹ are proven by Nemirovski et al. (2009) for (strongly) convex non-smooth problems with bounded sub-gradients, by Ghadimi and Lan (2012) for (strongly) convex smooth problems, and by Li and Orabona (2020) for smooth non-convex problems. These results are achieved for the same methods that are optimal in terms of the in-expectation convergence. However, when the noise has just a finite variance, some algorithmic changes seem to be necessary, e.g., as it is shown in (Sadiev et al., 2023, Section 2), standard SGD has provably bad (inverse-power instead of $\text{poly}(\log(1/\epsilon))$) dependence on the confidence level ϵ in this case.

A popular tool for overcoming this issue is gradient clipping, i.e., the application of the clipping operator to the gradient estimator. A version of gradient clipping is used by Nazin et al. (2019) who derive the first (non-accelerated) high-probability complexity bounds for smooth (strongly) convex problems on compact domains with logarithmic dependence on ϵ under bounded variance assumption. Accelerated results are obtained by Davis et al. (2021) and Gorbunov et al. (2020) for smooth strongly convex and smooth convex problems respectively. Gorbunov et al. (2021) generalize these results to the case of problems with Hölder continuous gradients.

State-of-the-art high-probability complexity bounds are derived under bounded (central) β -th moment assumption (1). The first work in this direction is (Cutkosky and Mehta, 2021) where the authors derived optimal (up to logarithmic factors) bounds in the smooth non-convex regime with the additional assumption of boundedness of the gradients. Without this assumption, a worse bound is derived by Sadiev et al. (2023), and the optimal one is obtained by Nguyen et al. (2023b). For (strongly) convex problems the results for clipped-SGD and its accelerated version clipped-SSTM (Gorbunov et al., 2020) are derived by Sadiev et al. (2023). Up to logarithmic factors, these results match the known lower bounds in the strongly convex case (Zhang et al., 2020). Nguyen et al. (2023a) improve the logarithmic factors in the upper bounds from Sadiev et al. (2023); Nguyen et al. (2023b). Recently, the generalization of the results from (Sadiev et al., 2023) to the case of composite and distributed optimization were obtained by Gorbunov et al. (2023).

Other results under heavy-tailed noise. Although in our work we primarily focus on high-probability convergence results, we briefly discuss here

¹Such results establish upper bounds for the number of oracle calls needed for a method to find point x such that $f(x) - f(x_*) \leq \epsilon$ or $\| \nabla f(x) \|^2 \leq \epsilon$ or $\| \nabla f(x) \|^2 \leq \epsilon$ are less than n with probability at least $1 - \delta$, where x_* is a solution of (2).

other existing works devoted to the convergence of stochastic methods under heavy-tailed noise assumption. For convex functions with bounded gradients, Nemirovskij and Yudin (1983) show $O(K^{-(1-\epsilon)})$ in-expectation convergence rate for Mirror Descent and Vural et al. (2022) propose an extension of this result for uniformly convex functions. For strongly convex functions with bounded gradients, Zhang et al. (2020) show $O(K^{-2(1-\epsilon)})$ in-expectation rate of convergence. In the smooth non-convex case $O(K^{-2(1-\epsilon)/3})$ in-expectation convergence rate is achieved by Zhang et al. (2020) who also derive a matching lower bound.

In the case of the noise with symmetric density function and bounded first moment, Jakovetić et al. (2023) derive $O(K^{-1})$ in-expectation convergence rate for β -strongly convex L -smooth functions and SGD-type methods with general non-linearities. However, parameter β is proportional to ϵ^{-p} in the worst case. Therefore, this rate can be much slower than $O(K^{-2(1-\epsilon)})$ for ill-conditioned/large-scale problems though Jakovetić et al. (2023) do not assume (1) and consider general class of non-linearities. Our analysis also does not rely on (1), but we additionally allow non-symmetric noise distributions and do not assume the existence of the finite first moment of the noise.

Median estimates. Median, median of means, and smoothed median were extensively used in the problems of robust mean estimation and robust machine learning (see, for instance, (Nemirovskij and Yudin, 1983; Minsker, 2015; Devroye et al., 2016; Lugosi and Mendelson, 2019b, 2020; Lecué and Lerasle, 2020; Cherapanamjeri et al., 2022)). A reader is referred to a comprehensive survey of Lugosi and Mendelson (2019a) on this topic. Usually, the authors use median of means or its modifications to get sub-Gaussian rates of convergence, assuming the existence of only two moments. In Cherapanamjeri et al. (2022), the authors went further and derived a minimax optimal upper bound in the problem of mean estimation when observations have finite moments of order $2(1; 2]$. However, the authors faced the same problem as Zhang et al. (2020): the rate of convergence became very slow when approached 1. This happens, because the family of distributions of interest is extremely large if one assumes the existence of β -th moment only. In our paper, we exploit the special noise structure, described in Section 2. Under Assumption 2.1, we derive new non-asymptotic bounds on the performance of the smoothed median of means, which do not deteriorate even if the underlying density has quite heavy tails.

4 SMOOTHED MEDIAN OF MEANS AND ITS PROPERTIES

In this section, we describe how to get reliable gradient estimates from noisy stochastic gradients given by the first-order oracle. Let us start with a simple example. Fix an arbitrary $x \in \mathbb{R}^d$ and assume that the noise $\epsilon = r f(x) - r f(x) \in \mathbb{R}^d$ has a symmetric absolutely continuous distribution. For any $j \in \{1, \dots, d\}$, let $p_j(u)$ be the marginal density of ϵ_j , the j -th component of ϵ . Then the following proposition holds true.

Proposition 4.1. Fix any $j \in \{1, \dots, d\}$ and assume that the marginal density of ϵ_j is symmetric, that is, $p_j(u) = p_j(-u)$ for all $u \in \mathbb{R}$. Suppose that there exist positive numbers B_j and β_j , such that

$$p_j(u) \leq \frac{B_j}{1 + |u|^{\beta_j + 1}}; \quad \text{for all } u \in \mathbb{R}.$$

Let $\epsilon_j^{(1)}, \dots, \epsilon_j^{(2m+1)}$ be independent copies of ϵ_j . If $m > 3/\beta_j$, then $E \text{Med}(\epsilon_j^{(1)}, \dots, \epsilon_j^{(2m+1)}) = 0$ and $E \text{Med}(\epsilon_j^{(1)}, \dots, \epsilon_j^{(2m+1)})^2$ is finite.

The proof of the proposition with an explicit bound on the variance of $\text{Med}(\epsilon_j^{(1)}, \dots, \epsilon_j^{(2m+1)})$ is postponed to Appendix. Proposition 4.1 shows that, despite the heavy tails of the underlying density p , $m > \max_j 3/\beta_j$: $16j \leq dg$ oracle calls are enough to produce an unbiased estimate $\hat{g}(x) = \text{Med}(r f_1(x); \dots; r f_m(x))$ of $r f(x)$ with a finite variance. After that, we can use the standard clipping technique to solve the optimization problem (2).

Unfortunately, the symmetry assumption, which played the central role in Proposition 4.1, is rather restrictive. To deal with asymmetric distributions, we use more sophisticated gradient estimates, based on smoothed median of means.

Definition 4.2. Let ϵ be a random element in \mathbb{R}^d and let $\beta > 0$ be an arbitrary number. For any positive integers m and n , the smoothed median of means $\text{SMoM}_n(\epsilon; \beta)$ is defined as follows:

$$\text{SMoM}_n(\epsilon; \beta) = \text{Med}(\epsilon^{(1)}, \dots, \epsilon^{(2m+1)});$$

where, for each $j \in \{0, \dots, 2m\}$,

$$\epsilon^{(j)} = \text{Med}(\epsilon^{(j+1)}, \dots, \epsilon^{(j+1+n)}) + \epsilon^{(j+1)};$$

$\epsilon^{(1)}, \dots, \epsilon^{(2m+1)}$ are i.i.d. copies of ϵ , and $\epsilon^{(1)}, \dots, \epsilon^{(2m+1)} \in \mathbb{N}(0; I_d)$ are independent standard Gaussian random vectors.

Let us briefly describe the idea behind our approach. Assuming that $r f(x) - r f(x)$ at a point $x \in \mathbb{R}^d$ has a density $p(u)$, we represent the latter in the following form:

$$p(u) = s(u) + r(u);$$

where $s(u) = (p(u) + p(-u))/2$ is a symmetric part and $r(u) = (p(u) - p(-u))/2$ is an antisymmetric remainder. If the tails of the remainder $r(u)$ are much lighter than the ones of $p(u)$, we can make oracle calls at the point x and take the average of $r f_1(x); \dots; r f_n(x)$. If n is large enough, then the distribution of $\text{Mean}(r f_1(x); \dots; r f_n(x))$ is almost symmetric. Hence, we can use the same trick as in Proposition 4.1 to get an estimate of $r f(x)$ with a finite variance. We add small Gaussian noise to ensure that the density of our estimate is infinitely differentiable, as we need it for technical purposes. Note that, in general, the expectation of $\text{SMoM}_n(r f(x) - r f(x); \cdot)$ is not equal to zero, and it is a challenging task to show that it is sufficiently small.

Before we move to the heavy-tailed setup, let us illustrate the efficiency of our approach in the case when the stochastic gradients have a finite second moment.

Lemma 4.3. Assume that stochastic gradient $r f(x)$ at a point $x \in \mathbb{R}^d$ has an absolutely continuous distribution and a finite second moment $E(r f(x) - r f(x))(r f(x) - r f(x))^T = \Sigma$. Then, for any positive integer m and n , it holds that

$$E \text{SMoM}_n(r f(x); r f(x))^2 \leq \frac{6 \cdot 4(2m+1)}{n} \text{Tr}(\Sigma) + 2^d d$$

If, in addition, $m > 3$ and $n \geq mk$, then

$$E \text{SMoM}_n(r f(x); r f(x)) \leq \frac{m^p}{n} \overline{\text{Tr}(\Sigma^2)}$$

Remark 4.4. Lemma 4.3 also yields that

$$E \text{SMoM}_n(r f(x); r f(x)) \leq \frac{p}{m} \frac{\text{Tr}(\Sigma)}{n} + \frac{p}{d}$$

Though this bound is enough for our purposes, nevertheless, we find it useful to prove a dimension-free $O(1/n)$ upper bound on expectation of the smoothed median of means, which follows from the analysis of impact of antisymmetric density part on the expectation (see Lemma A.1 in Appendix).

Lemma 4.3 shows that, if the noise vector has a finite second moment, then the smoothed median of means has a small controllable shift and bounded variance. In this case, it behaves similarly to clipping. However, if we deal with heavy-tailed noise, the standard clipping technique fails, while the smoothed median of means still has a small bias and finite variance. We proceed with the main result of this section.

Lemma 4.5. Assume that the stochastic gradient $r f(x)$ at a point $x \in \mathbb{R}^d$ has an absolutely continuous distribution. Suppose that, for any $j \in \{1; \dots; d\}$

and any $x \in \mathbb{R}^d$, the density of $r f_j(x) - r f_j(x)$ meets Assumption 2.1. Then, if $m > 2 + 3 = j$ and $2n > (2 - m^2)M_j$ for all $j \in \{1; \dots; d\}$, it holds that

$$E \text{SMoM}_n(r f(x); r f(x))^2 \leq \frac{m(1 + 2^d)}{2n} \sum_{j=1}^d \frac{M_j^2}{n} + \sum_{j=1}^d \frac{2^j B_j}{n^{j-1}} + \sum_{j=1}^d \frac{B_j M_j}{n^j} \quad \#$$

and

$$E \text{SMoM}_n(r f(x); r f(x)) \leq \frac{m(1 + 2^d)}{2n} \sum_{j=1}^d \frac{M_j^2}{n} + \frac{m}{2n} \sum_{j=1}^d \frac{2^j B_j}{n^{j-1}}$$

Remark 4.6. The bounds on $E \text{SMoM}_n(r f(x); r f(x))$ in Lemma 4.3 and Lemma 4.5 rely on the fact that convolution with the Gaussian density is infinitely differentiable. However, if $s_1; \dots; s_d$ are also sufficiently smooth, then one can take $\epsilon = 0$ and apply a similar technique as in Lemma A.1 and Lemma A.4.

The proof of Lemma 4.5 is moved to Appendix. It shows that the bias of the smoothed median of means decays rapidly with the growth of the batch size, though the noise may have extremely heavy tails. The reason for that is the special noise structure guaranteed by Assumption 2.1. This favourable property allows us to obtain faster rates of convergence in heavy-tailed stochastic convex optimization problems.

5 MAIN RESULTS FOR STOCHASTIC OPTIMIZATION

In view of the results of the previous section, Assumption 2.1 allows constructing an estimator with bounded bias and variance using the smoothed median of means. Since in the analysis of the stochastic first-order methods we only use these two properties, we formulate them as a separate assumption for convenience.

Assumption 5.1. There exists $N \in \mathbb{N}$, aggregation rule R and (possibly dependent on N) constants $b > 0$; $\gamma > 0$ such that for an $x \in \mathbb{R}^d$ i.i.d. samples $r f_1(x); \dots; r f_N(x)$ from the oracle $G(x)$ satisfy the following relations:

$$\| \frac{1}{N} \sum_{i=1}^N E[r f_i(x)] - r f(x) \| \leq b; \quad (7)$$

$$E \| r f(x) - E[r f(x)] \|^2 \leq \gamma; \quad (8)$$

where $r f(x) = R(r f_1(x); \dots; r f_N(x))$ and expectations are taken w.r.t. $r f_1(x); \dots; r f_N(x)$.

We emphasize once again that Assumption 5.1 holds whenever Assumption 2.1 is satisfied. Indeed, we can take $r f(x) = \text{SMoM}_n r f(x)$; with parameters m, n , and β , satisfying the conditions of Lemma 4.5. In this case, the batch size is equal to $n = (2m + 1)n$.

5.1 Convergence of clipped-SGD

We start with clipped-SGD defined as follows:

$$x^{k+1} = x^k - \eta \text{clip}(r f_k(x^k); \kappa); \quad (9)$$

where $r f_k(x^k)$ is an estimator satisfying Assumption 5.1 sampled independently from previous iterations. Below we formulate the main convergence result for clipped-SGD in the quasi-convex case.

Theorem 5.2. Let Assumptions 2.2 and 2.4 with $b = 0$ hold on $Q = B_{2R}(x)$, where $R > \kappa x^0$. Suppose that $r f_k(x^k)$ satisfies Assumption 5.1 with parameters $b_k; \kappa$ for $k = 0; 1; \dots; K$, $K > 0$ and $\kappa = (\min_{k=0;1;\dots;K} \frac{1}{L_k}; R = \frac{1}{\kappa A}; R = \frac{1}{\kappa A}; R = \frac{1}{\kappa A}; R = \frac{1}{\kappa A})$; $\kappa = (\frac{1}{R=A})$; where $A = \ln(4(K+1))$ and $b = \max_{k=0;1;\dots;K} b_k$, $\kappa = \max_{k=0;1;\dots;K} \kappa$. Then the iterates produced by clipped-SGD after K iterations with probability at least $1 - \epsilon$ satisfy

$$f(\bar{x}^K) - f(x^*) = \mathcal{O} \left(\max \left\{ \frac{LR^2}{K}; \frac{R}{K}; bR \right\} \right);$$

$$\text{where } \bar{x} = \frac{1}{K+1} \sum_{k=0}^K x^k.$$

The rate of convergence in the above result matches (up to logarithmic factors) the best-known one for clipped-SGD under bounded variance assumption (Gorbunov et al., 2021). Due to systematic bias bounded by the method reaches only $\mathcal{O}(bR)$ error after a sufficiently large number of steps. When the bias is just bounded and cannot be controlled, this situation is standard (Devolder et al., 2014). The proof of the above result follows the ones given in (Gorbunov et al., 2021; Sadiev et al., 2023): using the induction argument, we show that under a proper choice of parameters, the iterates stay in a bounded set with high probability, which allows us to apply standard Bernstein inequality (see Lemma B.1). In particular, this proof technique differs from the standard ones that rely on the boundedness of the noise (Rakhlin et al., 2011) or on the assumption that the noise is sub-Gaussian (Harvey et al., 2019).

However, in our setup, we can control the bias. For example, if the distribution is symmetric and has a bounded moment of the order β for some $\beta > 0$, then according to Proposition 4.1, it is sufficient to use coordinate-wise median estimator to get $r f(x)$ satisfying Assumption 5.1 with $b = 0$ and $\kappa = \frac{1}{d(2m+1)} \frac{1}{\beta} \frac{1}{j^{2\beta}}$ (see (18)) using $\mathcal{O}(1/\beta)$ samples of $r f(x)$. In this case, we have the following result.

Corollary 5.3 (Symmetric noise). Let the assumptions of Theorem 5.2 hold and for all $x \in \mathbb{R}^d$ the noise $\xi = r f(x) - f(x)$ satisfies the conditions from Proposition 4.1. Then the iterates produced after K iterations of clipped-SGD with $r f_k(x^k)$ being a coordinate-wise median of $2m+1$ samples $r f(x^k)$ with $m > \max_{j=1;\dots;d} \frac{1}{\beta} \frac{1}{j^{2\beta}}$ and $\kappa = \frac{1}{d(2m+1)} \frac{1}{\beta} \frac{1}{j^{2\beta}}$ and where $A = \ln(4(K+1))$ with probability at least $1 - \epsilon$ satisfy

$$f(\bar{x}^K) - f(x^*) = \mathcal{O} \left(\max \left\{ \frac{LR^2}{K}; \frac{R}{K} \right\} \right)$$

and the overall number of stochastic oracle calls equals $(2m+1)K = \mathcal{O}(K \max_{j=1;\dots;d} \frac{1}{\beta} \frac{1}{j^{2\beta}})$.

This result implies that as long as the distribution is symmetric, its tails can be even heavier than the ones of Cauchy distribution, i.e., moments of order larger than β for some $\beta \in (0; 1]$ can be unbounded, but clipped-SGD with coordinate-wise median estimator inside will still converge as in the case when the stochastic gradients are unbiased and have bounded variance. We emphasize that the existing state-of-the-art high probability convergence results (Sadiev et al., 2023; Nguyen et al., 2023b,a) have a slower decreasing main term (of the order $\mathcal{O}(K^{-\beta})$) and are derived for much lighter tails. However, in contrast to Corollary 5.3, the mentioned results do not rely on the symmetry.

Finally, we consider the general case, when the noise satisfies Assumption 2.1, i.e., the noise also has a non-symmetric component. In this case, Lemma 4.5 implies that the smoothed median of means gives an estimator $r f(x)$ satisfying Assumption 5.1 with

$$b = \mathcal{O}(C_m); \quad \kappa = \mathcal{O}(d(1 + \frac{1}{\beta}) + D); \quad (10)$$

$$C = \frac{(1 + \frac{1}{\beta})^d}{2} \sum_{j=1}^d \frac{X_j^d}{M_j^2} + \frac{1}{2} \sum_{j=1}^d \frac{X_j^d}{M_j^2} \frac{2^j B_j}{n^{j-1}}; \quad (11)$$

$$D = \sum_{j=1}^d \frac{X_j^d}{n} + \frac{2^j B_j}{j n^{j-1}} + \sum_{j=1}^d \frac{X_j^d}{n^j} \frac{B_j M_j}{n^{j-1}} \frac{2^{j-1}}{j^{2\beta}}; \quad (12)$$

using $\mathcal{O}(n)$ samples $r f(x)$ (when $m = \mathcal{O}(1)$). Together with Theorem 5.2 this implies the following result.

Corollary 5.4 (General noise) Let the assumptions of Theorem 5.2 hold and for all $x \in \mathbb{R}^d$ the noise

$= r f(x) - r f(x)$ satisfies Assumption 2.1. Then clipped-SGD with $r f_k(x^k)$ being the smoothed median of means of $O(n)$ samples $r f(x^k)$ and $k = \dots$

$$\Theta \max \left(\frac{LR^2}{K}; \frac{p}{K}; \frac{(1+\epsilon)CR}{2n} \right);$$

where $n > (\max_{j \in [d]} M_j = 2)$ and C, D are defined in (11)-(12). The overall number of oracle calls equals $O(nK)$.

Due to the heavy-tailedness of the symmetric part of the noise distribution and the presence of bias, the variance term does not necessarily improve with the growth of the number of samples. However, the bias term still can be smaller than any predefined level via increasing n . When the bias is large, i.e., the noise is sufficiently non-symmetric, then one can take $n = \Theta(\epsilon^{-1})$ and $K = \Theta(\epsilon^{-2})$ to guarantee $f(x^K) - f(x) \leq \epsilon$ with probability at least $1 - \epsilon$, i.e., the total oracle complexity is $\Theta(\epsilon^{-3})$. However, when the bias is small, i.e., $\max_{j \in [d]} M_j$ are sufficiently small ², then for $\epsilon = (1 + \epsilon)CR = 2$ one can achieve even $K = \Theta(\epsilon^{-2})$ total oracle complexity that matches (up to logarithmic factors and constants related to the variance) the main term in the optimal complexities under bounded variance assumption (Gorbunov et al., 2020). However, in contrast to the existing results, we do not require the noise to have a finite first moment.

The following theorem gives a convergence rate for clipped-SGD in the quasi-strongly convex case.

Theorem 5.5. Let Assumptions 2.2 and 2.4 with $\epsilon > 0$ hold on $Q = B_{2R}(x)$, where $R > \epsilon^{-1} \|x - x^0\|$. Suppose that $r f_k(x^k)$ satisfies Assumption 5.1 with parameters $b_k; k$ for $k = 0; 1; \dots; K$, $K > 0$ and $k = \dots$

²In practice, this can happen when a non-symmetric noise is added to the stochastic gradients with symmetric noise, e.g., this can happen in some mechanisms for ensuring differential privacy such as the one from (Guo et al., 2023) (the non-symmetric part of the noise in the resulting vector after averaging over multiple clients can have a small variance when the number of clients is large, which is typical for modern Federated Learning applications (Kairouz et al., 2021)).

K iterations with probability at least $1 - \epsilon$ satisfies

$$\|x^K - x^0\|^2 \leq \Theta \max \left(R^2 \exp \left(-\frac{K}{L \ln K} \right); \frac{2}{K}; \frac{bR}{\epsilon} \right);$$

Similarly to the convex case, in the case of symmetric noise with bounded ϵ -th moment for some $\epsilon > 0$, the above result matches the best-known one for clipped-SGD under bounded variance and strong convexity (Gorbunov et al., 2021). In particular, for such noise distributions, condition (1) is not necessarily satisfied, and even if it is satisfied, our rate $\Theta(K^{-1})$ is better than the lower bound $(K^{-2(1-\epsilon)})$ under condition (1) for $\epsilon \in (1/2; 1)$ (Zhang et al., 2020). However, it is worth mentioning that we do rely on the symmetry of the noise distribution to achieve this rate, while the lower bound holds for any distributions satisfying (1).

In the non-symmetric case, we combine Theorem 5.5 with Lemma 4.5 and get the following result.

Corollary 5.6 (General noise) Let the assumptions of Theorem 5.5 hold and for all $x \in \mathbb{R}^d$ the noise $= r f(x) - r f(x)$ satisfies Assumption 2.1. Then the iterates produced after K iterations of clipped-SGD with $r f_k(x^k)$ being the smoothed median of means of $O(n)$ samples $r f(x^k)$ and $k = \dots$

$$\|x^K - x^0\|^2 \leq \Theta \left(R^2 \exp \left(-\frac{K}{L \ln K} \right) + \Theta \max \left(\frac{(1+\epsilon)d + D}{K}; \frac{(1+\epsilon)CR}{2n} \right) \right);$$

where $n > (\max_{j \in [d]} M_j = 2)$ and C, D are defined in (11)-(12). The overall number of oracle calls equals $O(nK)$.

Taking $n = \Theta(\epsilon^{-1})$ and $K = \Theta(\epsilon^{-1})$, one can guarantee $\|x^K - x^0\|^2 \leq \epsilon$ with probability at least $1 - \epsilon$, i.e., the total oracle complexity is $\Theta(\epsilon^{-2})$. This result is worse than the one for the case of symmetric distribution, but it still does not require the existence of the expectation of the noise and is better than $\Theta(\epsilon^{-2(1-\epsilon)})$, which is known to be optimal (up to logarithmic factors) under assumption (1), when $\epsilon < 1/2$.

5.2 Convergence of clipped-SSTM

Next, we consider an accelerated variant of clipped-SGD called clipped-SSTM (Gorbunov et al., 2020; Gasnikov

and Nesterov, 2016):

$$x^{k+1} = \frac{A_k y^k + a_{k+1} z^k}{A_{k+1}}; \quad (13)$$

$$z^{k+1} = z^k \text{clip}(r f_k(x^{k+1}); a_{k+1}); \quad (14)$$

$$y^{k+1} = \frac{A_k y^k + a_{k+1} z^{k+1}}{A_{k+1}}; \quad (15)$$

where $z^0 = y^0 = x^0$, $a_0 = A_0$, $a_{k+1} = \frac{k+2}{2aL}$ for some parameter $a_{k+1} > 0$, $A_{k+1} = A_k + a_{k+1}$, and $r f_k(x^{k+1})$ is an estimator satisfying Assumption 5.1 sampled independently from previous iterations. Below we formulate the main convergence result for clipped-SSTM in the convex case.

Theorem 5.7. Let Assumptions 2.2 and 2.3 with $\sigma = 0$ hold on $Q = B_{3R}(x)$, where $R > \kappa x^0 - x^k$. Suppose that $r f_k(x^{k+1})$ satisfies Assumption 5.1 with parameters b_k, κ_k for $k = 0; 1; \dots; K$, $K > 0$ and $a = (\min_k A^2; (K+1)^{3-2p} \bar{A}_{LR}; b(K+2)^2 \bar{L}g)$; $\kappa = (\bar{R} = \kappa_{k+1} A)$; where $A = \ln(4(K+1))$ and $b = \max_{k=0; 1; \dots; K} b_k$, $\bar{L} = \max_{k=0; 1; \dots; K} \kappa_k$. Then the iterates produced by clipped-SSTM after K iterations with probability at least $1 - \epsilon$ satisfy

$$f(y^K) - f(x) = \Theta \max \left\{ \frac{LR^2}{K^2}; \frac{R}{K}; bR \right\};$$

As expected for accelerated methods, the above bound has a better $\Theta(K^{-2})$ deterministic term in contrast to the $\Theta(K^{-1})$ corresponding term in the upper bound for clipped-SGD. When the noise is symmetric and has bounded σ -th moment for some $\sigma > 0$ (not necessarily larger than 1), then one can construct an estimator with $b = 0$ and finite \bar{L} (see Proposition 4.1). In this case, the result matches (up to logarithmic factors) the optimal ones derived under bounded variance (Gorbunov et al., 2020) or sub-Gaussian noise (Ghadimi and Lan, 2012) assumptions. The improvement of the deterministic part can be utilized when parallel independent computations of the estimator are possible with marginal overheads (e.g., on communications/aggregation of the results of parallel computations).

Finally, for non-symmetric noise distributions satisfying Assumption 2.1, Theorem 5.7 with Lemma 4.5 imply the following result.

Corollary 5.8 (General noise) Let the assumptions of Theorem 5.7 hold and for all $x \in \mathbb{R}^d$ the noise $\xi = r f(x) - f(x)$ satisfies Assumption 2.1. Then clipped-SSTM with $r f_k(x^k)$ being the smoothed median of means of $O(n)$ samples $r f(x^k)$ and $a = (\min_k A^2; (K+1)^{3-2p} \bar{A}_{LR}; b(K+2)^2 \bar{L}g)$; $\kappa = (\bar{R} = \kappa_{k+1} A)$; where $A = \ln(4(K+1))$ and b and \bar{L} defined in (10), with probability at least $1 - \epsilon$ after K

iterations ensures that $f(y^K) - f(x)$ equals

$$\Theta \max \left\{ \frac{LR^2}{K^2}; \frac{R}{K}; \frac{(1+\sigma^2)d + DR}{2n}; \frac{(1+\sigma)CR}{2n} \right\};$$

where $n > (\max_{j \in [d]} M_j = 2)$ and C, D are defined in (11)-(12). The overall number of oracle calls equals $O(nK)$.

When the non-symmetric part is large, then the same comments are valid as the ones we make after Corollary 5.4. However, when the non-symmetric part is small, then there are regimes when the effect of acceleration is noticeable (for small enough σ).

For the strongly convex problems, we consider a restarted version of clipped-SSTM. We provide the results for this method in Appendix C.2.

6 NUMERICAL EXPERIMENTS

In this section, we illustrate the performance of clipped-SGD combined with the median and smoothed median of means on a simple quadratic problem:

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} x^T A x; \quad (16)$$

where $A \in \mathbb{R}^{d \times d}$ is a randomly generated symmetric positive definite matrix. The code of numerical experiments is available on GitHub³. We consider stochastic gradients of the form $r f(x) = Ax + \xi$, where ξ is an artificial noise following one of the following distributions.

Example 1. Cauchy distribution with the density $p_c(x) = \frac{1}{(1+x^2)^2}$.

Example 2. The mixture of Cauchy and exponential distributions with the density $p(x) = 0.7 p_c(x) + 0.3 e^{-(x+1)} 1_{x > -1} g$.

Example 3. The mixture of Cauchy and Pareto distributions with the density $p(x) = 0.7 p_c(x) + 0.3 \frac{3}{(x+1.5)^4} 1_{x > -1.5} g$.

The experiments check the ability of median and smoothed median of means to deal with symmetric and asymmetrical heavy-tailed noise. We consider two examples of asymmetrical distributions with rapidly (Example 2) and slowly (Example 3) decaying antisymmetric part to examine its influence on the performance of optimization procedures.

We compare the following baselines:

- clipped-MB-SGD
- clipped-SGD where clipping is taken after mini-batching/averaging;

³https://github.com/Kutuz4/AISTATS2024_SMoM

Figure 1: Dependence of the mean error on the oracle calls number with a 95th and 5th percentile bounds.

- ^ MB-clipped-SGD (mini-batched clipped-SGD where averaging is taken after clipping);
- ^ Med-MB-SGD (mini-batched SGD with median instead of averaging);
- ^ clipped-Med-MB-SGD (mini-batched SGD with median instead of averaging and clipping operation after median);
- ^ SMoM-MB-SGD (mini-batched SGD with the smoothed median of means);
- ^ clipped-SMoM-MB-SGD (mini-batched SGD with clipping of the smoothed median of means).

For all methods, except for SMoM-MB-SGD and clipped-SMoM-MB-SGD the batch size is 5, while for SMoM-MB-SGD and clipped-SMoM-MB-SGD we took $\text{SMoM}_{m,n}$ with $m = n = 2$. We have chosen $x_0 = 8 \cdot \bar{d} (1; 1; 1; \dots; 1)^>$, where $d = 50$, as a starting point, launched all the methods 50 times and computed the average errors. The results are displayed in Figure 1.

In the case of a symmetric distribution, Med-MB-SGD and clipped-Med-MB-SGD perform better than clipped-SMoM-MB-SGD due to lower oracle calls count for one iteration. However, as we expected, Med-MB-SGD and clipped-Med-MB-SGD cannot achieve high accuracy in the case of asymmetric distributions due to the presence of a bias, while the smoothed median of means successfully adapts to this situation. Suddenly, averaging gradients after clipping has good performance in the asymmetric case, but it still converges slower compared to clipped-SMoM-MB-SGD. In terms of the number of steps, clipped-SMoM-MB-SGD converges much faster than other methods on asymmetric noise because it needs $(2m + 1)n = 10$ oracle calls on each iteration. We also see that SMoM-MB-SGD has a similar convergence rate to clipped-SMoM-MB-SGD in the case of distributions with less heavy tails.

7 CONCLUSION

In this work, we show that under some structural assumptions on the noise distribution with heavy tails, one can achieve faster convergence in solving of stochastic optimization problems. The key instrument we use is the smoothed median of means, which provably has a small bias and a finite variance for quite a wide class of distributions. Although our results are given for smooth convex/strongly convex problems, using similar technique, one can derive high-probability convergence results for smooth non-convex problems (Sadiev et al., 2023; Nguyen et al., 2023b), non-smooth convex and strongly convex problem (Gorbunov et al., 2021), variational inequalities under some structured non-monotonicity assumptions (Gorbunov et al., 2022), and composite and distributed optimization problems (Gorbunov et al., 2023). One can also improve the logarithmic factors in our results using the technique from (Nguyen et al., 2023a).

Acknowledgements

The work of A. Gasnikov was supported by a grant for research centers in the field of artificial intelligence, provided by the Analytical Center for the Government of the Russian Federation in accordance with the subsidy agreement (agreement identifier 000000D730321P5Q0002) and the agreement with the Ivannikov Institute for System Programming of the Russian Academy of Sciences dated November 2, 2021 No. 70-2021-00142.

References

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. (2016). Deep learning with differential privacy. In Proceedings of

- the 2016 ACM SIGSAC conference on computer and communications security, pages 308 318.
- Bennett, G. (1962). Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33 45.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems* 33:1877 1901.
- Chen, X., Wu, S. Z., and Hong, M. (2020). Understanding gradient clipping in private sgd: A geometric perspective. *Advances in Neural Information Processing Systems* 33:13773 13782.
- Cherapanamjeri, Y., Tripuraneni, N., Bartlett, P., and Jordan, M. (2022). Optimal mean estimation without a variance. In *Conference on Learning Theory*, pages 356 357. PMLR.
- Cutkosky, A. and Mehta, H. (2021). High-probability bounds for non-convex stochastic optimization with heavy tails. *Advances in Neural Information Processing Systems* 34:4883 4895.
- Davis, D., Drusvyatskiy, D., Xiao, L., and Zhang, J. (2021). From low probability to high confidence in stochastic convex optimization. *The Journal of Machine Learning Research* 22(1):2237 2274.
- Devolder, O., Glineur, F., and Nesterov, Y. (2014). First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146:37 75.
- Devroye, L., Lerasle, M., Lugosi, G., and Oliveira, R. I. (2016). Sub-Gaussian mean estimators. *The Annals of Statistics*, 44(6):2695 2725.
- Dzhaparidze, K. and Van Zanten, J. (2001). On Bernstein-type inequalities for martingales. *Stochastic processes and their applications* 93(1):109 117.
- Feller, W. (1971). *An introduction to probability theory and its applications*. Vol. II. Second edition. John Wiley & Sons Inc., New York.
- Freedman, D. A. et al. (1975). On tail probabilities for martingales. *the Annals of Probability*, 3(1):100 118.
- Gasnikov, A. and Nesterov, Y. (2016). Universal fast gradient method for stochastic composite optimization problems. *arXiv preprint arXiv:1604.05275*.
- Ghadimi, S. and Lan, G. (2012). Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization I: A generic algorithmic framework. *SIAM Journal on Optimization* , 22(4):1469 1492.
- Ghadimi, S. and Lan, G. (2013). Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization* , 23(4):2341 2368.
- Gorbunov, E., Danilova, M., Dobre, D., Dvurechenskii, P., Gasnikov, A., and Gidel, G. (2022). Clipped stochastic methods for variational inequalities with heavy-tailed noise. *Advances in Neural Information Processing Systems* 35:31319 31332.
- Gorbunov, E., Danilova, M., and Gasnikov, A. (2020). Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing Systems* 33:15042 15053.
- Gorbunov, E., Danilova, M., Shibaev, I., Dvurechensky, P., and Gasnikov, A. (2021). Near-optimal high probability complexity bounds for non-smooth stochastic optimization with heavy-tailed noise. *arXiv preprint arXiv:2106.05958*
- Gorbunov, E., Sadiev, A., Danilova, M., Horváth, S., Gidel, G., Dvurechensky, P., Gasnikov, A., and Richtárik, P. (2023). High-probability convergence for composite and distributed stochastic minimization and variational inequalities with heavy-tailed noise. *arXiv preprint arXiv:2310.01860*.
- Guo, C., Chaudhuri, K., Stock, P., and Rabbat, M. (2023). Privacy-aware compression for federated learning through numerical mechanism design. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J., editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 11888 11904. PMLR.
- Harvey, N. J., Liaw, C., and Randhawa, S. (2019). Simple and optimal high-probability bounds for strongly-convex stochastic gradient descent. *arXiv preprint arXiv:1909.00843*
- Indritz, J. (1961). An inequality for hermite polynomials. *Proceedings of the American Mathematical Society*, 12(6):981 983.
- Jakovetić, D., Bajović, D., Sahu, A. K., Kar, S., Milošević, N., and Stamenković, D. (2023). Nonlinear gradient mappings and stochastic optimization: A general framework with applications to heavy-tail noise. *SIAM Journal on Optimization* , 33(2):394 423.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. (2021). *Advances and open problems in federated learning*. *Foundations and Trends® in Machine Learning*, 14(1 2):1 210.
- Karimireddy, S. P., He, L., and Jaggi, M. (2021). Learning from history for byzantine robust optimization. In *International Conference on Machine Learning*, pages 5311 5319. PMLR.

- Lecué, G. and Lerasle, M. (2020). Robust machine learning by median-of-means: Theory and practice. *The Annals of Statistics*, 48(2):906–931.
- Li, X. and Orabona, F. (2020). A high probability analysis of adaptive sgd with momentum. arXiv preprint arXiv:2007.14294.
- Lugosi, G. and Mendelson, S. (2019a). Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics* 19(5):1145–1190.
- Lugosi, G. and Mendelson, S. (2019b). Sub-Gaussian estimators of the mean of a random vector. *The Annals of Statistics*, 47(2):783–794.
- Lugosi, G. and Mendelson, S. (2020). Risk minimization by median-of-means tournaments. *Journal of the European Mathematical Society* 22(3):925–965.
- Mai, V. V. and Johansson, M. (2021). Stability and convergence of stochastic gradient clipping: Beyond lipschitz continuity and smoothness. In *International Conference on Machine Learning* pages 7325–7335. PMLR.
- Minsker, S. (2015). Geometric median and robust estimation in Banach spaces. *Bernoulli*, 21(4):2308–2335.
- Nazin, A. V., Nemirovsky, A. S., Tsybakov, A. B., and Juditsky, A. B. (2019). Algorithms of robust stochastic optimization based on mirror descent method. *Automation and Remote Control*, 80:1607–1627.
- Necoara, I., Nesterov, Y., and Glineur, F. (2019). Linear convergence of first order methods for non-strongly convex optimization. *Mathematical Programming*, 175:69–107.
- Nemirovski, A. S., Juditsky, A. B., Lan, G., and Shapiro, A. (2009). Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609.
- Nemirovskij, A. S. and Yudin, D. B. (1983). Problem complexity and method efficiency in optimization.
- Nguyen, T. D., Ene, A., and Nguyen, H. L. (2023a). Improved convergence in high probability of clipped gradient methods with heavy tails. arXiv preprint arXiv:2304.01119.
- Nguyen, T. D., Nguyen, T. H., Ene, A., and Nguyen, H. L. (2023b). High probability convergence of clipped-sgd under heavy-tailed noise. arXiv preprint arXiv:2302.05437.
- OpenAI (2023). GPT-4 technical report.
- Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *International conference on machine learning* pages 1310–1318. Pmlr.
- Rakhlin, A., Shamir, O., and Sridharan, K. (2011). Making gradient descent optimal for strongly convex stochastic optimization. arXiv preprint arXiv:1109.5647.
- Sadiev, A., Danilova, M., Gorbunov, E., Horváth, S., Gidel, G., Dvurechensky, P., Gasnikov, A., and Richtárik, P. (2023). High-probability bounds for stochastic optimization and variational inequalities: the case of unbounded variance. In *Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J., editors, Proceedings of the 40th International Conference on Machine Learning volume 202 of Proceedings of Machine Learning Research* pages 29563–29648. PMLR.
- Vural, N. M., Yu, L., Balasubramanian, K., Volgushev, S., and Erdogdu, M. A. (2022). Mirror descent strikes again: Optimal stochastic convex optimization under in nite noise variance. In *Conference on Learning Theory*, pages 65–102. PMLR.
- Zhang, J., He, T., Sra, S., and Jadbabaie, A. (2019). Why gradient clipping accelerates training: A theoretical justification for adaptivity. arXiv preprint arXiv:1905.11881.
- Zhang, J., Karimireddy, S. P., Veit, A., Kim, S., Reddi, S., Kumar, S., and Sra, S. (2020). Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems* 33:15383–15393.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes, see Sections 2, 4 and 5.](#)
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes, see our results in Sections 4 and 5.](#)
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes, we provide the source code with supplementary materials.](#)
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes.](#)
 - (b) Complete proofs of all theoretical results. [Yes, the proofs are collected in Appendix.](#)
 - (c) Clear explanations of any assumptions. [Yes.](#)
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes.](#)
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes, see Section 6 and Appendix.](#)
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes.](#)
 - (d) A description of the computing infrastructure used. [Not Applicable. We do not use any computing infrastructure, except for an ordinary laptop.](#)
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. [Not Applicable.](#)
 - (b) The license information of the assets, if applicable. [Not Applicable.](#)
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Not Applicable.](#)
 - (d) Information about consent from data providers/curators. [Not Applicable.](#)
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable.](#)
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
 - (a) The full text of instructions given to participants and screenshots. [Not Applicable.](#)
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable.](#)
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable.](#)

Contents

1	INTRODUCTION	1
2	SETUP AND NOTATION	2
3	RELATED WORK	3
4	SMOOTHED MEDIAN OF MEANS AND ITS PROPERTIES	4
5	MAIN RESULTS FOR STOCHASTIC OPTIMIZATION	5
5.1	Convergence of clipped-SGD	6
5.2	Convergence of clipped-SSTM	7
6	NUMERICAL EXPERIMENTS	8
7	CONCLUSION	9
A	SMOOTHED MEDIAN OF MEANS ESTIMATOR	14
A.1	Proof of Proposition 4.1	14
A.2	Proof of Lemma 4.3	15
A.3	Proof of Lemma 4.5	19
A.4	Technical results	23
B	Proofs for clipped-SGD	37
B.1	Auxiliary Results	37
B.2	Quasi-Convex Case	37
B.3	Quasi-Strongly Convex Case	42
C	Proofs for clipped-SSTM	50
C.1	Convex Case	50
C.2	Strongly Convex Case	56
D	PROPERTIES OF HERMITE POLYNOMIALS	59
E	NUMERICAL EXPERIMENTS: ADDITIONAL DETAILS	59

A SMOOTHED MEDIAN OF MEANS ESTIMATOR

A.1 Proof of Proposition 4.1

Let P_j be the cumulative distribution function of ξ_j :

$$P_j(t) = \int_1^{Zt} p_j(u) du; \quad \text{for all } t \geq \mathbb{R}.$$

Then the probability density of the median $\text{Med}(\xi_{j,1}, \dots, \xi_{j,2m+1})$ is given by

$$(2m+1) \binom{2m}{m} P_j(t)^{m-1} (1-P_j(t))^m p_j(t).$$

Let us prove that the variance is bounded. To be more precise, we are going to show that

$$\int_1^{Zt} t^2 P_j(t)^{m-1} (1-P_j(t))^m p_j(t) dt \leq \frac{4B_j}{j} \frac{1}{t^{j-1}} \leq 4 \frac{B_j}{j} t^{1-j}.$$

Since, according to the conditions of the proposition, $p_j(u) \leq B_j (1-u)^{1-j}$ for any $u \geq \mathbb{R}$, it holds that

$$1 - P_j(t) \leq \int_t^{Zt} \frac{B_j}{u^{1+j}} du = \frac{B_j}{j} \frac{1}{t^j} \quad \text{for any } t > 1.$$

Similarly, for any $t \leq 1$, we have

$$P_j(t) \leq \frac{B_j}{j} \frac{1}{t^j}.$$

Hence, for any $t \geq \mathbb{R}$, $P_j(t) \leq 1 - P_j(t)$ satisfies the inequality

$$P_j(t) \leq 1 - P_j(t) \leq \begin{cases} 4 \frac{B_j}{j} t^{1-j} & \text{if } jt > 1; \\ 4 \frac{B_j}{j} & \text{otherwise.} \end{cases} \quad (17)$$

This implies that the integral of $t^2 P_j(t)^{m-1} (1-P_j(t))^m p_j(t)$ over the real line does not exceed

$$\begin{aligned} \int_1^{Zt} t^2 P_j(t)^{m-1} (1-P_j(t))^m p_j(t) dt &\leq \sup_{t \in \mathbb{R}} t^2 P_j(t)^{m-1} (1-P_j(t))^m \int_1^{Zt} p_j(u) du \\ &= \sup_{t \in \mathbb{R}} t^2 P_j(t)^{m-1} (1-P_j(t))^m \\ &= \max \left(\sup_{jt \leq 1} t^2 P_j(t)^{m-1} (1-P_j(t))^m; \sup_{jt > 1} t^2 P_j(t)^{m-1} (1-P_j(t))^m \right) \\ &\leq \max \left(4 \frac{B_j}{j}; \sup_{jt > 1} t^2 P_j(t)^{m-1} (1-P_j(t))^m \right). \end{aligned}$$

We use (17) to bound the supremum in the right-hand side:

$$\begin{aligned} \sup_{jt > 1} t^2 P_j(t)^{m-1} (1-P_j(t))^m &\leq \sup_{jt > 1} t^2 \left(\frac{1}{4} \wedge \frac{B_j}{jt^j} \right)^m \\ &\leq \max \left(\sup_{16jt \leq (4B_j/j)^{1-j}} t^2 4^{-m}; \sup_{jt > (4B_j/j)^{1-j}} t^2 \frac{B_j}{j^m jt^j} \right). \end{aligned}$$

If $m_j > 2$, then we have

$$\sup_{|t| > (4B_j)^{1/m_j}} t^2 \frac{B_j}{|t|^{m_j}} = \frac{4B_j}{j} 4^{-m_j};$$

and, hence,

$$\sup_{|t| > 1} t^2 P_j(t)^{m-1} P_j(t)^m \leq 6 \sup_{|t| > 1} t^2 \frac{1}{4} \wedge \frac{B_j}{|t|^{m_j}} \leq 6 \left(1 - \frac{4B_j}{j}\right)^{2m_j} 4^{-m_j};$$

Thus, we obtain that

$$\int_1^{\infty} t^2 P_j(t)^{m-1} P_j(t)^m p_j(t) dt \leq 6 \sup_{t \in \mathbb{R}} t^2 P_j(t)^{m-1} P_j(t)^m \leq 6 \left(1 - \frac{4B_j}{j}\right)^{2m_j} 4^{-m_j};$$

It only remains to note that

$$\frac{2m}{m} = \frac{(2m)!}{m! m!} = \prod_{j=1}^m \frac{2j}{j} = \prod_{j=1}^m \frac{2j}{j} \leq 2^m \cdot 2^m = 4^m$$

to derive the desired bound

$$E \text{Med}(\xi_{j,1}; \dots; \xi_{j,2m+1})^2 \leq 6(2m+1) \left(1 - \frac{4B_j}{j}\right)^{2m_j}; \tag{18}$$

Concerning the expectation of $\text{Med}(\xi_{j,1}; \dots; \xi_{j,2m+1})$, we point out that it is finite, because $\text{Med}(\xi_{j,1}; \dots; \xi_{j,2m+1})$ has a finite second moment. Moreover, due to the symmetry of p_j , we have $P_j(-t) = 1 - P_j(t)$ and, thus,

$$(-t)P_j(-t)^{m-1} P_j(-t)^m p_j(-t) = t(1 - P_j(t))^{m-1} P_j(t)^m p_j(t) \text{ for all } t \in \mathbb{R}.$$

Hence, it holds that

$$E \text{Med}(\xi_{j,1}; \dots; \xi_{j,2m+1}) = \int_1^{\infty} t P_j(t)^{m-1} P_j(t)^m p_j(t) dt = 0;$$

A.2 Proof of Lemma 4.3

Let $\lambda_{11}; \dots; \lambda_{dd}$ be the diagonal elements of Σ . Denote the difference $r f(x) - r f(x)$ by $\Delta = (\lambda_{11}; \dots; \lambda_{dd})^\top$. It is enough to show that

$$j E \text{SMoM}_{M,n}(\xi_j) \leq 6(2m+1) \frac{j}{2n} \frac{r}{1 + \frac{k}{2n}} \frac{k}{n} \left(\frac{4m}{(2m-1)^2 e} + \frac{m(2m-1)}{(2m-2)^2 e} \frac{k}{2n} + 32 \frac{mk}{2n} \right)^{2\#}$$

$$\cdot (2m+1) \frac{j}{n} \frac{r}{1 + \frac{k}{2n}} \frac{k}{n} \left(1 + \frac{mk}{2n} + \frac{mk}{2n} \right)^{2\#}$$

and

$$E \text{SMoM}_{M,n}(\xi_j) \leq 6(2m+1) \frac{j}{n} + 2^2;$$

for all $j \in \{1; \dots; dg\}$. We start with the upper bound on the second moment. We split the rest of the proof into several steps for convenience.

Step 1: bound on the second moment. For a fixed $j \in \{1, \dots, d\}$, let $p_j(u)$ be the marginal density of j and let

$$F(t) = \int_1^{Z^1} t - \frac{u}{n} p_j^n(u) du$$

stand for the cumulative distribution function of $\text{Mean}(\{j_1, \dots, j_n\}) + j$, where j_1, \dots, j_n are i.i.d. copies of j . Then the density of $\text{SMoM}_m(j; \cdot)$ is equal to

$$(2m + 1) \int_1^{Z^1} t^m F(t)^{m-1} F'(t) dt$$

If we manage to prove that

$$\sup_{t \in \mathbb{R}} t^2 F(t)^{m-1} F'(t) \leq \frac{4(j_1 + \dots + j_n)^2}{4^m};$$

then we immediately obtain

$$\int_1^{Z^1} t^2 F(t)^{m-1} F'(t) dt \leq \sup_{t \in \mathbb{R}} t^2 F(t)^{m-1} F'(t) \int_1^{Z^1} F^0(u) du = \sup_{t \in \mathbb{R}} t^2 F(t)^{m-1} F'(t) \cdot \frac{1}{m};$$

Let j_1, \dots, j_n be independent copies of j and let $j \sim \mathcal{N}(0, 1)$ be a Gaussian random variable, which is independent of j_1, \dots, j_n . Then, according to the definition of $F(t)$, for any $t \in \mathbb{R}$, it holds that

$$1 - F(t) = \mathbb{P}\left(\frac{j_1 + \dots + j_n}{n} + j > t\right);$$

Since j_1, \dots, j_n and j have finite variance, we can apply Chebyshev's inequality to derive an upper bound on the right tail of $F(t)$:

$$1 - F(t) \leq \frac{\mathbb{E}(\frac{j_1 + \dots + j_n}{n} + j)^2}{t^2} = \frac{j_1^2 + \dots + j_n^2 + j^2}{nt^2} = \frac{j_1^2 + \dots + j_n^2}{nt^2} + \frac{1}{t^2} \text{ for all } t > 0.$$

Similarly, for any $t < 0$, we have

$$F(t) \leq \frac{j_1^2 + \dots + j_n^2}{nt^2} + \frac{1}{t^2};$$

Combining these bounds with the inequality $F(t)(1 - F(t)) \leq \frac{1}{4}$, which holds for any $t \in \mathbb{R}$, we deduce that

$$F(t) \leq \frac{1}{4} \wedge \left(\frac{j_1^2 + \dots + j_n^2}{nt^2} + \frac{1}{t^2} \right) = \begin{cases} \frac{1}{4} & \text{if } |t| \leq \frac{1}{\sqrt{j_1^2 + \dots + j_n^2 + 1}} \\ \frac{j_1^2 + \dots + j_n^2}{nt^2} + \frac{1}{t^2} & \text{otherwise.} \end{cases} \quad (19)$$

Hence, for any $m > 1$,

$$\sup_{t \in \mathbb{R}} t^2 F(t)^{m-1} F'(t) \leq \max_{|t| \leq \frac{1}{\sqrt{j_1^2 + \dots + j_n^2 + 1}}} \frac{t^2}{4^m} + \max_{|t| > \frac{1}{\sqrt{j_1^2 + \dots + j_n^2 + 1}}} t^2 \left(\frac{j_1^2 + \dots + j_n^2}{nt^2} + \frac{1}{t^2} \right)^{m-1} = \frac{j_1^2 + \dots + j_n^2 + 1}{4^{m-1}};$$

as we announced. This implies that

$$\mathbb{E} \text{SMoM}_m(j; \cdot)^2 \leq (2m + 1) \int_1^{Z^1} t^2 F(t)^{m-1} F'(t) dt \leq (2m + 1) \frac{1}{m} \frac{j_1^2 + \dots + j_n^2 + 1}{4^{m-1}};$$

Similarly to the proof of Proposition 4.1, we use the inequality

$$\frac{2m}{m} = \frac{(2m)!}{m! m!} = \prod_{j=1}^m \frac{2j}{j} = \prod_{j=1}^m \frac{2j-1}{j} \leq 4^m;$$

which yields that

$$E \text{SMoM}_n(j;)^2 \leq \frac{4(2m+1)}{n} \left(\frac{jj}{n} + \frac{2}{n} \right) :$$

Step 2: bound on the expectation. The rest of the proof is devoted to an upper bound on the expectation of $\text{SMoM}_n(j;)$. One could apply the Cauchy-Schwarz inequality to show that $E \text{SMoM}_n(j;)$ decreases with the growth of n . However, we are going to prove a stronger bound. Our approach is based on decomposition of $p_j(u)$ into the sum of symmetric and antisymmetric part:

$$p_j(u) = s_j(u) + r_j(u); \quad \text{where } s_j(u) = \frac{p_j(u) + p_j(-u)}{2} \quad \text{and} \quad r_j(u) = \frac{p_j(u) - p_j(-u)}{2}.$$

If r_j was equal to zero, we could say that $E \text{SMoM}_n(j;) = 0$ as well. However, in a general situation, r_j has some impact on the mean of $\text{SMoM}_n(j;)$. To quantify it, we compare the integrals

$$\int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt \quad \text{and} \quad \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt;$$

where G is a cumulative distribution function defined as

$$G(t) = \int_0^t \frac{u}{n} s_j^n(u) du;$$

In other words, G corresponds to the CDF of $\text{Mean}(j_{1,1}; \dots; j_{1,n}) + j$, where $j_{1,1}; \dots; j_{1,n}$ are i.i.d. copies of j and $j \sim N(0, 1)$, in the symmetric case. We are going to show that r_j has minor influence on the expectation of the smoothed median of means if n and n are sufficiently large.

First, let us show that the cumulative distribution functions $F(t)$ and $G(t)$ are close to each other. It is straightforward to check that

$$\sup_{x \in \mathbb{R}} |f(x) - g(x)| \leq \frac{1}{2} \sup_{x \in \mathbb{R}} |x - e^{-x^2/2}| \leq \frac{1}{2} \sup_{y=1}^{y^2=2} |y - e^{-y^2/2}| = \frac{1}{2} \left(\frac{1}{2} - \frac{1}{2} \right); \quad (20)$$

Then Lemma A.1 (see Appendix A.4 below) implies that

$$|F(t) - G(t)| \leq \frac{1}{2} \frac{jj}{2n} \quad \text{for all } t \in \mathbb{R}. \quad (21)$$

In view of (19), for any $m > 3$ it holds that

$$|t| F(t)^{m-1} F(t)^m \neq 0 \quad \text{and} \quad |t| G(t)^{m-1} G(t)^m \neq 0 \quad \text{as } t \neq 0.$$

Then, according to Lemma A.2, we have

$$\begin{aligned} & \int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq \frac{1}{2} \frac{jj}{2n} \int_0^1 F(t)^{m-1} F(t)^m dt + \frac{m}{16e} \frac{jj}{2n} \int_0^1 F(t)^{m-1} |F(t)^m - G(t)^m| dt \\ & \quad + \frac{m}{2n} \frac{jj}{2n} \sup_{t \in \mathbb{R}} |t| \max_{2[F(t) \wedge G(t); F(t) - G(t)]} |F(t) - G(t)|^{m-2} (1 - |F(t) - G(t)|)^{m-2} : \end{aligned}$$

Let us remind the reader that the CDF $F(t)$ satisfies the inequalities

$$1 - F(t) \leq \frac{jj}{nt^2} + \frac{2}{t^2} \quad \text{for all } t > 0 \quad \text{and} \quad F(t) \leq \frac{jj}{nt^2} + \frac{2}{t^2} \quad \text{for all } t < 0;$$

Due to the Chebyshev inequality, a similar bound holds for $G(t)$:

$$\mathbb{1} - G(t) \leq \frac{jj}{nt^2} + \frac{2}{t^2} \quad \text{for all } t > 0 \quad \text{and} \quad G(t) \leq \frac{jj}{nt^2} + \frac{2}{t^2} \quad \text{for all } t < 0:$$

This yields that

$$\begin{aligned} & \int_{-Z^1}^Z F(t)^{m-1} F(t)^m F^0(t) dt - \int_{-Z^1}^Z G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq \frac{1}{2} \frac{1}{2e} \frac{jj}{2n} \int_{-Z^1}^Z \frac{1}{4} \wedge \frac{jj = n + 2}{t^2} dt + \frac{m}{16e} \frac{jj}{2n} \int_{-Z^1}^Z \frac{1}{4} \wedge \frac{jj = n + 2}{t^2} dt \\ & \quad + \frac{m}{2n} \frac{jj}{2} \sup_{t \in 2R} |t| \int_{-Z^1}^Z \frac{1}{4} \wedge \frac{jj = n + 2}{t^2} dt : \end{aligned}$$

Applying Proposition A.3, we obtain that

$$\begin{aligned} & \int_{-Z^1}^Z F(t)^{m-1} F(t)^m F^0(t) dt - \int_{-Z^1}^Z G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq \frac{1}{2} \frac{1}{2e} \frac{jj}{2n} \frac{2m}{(2m-1)4^{m-1}} \frac{1}{r} \frac{1}{2 + \frac{jj}{n}} \\ & \quad + \frac{m}{16e} \frac{jj}{2n} \frac{2m-1}{(2m-2)4^{m-2}} \frac{1}{r} \frac{1}{2 + \frac{jj}{n}} \\ & \quad + \frac{m}{2n} \frac{jj}{2} \sup_{t \in 2R} |t| \int_{-Z^1}^Z \frac{1}{4} \wedge \frac{jj = n + 2}{t^2} dt : \end{aligned}$$

Similarly to Step 1, we can prove that

$$\left(\sup_{t \in 2R} |t| \int_{-Z^1}^Z \frac{1}{4} \wedge \frac{jj = n + 2}{t^2} dt \right) = \frac{2^p}{4^{m-2}} \frac{jj = n}{2 + \frac{jj}{n}};$$

and then it holds that

$$\begin{aligned} & \int_{-Z^1}^Z F(t)^{m-1} F(t)^m F^0(t) dt - \int_{-Z^1}^Z G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq \frac{m}{(2m-1)4^{m-1}} \frac{1}{2e} \frac{jj}{2n} + \frac{m(2m-1)}{(2m-2)4^m e} \frac{jj}{2n} + \frac{32}{4^m} \frac{m}{2n} \frac{jj}{2} \frac{1}{r} \frac{1}{2 + \frac{jj}{n}} : \end{aligned}$$

Taking into account

$$\frac{2m}{m} = \frac{(2m)!}{m! m!} = \prod_{j=1}^m \frac{2j}{j} \prod_{j=1}^m \frac{2j}{j} \leq 4^m;$$

we immediately obtain that

$$\begin{aligned} \mathbb{E} \text{SMoM}_n(j; \cdot) & \leq (2m+1) \frac{jj}{n} \frac{1}{1 + \frac{jj}{2n}} \frac{4m}{(2m-1)2e} + \frac{m(2m-1)}{(2m-2)e} \frac{jj}{2n} + 32 \frac{m}{2n} \frac{jj}{2} \frac{1}{r} \frac{1}{2 + \frac{jj}{n}} \\ & \leq (2m+1) \frac{jj}{n} \frac{1}{1 + \frac{k}{2n}} \frac{4m}{(2m-1)2e} + \frac{m(2m-1)}{(2m-2)e} \frac{k}{2n} + 32 \frac{mk}{2n} \frac{k}{2} \frac{1}{r} \frac{1}{2 + \frac{k}{n}} \\ & \leq (2m+1) \frac{jj}{n} \frac{1}{1 + \frac{k}{2n}} \left(1 + \frac{mk}{2n} \right) + \frac{mk}{2n} \frac{k}{2} \frac{1}{r} \frac{1}{2 + \frac{k}{n}} : \end{aligned}$$

A.3 Proof of Lemma 4.5

Let $\mathbf{a} = (a_1, \dots, a_d)^T$ stand for the difference $r f(x) - r f(x)$ and let us show that, for any $j \in \{1, \dots, d\}$, it holds that

$$j E S M o M_n(\mathbf{a}; \mathbf{a}) \leq \frac{m M_j}{2n} \left(1 + \frac{2 B_j}{n^{j-1}} \right)^{1=j \#}$$

and

$$E S M o M_n(\mathbf{a}; \mathbf{a})^2 \leq m \left(1 + \frac{M_j}{n} + \frac{2 B_j}{j n^{j-1}} + \frac{B_j M_j}{n^j} \right)^{2=(j+1) \#} :$$

From now on, we fix an arbitrary $j \in \{1, \dots, d\}$. Similarly to the proof of Lemma 4.3, the core idea is to compare the cumulative distribution functions

$$F(t) = \int_0^t \frac{u}{n} p_j^n(u) du \quad \text{and} \quad G(t) = \int_0^t \frac{u}{n} \varsigma_j^n(u) du :$$

The first one is directly related to the density of $S M o M_n(\mathbf{a}; \mathbf{a})$, which is equal to

$$(2m+1) \binom{2m}{m} F(t)^{m-1} (1-F(t))^m F'(t) :$$

However, the proof of Lemma 4.5 is far more technical. The main obstacle is that we cannot use Chebyshev's inequality to specify the rate of decay of $F(t) - F(t)$ and of $G(t) - G(t)$ as t approaches infinity. Instead, we prove the following non-trivial result (see Lemma A.4 below): if Assumption 2.1 holds and $2M_j \leq n^2$, then, for any $t \in \mathbb{R}$ it holds that

$$|F(t) - G(t)| \leq \frac{M_j}{nt} \left(1 + \frac{B_j}{n^{j-1} |t|^{j-1}} \right) :$$

Combining this result with the bound on the second derivative of (20) and Lemma A.1, we obtain that

$$|F(t) - G(t)| \leq \frac{M_j}{2n} \left(1 + \frac{B_j}{n^{j-1} |t|^{j+1}} \right) : \tag{22}$$

Despite the simple statement, the proof of Lemma A.4 is quite technical. A reader can find it in Appendix A.4. With the bound (22) at hand, the proof of Lemma 4.5 is relatively simple. For convenience, we divide it into several steps.

Step 1: a bound on the tails of G . The goal of this step is to specify the rate of decay of $G(t) - G(t)$ as t tends to infinity. First, consider the case $t > 0$. By the definition of $G(t)$, it holds that

$$\begin{aligned} 1 - G(t) &= \int_t^\infty \frac{u}{n} \varsigma_j^n(u) du \\ &= n \int_1^\infty (y) \varsigma_j^n(nt + ny) dy \\ &= n \int_1^{t=2} (y) \varsigma_j^n(nt + ny) dy + n \int_{t=2}^\infty (y) \varsigma_j^n(nt + ny) dy : \end{aligned} \tag{23}$$

If $y \geq t=2$, then

$$(y) \leq 1 \quad (t=2) \leq \exp \left(-\frac{t^2}{8} \right) ;$$

and we have

$$\int_1^{Z=2} n^{-1} \int_0^y \xi^n(nt + ny) dy \leq \exp\left\{-\frac{t^2}{8}\right\} \int_1^{Z=2} \xi^n(nt + ny) n dy \leq \exp\left\{-\frac{t^2}{8}\right\} : \quad (24)$$

Otherwise, if $y > t=2$, then, due to Assumption 2.1, it holds that

$$\int_{t=2}^{Z=1} n^{-1} \int_0^y \xi^n(nt + ny) dy \leq \int_{t=2}^{Z=1} \frac{B_j n^2}{n^{(1+j)=i} + n^{1+j}(t+y)^{1+j}} dy \leq \int_0^{Z=1} \frac{B_j}{n^{i-1}(t=2+v)^{1+j}} dv = \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} : \quad (25)$$

Plugging the inequalities (24), (25) into (23), we obtain that

$$1 - G(t) \leq \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \exp\left\{-\frac{t^2}{8}\right\} \quad \text{for all } t > 0.$$

Similarly, we can prove that

$$G(t) \leq \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \exp\left\{-\frac{t^2}{8}\right\} \quad \text{for all } t < 0.$$

Hence, for any $t \in \mathbb{R}$, we have

$$|G(t) - 1| \leq \min\left\{\frac{1}{4}, \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \exp\left\{-\frac{t^2}{8}\right\}\right\} : \quad (26)$$

Step 2: bound on the second moment. The second moment of $\text{SMoM}_n(j; \cdot)$ satisfies the inequality

$$\begin{aligned} E \text{SMoM}_n(j; \cdot)^2 &= (2m+1) \int_1^{Z=1} t^{2m} F(t)^{m-1} (1-F(t))^m F^0(t) dt \\ &\leq (2m+1) \int_1^{Z=1} t^{2m} F(t)^{m-1} (1-F(t))^m dt : \end{aligned}$$

On the other hand, (22) and (26) imply that

$$|F(t) - 1| \leq \min\left\{\frac{1}{4}, \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \exp\left\{-\frac{t^2}{8}\right\} + \frac{M_j}{nt}\right\} + \frac{B_j}{j n^{i-1} t^{-j}} \quad \text{for all } t \in \mathbb{R}. \quad (27)$$

This yields

$$\int_1^{Z=1} t^{2m} F(t)^{m-1} (1-F(t))^m dt \leq 4^{-m} \int_1^{Z=1} t^{2m} \left(1 + \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \frac{M_j}{nt}\right)^{2m} dt \leq \frac{M_j^2}{n} + \frac{2^{-j} B_j}{j n^{i-1}} + \frac{B_j M_j}{n^{i-1}} \quad \text{if } m_j > 2:$$

Since

$$\frac{2m}{m} = \frac{(2m)!}{m! m!} = \prod_{j=1}^m \frac{2j}{j} = \prod_{j=1}^m \frac{2j}{j} \leq 4^m;$$

we obtain that

$$E \text{SMoM}_n(j; \cdot)^2 \leq m \int_1^{Z=1} t^{2m} \left(1 + \frac{2^{-j} B_j}{j n^{i-1} t^{-j}} + \frac{M_j}{nt}\right)^{2m} dt \leq m \left(\frac{M_j^2}{n} + \frac{2^{-j} B_j}{j n^{i-1}} + \frac{B_j M_j}{n^{i-1}}\right) :$$

Step 3: bound on the expectation. It remains to bound the expectation of $SMoM_n(j; \cdot)$. For this purpose, we use Lemma A.2, which yields that

$$\begin{aligned}
 & (2m+1) \int_0^1 \int_{\mathbb{R}} |jESMoM_n(j; \cdot)| \\
 &= \int_0^1 \int_{\mathbb{R}} |F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 \int_{\mathbb{R}} |G(t)^{m-1} G(t)^m G^0(t) dt \\
 & \leq \int_0^1 \int_{\mathbb{R}} |G(t)^{m-1} G(t)^m F(t) - G(t)^{m-1} G(t)^m F(t)| dt + \frac{m}{2} \int_0^1 \int_{\mathbb{R}} |G(t)^{m-1} G(t)^m F(t) - G(t)^{m-1} G(t)^m F(t)|^2 dt \\
 & \quad + m^2 \sup_{t \in \mathbb{R}} |jt| F(t) - G(t)^2 \max_{2[F(t) \wedge G(t); F(t) - G(t)]} |m-2| (1 - |j|)^{m-2} :
 \end{aligned}$$

Note that the requirement

$$|jt|G(t)^{m-1} G(t)^m F(t) - G(t) \leq 0 \text{ and } |jt|G(t)^{m-1} G(t)^m F(t) - G(t)^2 \leq 0 \quad \forall t \in \mathbb{R}$$

is satisfied, because of the inequalities (22), (26) and the conditions of the lemma. These inequalities also imply that

$$\begin{aligned}
 & \int_0^1 \int_{\mathbb{R}} |F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 \int_{\mathbb{R}} |G(t)^{m-1} G(t)^m G^0(t) dt \\
 & \leq \int_0^1 \int_{\mathbb{R}} \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \exp\left(-\frac{t^2}{8} \right)^m |F(t) - G(t)| dt \right. \\
 & \quad \left. + \frac{m}{2} \int_0^1 \int_{\mathbb{R}} \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \exp\left(-\frac{t^2}{8} \right)^{m-1} |F(t) - G(t)|^2 dt \right) \right. \\
 & \quad \left. + m^2 \sup_{t \in \mathbb{R}} |jt| F(t) - G(t)^2 \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \frac{M_j}{nt} \left(1 + \frac{B_j}{|n_j - 1| |jt|} \right) + \exp\left(-\frac{t^2}{8} \right)^{m-2} \right) :
 \end{aligned}$$

Taking into account the bound (22) on the absolute value of the difference $F(t) - G(t)$, we obtain that

$$\begin{aligned}
 & \int_0^1 \int_{\mathbb{R}} |F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 \int_{\mathbb{R}} |G(t)^{m-1} G(t)^m G^0(t) dt \\
 & \leq \frac{M_j}{2n} \int_0^1 \int_{\mathbb{R}} \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \exp\left(-\frac{t^2}{8} \right)^m dt \right. \\
 & \quad \left. + \frac{m}{2} \frac{M_j}{2n} \int_0^1 \int_{\mathbb{R}} \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \exp\left(-\frac{t^2}{8} \right)^{m-1} dt \right) \right. \\
 & \quad \left. + m^2 \frac{M_j}{2n} \sup_{t \in \mathbb{R}} |jt| \left(\frac{1}{4} \wedge \frac{2|B_j|}{|n_j - 1| |jt|} + \frac{M_j}{nt} \left(1 + \frac{B_j}{|n_j - 1| |jt|} \right) + \exp\left(-\frac{t^2}{8} \right)^{m-2} \right) :
 \end{aligned}$$

The inequality $(a + b)^k \leq 2^k (a^k + b^k)$, which holds for any $k > 1$ and positive a and b , implies that

$$\begin{aligned} & \int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt + \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq \frac{M_j}{2n} \int_0^1 \frac{1}{4} \wedge \frac{2^{j+1} B_j}{j n^{j-1} |t|^{j-1}} dt + \frac{M_j}{2n} \int_0^1 \frac{1}{4} \wedge 2 \exp\left(-\frac{t^2}{8}\right)^m dt \\ & + \frac{m}{2} \frac{M_j}{2n} \int_0^1 \frac{1}{4} \wedge \frac{2^{j+1} B_j}{j n^{j-1} |t|^{j-1}} dt + \frac{m}{2} \frac{M_j}{2n} \int_0^1 \frac{1}{4} \wedge 2 \exp\left(-\frac{t^2}{8}\right)^{m-1} dt \\ & + m^2 \frac{M_j}{2n} \sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{2^{j+1} B_j}{j n^{j-1} |t|^{j-1}} \right)^{m-2} + m^2 \frac{M_j}{2n} \sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{4M_j}{n |t|} \right)^{m-2} \\ & + m^2 \frac{M_j}{2n} \sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{8M_j B_j}{n^j |t|^{j+1}} \right)^{m-2} + m^2 \frac{M_j}{2n} \sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge 8 \exp\left(-\frac{t^2}{8}\right)^{m-2} \right) : \end{aligned}$$

Due to Proposition A.3, it holds that

$$\int_0^1 \frac{1}{4} \wedge \frac{2^{j+1} B_j}{j n^{j-1} |t|^{j-1}} dt + \int_0^1 \frac{1}{4} \wedge 2 \exp\left(-\frac{t^2}{8}\right)^m dt \leq 4^m (1 + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} :$$

Since

$$\begin{aligned} & \left(\sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{2^{j+1} B_j}{j n^{j-1} |t|^{j-1}} \right)^{m-2} \right) + \left(\sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{4M_j}{n |t|} \right)^{m-2} \right) \\ & + \left(\sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge \frac{8M_j B_j}{n^j |t|^{j+1}} \right)^{m-2} \right) + \left(\sup_{t \in \mathbb{R}} |t| \left(\frac{1}{4} \wedge 4 \exp\left(-\frac{t^2}{8}\right)^{m-2} \right) \right) \\ & \leq 4^m (1 + \frac{M_j}{n} + \frac{M_j B_j}{n^j})^{1-(j+1)} + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} ; \end{aligned}$$

we conclude that

$$\begin{aligned} & \int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt + \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt \\ & \leq 4^m \frac{M_j}{2n} (1 + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} \\ & + m^2 4^m \frac{M_j}{2n} (1 + \frac{M_j}{n} + \frac{M_j B_j}{n^j})^{1-(j+1)} + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} : \end{aligned}$$

Then it holds that

$$\begin{aligned} & j \text{ESMoM}_n(j; j) \leq (2m+1) \frac{2m}{m} 4^m \frac{M_j}{2n} (1 + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} \\ & + (2m+1) \frac{2m}{m} m^2 4^m \frac{M_j}{2n} (1 + \frac{M_j}{n} + \frac{M_j B_j}{n^j})^{1-(j+1)} + \frac{2^{j+1} B_j}{j n^{j-1}})^{\frac{1}{j}} \\ & \leq \frac{mM_j}{2n} (1 + \frac{2^j B_j}{n^{j-1}})^{\frac{1}{j}} : \end{aligned}$$

The proof is finished.

A.4 Technical results

Lemma A.1. Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a twice differentiable function with uniformly bounded second derivative:

$$|h''(x)| \leq H \quad \text{for all } x \in \mathbb{R}.$$

Let $q(x)$ be a probability density of a random variable and let

$$s(x) = \frac{q(x) + q(-x)}{2} \quad \text{and} \quad a(x) = \frac{q(x) - q(-x)}{2}$$

stand for its symmetric and antisymmetric parts, respectively. Assume that there exists $M > 0$ such that the function $s(x)$ fulfils

$$\int_{-1}^1 x s(x) dx = 0 \quad \text{and} \quad \int_{-1}^1 x^2 s(x) dx \leq M;$$

Then, for any positive integer n , it holds that

$$\int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) q^n(x) dx - \int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) s^n(x) dx \leq \frac{HM}{2n};$$

provided that the integrals in the left-hand side converge.

Proof. Let X_1, \dots, X_n be i.i.d. random variables with the density $q(x)$. It is known that $(X_1 + \dots + X_n) \sim q^n(x)$. Then the integral

$$\int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) q^n(x) dx$$

admits a representation

$$\int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) q^n(x) dx = E h\left(\frac{X_1 + \dots + X_n}{n}\right) = \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) q(x_1) \dots q(x_n) dx_1 \dots dx_n;$$

Similarly, it holds that

$$\int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) s^n(x) dx = \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) s(x_1) \dots s(x_n) dx_1 \dots dx_n;$$

and thus,

$$\int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) q^n(x) dx - \int_{-\infty}^{\infty} h\left(\frac{x}{n}\right) s^n(x) dx = \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) q(x_1) \dots q(x_n) dx_1 \dots dx_n - \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) s(x_1) \dots s(x_n) dx_1 \dots dx_n;$$

Let us introduce

$$p_k(x_1, \dots, x_n) = \prod_{i=1}^k q(x_i) \prod_{i=k+1}^n s(x_i); \quad \text{where } k \in \{0, \dots, n\};$$

Then it holds that

$$\begin{aligned} & \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) q(x_1) \dots p(x_n) dx_1 \dots dx_n \\ &= \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) s(x_1) \dots s(x_n) dx_1 \dots dx_n \\ &= \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) p_0(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) p_1(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \sum_{k=1}^n \int_{-\infty}^{\infty} h\left(\frac{x_1 + \dots + x_n}{n}\right) p_k(x_1, \dots, x_n) dx_1 \dots dx_n; \end{aligned}$$

Let us fix any $k \in \{1, \dots, n\}$ and consider

$$Z \int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n$$

Note that, according to the definition of μ_k , we have

$$\mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) = \prod_{i=1}^{k-1} q(x_i) \prod_{i=k}^n (x_i - x_k)$$

Moreover, due to the Taylor's expansion with the Lagrange remainder term, it holds that

$$\int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n = \int_{\mathcal{H}} \frac{1}{n} \sum_{i \in [k]} x_i \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n + \frac{H x_k^2}{2n^2}$$

Since, according to the definition of μ_k and the conditions of the lemma,

$$\int_{\mathcal{H}} x_k \mu_k(x_1, \dots, x_n) dx_1 \dots dx_n = 0; \quad \int_{\mathcal{H}} x_k \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n = 0; \quad \text{and} \quad \int_{\mathcal{H}} x_k^2 \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n \leq M^2;$$

we have

$$\begin{aligned} Z \int_{\mathcal{H}} \frac{1}{n} \sum_{i \in [k]} x_i \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n &= 0; \\ Z \int_{\mathcal{H}} \frac{1}{n} \sum_{i \in [k]} x_i \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n &= 0; \end{aligned}$$

and then

$$\begin{aligned} Z \int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n \\ \leq \frac{H}{2n^2} \int_{\mathcal{H}} \sum_{i=1}^{k-1} q(x_i) \prod_{i=k}^n (x_i - x_k) dx_1 \dots dx_n \leq \frac{HM}{2n^2}. \end{aligned}$$

Finally, applying the triangle inequality, we obtain that

$$\begin{aligned} Z \int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} q(x_1) \dots q(x_n) dx_1 \dots dx_n \\ \leq Z \int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} (x_1) \dots (x_n) dx_1 \dots dx_n \\ \leq \sum_{k=1}^n Z \int_{\mathcal{H}} \frac{x_1 + \dots + x_n}{n} \mu_k(x_1, \dots, x_n) \mu_{k-1}(x_1, \dots, x_n) dx_1 \dots dx_n \\ \leq \sum_{k=1}^n \frac{HM}{2n^2} = \frac{HM}{2n}. \end{aligned}$$

Lemma A.2. Let F and G be any differentiable cumulative distribution functions, such that

$$\int_{\mathcal{H}} |G(t)^{m-1} - G(t)^m| F(t) - G(t) dt \leq 0 \quad \text{and} \quad \int_{\mathcal{H}} |G(t)^{m-1} - G(t)^m| F(t) - G(t)^2 dt \leq 0 \quad \text{as } t \rightarrow 1^-.$$

Then it holds that

$$\int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt$$

$$= \int_0^1 t G(t)^{m-1} G(t)^m F(t) G^0(t) dt + \frac{m}{2} \int_0^1 G(t)^{m-1} G(t)^m F(t) G(t)^2 dt$$

$$+ m^2 \sup_{t \in \mathbb{R}} |t| F(t) G(t)^2 \max_{2[F(t) \wedge G(t); F(t) - G(t)]} (1 - |t|)^{m-2} :$$

Proof. The proof is based on integration by parts. Let us define a function $\phi : [0; 1] \rightarrow \mathbb{R}$ as follows:

$$\phi(x) = \frac{1}{m^2} (x^m (1-x)^m)^{00} = \frac{1}{m} x^{m-2} (1-x)^{m-2} (m-1)(1-2x)^2 - 2x(1-x) :$$

Note that, for any $x \in [0; 1]$, we have

$$\frac{1}{2m} x^{m-2} (1-x)^{m-2} \phi(x) = \frac{1}{m} x^{m-2} (1-x)^{m-2} : \tag{28}$$

Due to Taylor's expansion with the Lagrange remainder term, for any $t \in \mathbb{R}$, there exists $\xi(t) \in [F(t) \wedge G(t); F(t) - G(t)]$, such that

$$F(t)^m - G(t)^m = m G(t)^{m-1} (F(t) - G(t)) + \frac{m^2}{2} \phi(\xi(t)) F(t) G(t)^2 :$$

Then it holds that

$$\int_0^1 t F(t)^{m-1} F(t)^m F^0(t) dt - \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt$$

$$= \int_0^1 t G(t)^{m-1} G(t)^m F^0(t) G^0(t) dt$$

$$+ m \int_0^1 t G(t)^{m-1} G(t)^m F(t) G^0(t) dt - \int_0^1 t G(t)^{m-1} G(t)^m G^0(t) dt$$

$$+ \frac{m^2}{2} \int_0^1 t \phi(\xi(t)) F(t) G(t)^2 F^0(t) dt : \tag{29}$$

Let us focus on the first term in the right-hand side. Integration by parts yields that

$$\int_0^1 t G(t)^{m-1} G(t)^m F^0(t) G^0(t) dt = \int_0^1 G(t)^{m-1} G(t)^m F(t) G(t) dt$$

$$- m \int_0^1 t G(t)^{m-1} G(t)^m F(t) G^0(t) dt :$$

Substituting this equality into (29), we obtain that

$$\begin{aligned}
 & \int_0^Z |F(t)|^{m-1} |F(t) - F^0(t)| dt + \int_0^Z |G(t)|^{m-1} |G(t) - G^0(t)| dt \\
 &= \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |G(t)| dt \\
 &+ m \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |2G(t) - F(t) - G(t) - F^0(t) - G^0(t)| dt \\
 &+ \frac{m^2}{2} \int_0^Z |G(t) - F(t)| |G(t) - F^0(t)| dt.
 \end{aligned} \tag{30}$$

Now we can apply integration by parts to the second term in the right-hand side of (30):

$$\begin{aligned}
 & m \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |2G(t) - F(t) - G(t) - F^0(t) - G^0(t)| dt \\
 &= \frac{m}{2} \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |2G(t) - F(t) - G(t)|^2 dt \\
 &+ \frac{m^2}{2} \int_0^Z |G(t) - F(t)| |G(t) - F^0(t)| |G(t) - G^0(t)| dt.
 \end{aligned}$$

Hence, we can rewrite the equality (30) in the following form:

$$\begin{aligned}
 & \int_0^Z |F(t)|^{m-1} |F(t) - F^0(t)| dt + \int_0^Z |G(t)|^{m-1} |G(t) - G^0(t)| dt \\
 &= \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |G(t)| dt \\
 &+ \frac{m}{2} \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |2G(t) - F(t) - G(t)|^2 dt \\
 &+ \frac{m^2}{2} \int_0^Z |G(t) - F(t)| |F^0(t) - G^0(t) - F(t) - G(t)|^2 dt.
 \end{aligned}$$

Then, due to (28) and the triangle inequality, it holds that

$$\begin{aligned}
 & \int_0^Z |F(t)|^{m-1} |F(t) - F^0(t)| dt + \int_0^Z |G(t)|^{m-1} |G(t) - G^0(t)| dt \\
 & \leq \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |G(t)| dt + \frac{m}{2} \int_0^Z |G(t)|^{m-1} |G(t) - F(t)| |2G(t) - F(t) - G(t)|^2 dt \\
 & + \frac{m^2}{2} \int_0^Z |G(t) - F(t)| |F^0(t) - G^0(t) - F(t) - G(t)|^2 dt.
 \end{aligned}$$

Since $\int_1^{Z^1} |2G(t) - 1| dt$,

$$\int_1^{Z^1} |F(t) - G(t)|^2 dt \leq \sup_{t \in \mathbb{R}^+} |F(t) - G(t)| \int_1^{Z^1} |F(t) - G(t)| dt \leq \sup_{t \in \mathbb{R}^+} |F(t) - G(t)|^2 \int_1^{Z^1} dt;$$

and, similarly,

$$\int_1^{Z^1} |G(t) - F(t)|^2 dt \leq \sup_{t \in \mathbb{R}^+} |G(t) - F(t)| \int_1^{Z^1} |G(t) - F(t)| dt \leq \sup_{t \in \mathbb{R}^+} |G(t) - F(t)|^2 \int_1^{Z^1} dt;$$

we finally obtain that

$$\begin{aligned} & \int_1^{Z^1} |F(t) - G(t)|^2 dt \leq \sup_{t \in \mathbb{R}^+} |F(t) - G(t)| \int_1^{Z^1} |F(t) - G(t)| dt \\ & \leq \sup_{t \in \mathbb{R}^+} |F(t) - G(t)| \int_1^{Z^1} |F(t) - G(t)| dt + \frac{m}{2} \int_1^{Z^1} |G(t) - F(t)|^2 dt \\ & + m^2 \sup_{t \in \mathbb{R}^+} |F(t) - G(t)|^2 \max_{2[F(t) \wedge G(t); F(t) - G(t)]} (1 - \max_{2[F(t) \wedge G(t); F(t) - G(t)]})^m : \end{aligned}$$

Proposition A.3. For any positive numbers a, k, α , and β , such that $k > \alpha + 1$, it holds that

$$\int_1^{Z^1} |t^{-\alpha} - \frac{a}{4t^k}| dt = \frac{k}{(\alpha + 1)(k - \alpha - 1)} \frac{2(4^{1-\alpha} - a)^{\alpha + 1}}{4^k}.$$

Proof. The proof follows from simple calculations:

$$\begin{aligned} \int_1^{Z^1} |t^{-\alpha} - \frac{a}{4t^k}| dt &= 2 \int_0^{Z^1} |t^{-\alpha} - \frac{a}{4t^k}| dt \\ &= \frac{2}{4^k} \int_0^{4Z^1/a} |t^{-\alpha} - 1| dt + 2a^k \int_{4Z^1/a}^{Z^1} |t^{-\alpha} - 1| dt \\ &= \frac{2(4^{1-\alpha} - a)^{\alpha + 1}}{(\alpha + 1)4^k} + \frac{2a^k (4^{1-\alpha} - a)^{\alpha + 1}}{(k - \alpha - 1)} \\ &= \frac{2(4^{1-\alpha} - a)^{\alpha + 1}}{4^k} \left(\frac{1}{\alpha + 1} + \frac{1}{k - \alpha - 1} \right) \\ &= \frac{k}{(\alpha + 1)(k - \alpha - 1)} \frac{2(4^{1-\alpha} - a)^{\alpha + 1}}{4^k}. \end{aligned}$$

Lemma A.4. Grant Assumption 2.1 and let

$$F(t) = \int_1^{Z^1} t^{-\frac{u}{n}} p_j^n(u) du = \int_{\mathbb{R}^n} t^{-\text{Mean}(u_1, \dots, u_n)} p_j(u_1) \dots p_j(u_n) du$$

and

$$G(t) = \int_1^{Z^1} t^{-\frac{u}{n}} \xi_j^n(u) du = \int_{\mathbb{R}^n} t^{-\text{Mean}(u_1, \dots, u_n)} \xi_j(u_1) \dots \xi_j(u_n) du$$

be two cumulative distribution functions. Assume that $2M_j \leq n^2$. Then, for any $t \in \mathbb{R}$ it holds that

$$jF(t) - G(t) \leq \frac{M_j}{nt} \left(1 + \frac{1}{t} + \frac{B_j}{n^{j-1}t^j} \right) :$$

Proof. Since $p_j(u) = \varsigma_j(u) + r_j(u)$, it holds that

$$p_j^n(u) = \sum_{k=0}^n \binom{n}{k} r_j^k \varsigma_j^{(n-k)}(u);$$

and then

$$\begin{aligned} F(t) - G(t) &= \int_0^t \frac{u}{n} p_j^n(u) du - \int_0^t \frac{u}{n} \varsigma_j^n(u) du \\ &= \sum_{k=1}^n \binom{n}{k} \int_0^t \frac{u}{n} r_j^k \varsigma_j^{(n-k)}(u) du \\ &= \sum_{k=1}^n \binom{n}{k} \int_{\mathbb{R}^n} t \text{Meas}(u_1, \dots, u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n : \end{aligned} \tag{31}$$

In the rest of the proof, we bound the summands in the right-hand side one by one. For readability, we split our derivations into several steps.

Step 1: Taylor's expansion. Let us fix an arbitrary $k \geq 1, \dots, n$ and consider

$$\int_{\mathbb{R}^n} t \text{Meas}(u_1, \dots, u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n :$$

Using Taylor's expansion with the integral remainder term, we rewrite the expression of interest in the form

$$\begin{aligned} &\int_{\mathbb{R}^n} t \text{Meas}(u_1, \dots, u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n \\ &= \int_{\mathbb{R}^n} \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{u_1}{n} \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{u_1}{n} \frac{u_1}{n} \\ &\quad + \frac{u_1^2}{n^2} \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{(1-v_1)u_1}{n} \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{u_1}{n} \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n : \end{aligned} \tag{32}$$

Note that

$$\int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{u_1}{n} \quad \text{and} \quad \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{u_1}{n}$$

do not depend on u_1 . Since, according to Assumption 2.1, it holds that

$$\int_0^1 r_j(u_1) du_1 = 0 \quad \text{and} \quad \int_0^1 u_1 r_j(u_1) du_1 = 0;$$

the right-hand side of (32) simplifies to

$$\begin{aligned} &\int_{\mathbb{R}^n} \int_0^t \text{Meas}(u_1, \dots, u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n \\ &= \int_0^1 \int_{\mathbb{R}^n} v_1 dv_1 \int_0^t \text{Meas}(u_1, \dots, u_n) + \frac{(1-v_1)u_1}{n} \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n \varsigma_j(u_i) du_1 \dots du_n : \end{aligned}$$

Repeating this trick ($k - 1$) more times, we obtain that

$$\begin{aligned} & \int_{\mathbb{R}^n} \text{Meas}(u_1; \dots; u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n s_j(u_i) du_1 \dots du_n \\ &= \int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^n} \text{Meas}(u_1; \dots; u_n) + \sum_{i=1}^k \frac{(1-v_i)u_i}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \prod_{i=k+1}^n s_j(u_i) ; \end{aligned}$$

where $(2k)$ stands for the $(2k)$ -th derivative of Meas . Due to the properties of convolution, the integral in the right-hand side is equal to

$$\int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^n} \text{Meas}(u_1; \dots; u_n) + \sum_{i=1}^k \frac{v_i u_i}{n} \prod_{i=k+1}^n \frac{u_i}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(u_i) ;$$

Substituting $u_{k+1} + \dots + u_n$ by y , we conclude that

$$\begin{aligned} & \int_{\mathbb{R}^n} \text{Meas}(u_1; \dots; u_n) \prod_{i=1}^k r_j(u_i) \prod_{i=k+1}^n s_j(u_i) du_1 \dots du_n \\ &= \int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^k} \text{Meas}(u_1; \dots; u_k) dy \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(y); \end{aligned} \quad (33)$$

Step 2: bound on the integral (33). We represent the expression (33) as a sum of two terms:

$$\begin{aligned} & \int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^k} \text{Meas}(u_1; \dots; u_k) dy \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(y) \\ &= \int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^k} \text{Meas}(u_1; \dots; u_k) dy \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(y) \\ &+ \int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^k} \text{Meas}(u_1; \dots; u_k) dy \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(y); \end{aligned}$$

and bound the former and the latter summands in the right-hand side separately. According to Lemma A.5 and A.6, it holds that

$$\int_0^1 v_1 dv_1 \dots \int_0^1 v_k dv_k \int_{\mathbb{R}^k} \text{Meas}(u_1; \dots; u_k) dy \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} s_j^{(n-k)}(y) \leq 6 \frac{p}{2} \frac{(2k-1)!}{2} \frac{2^{3+2j} B_j}{n^{j+1} j! j^{1+j}} + \exp \left(-\frac{t^2}{64} \right) \frac{M_j}{2^{2n^2}} ;$$

and

$$\int_0^{Z^1} \int_0^{Z^1} \int_0^Z \int_0^{Z^1} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \, du_1 \cdots du_k \, dy \quad (2k) \quad t \quad X^k \quad \frac{v_i u_i}{n} \quad \frac{y}{n} \quad Y^k \quad \frac{u_i^2 r_j(u_i)}{n^2} \quad \zeta_j^{(n-k)}(y)$$

$$6 \frac{p}{(2k-2)!} \frac{M_j}{2^{2n^2}}^k \frac{4k^2 - 2}{t^2} + p \frac{1}{2} \frac{2k}{t} :$$

Hence, we obtain that

$$\int_0^{Z^1} \int_0^{Z^1} \int_0^Z \int_0^{Z^1} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \, du_1 \cdots du_k \, dy \quad (2k) \quad t \quad X^k \quad \frac{v_i u_i}{n} \quad \frac{y}{n} \quad Y^k \quad \frac{u_i^2 r_j(u_i)}{n^2} \quad \zeta_j^{(n-k)}(y)$$

$$6 \frac{p}{(2k-1)!} \frac{2^{3+2j} B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} \frac{M_j}{2^{2n^2}}^k \quad (34)$$

$$+ p \frac{(2k-2)!}{(2k-2)!} \frac{M_j}{2^{2n^2}}^k \frac{4k^2 - 2}{t^2} + p \frac{1}{2} \frac{2k}{t} :$$

Step 3: final bound. Summing up the equalities (31), (33), and the equality (34), we obtain that

$$F(t) - G(t) \leq 6 \sum_{k=1}^X \frac{n^p}{k} \frac{(2k-1)!}{p^2} \frac{2^{3+2j} B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} \frac{M_j}{2^{2n^2}}^k$$

$$+ \sum_{k=1}^X \frac{n^p}{k} \frac{(2k-2)!}{(2k-2)!} \frac{4k^2 - 2}{t^2} + p \frac{1}{2} \frac{2k}{t} \frac{M_j}{2^{2n^2}}^k$$

$$+ \sum_{k=1}^X \frac{n^p}{k} \frac{p}{(2k)!} \frac{B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} + \frac{k^2}{t^2} + \frac{1}{t} \frac{M_j}{2^{2n^2}}^k :$$

Finally, Lemma A.7 implies that

$$F(t) - G(t) \leq \sum_{k=1}^X \frac{n^p}{k} \frac{p}{(2k)!} \frac{B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} + \frac{k^2}{t^2} + \frac{1}{t} \frac{M_j}{2^{2n^2}}^k$$

$$6 \frac{B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} + \frac{2^2}{t^2} + \frac{1}{t} \frac{2M_j}{2^n}$$

$$+ \frac{M_j}{nt} \left(1 + \frac{1}{t} + \frac{B_j}{n^{j+1} j t^{j+1}} \right) ;$$

whenever $2M_j \leq 6n^2$. The proof is finished.

Lemma A.5. Under Assumption 2.1, it holds that

$$\int_0^{Z^1} \int_0^{Z^1} \int_0^Z \int_0^{Z^1} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \, du_1 \cdots du_k \, dy \quad (2k) \quad t \quad X^k \quad \frac{v_i u_i}{n} \quad \frac{y}{n} \quad Y^k \quad \frac{u_i^2 r_j(u_i)}{n^2} \quad \zeta_j^{(n-k)}(y)$$

$$6 \frac{p}{(2k-1)!} \frac{2^{3+2j} B_j}{n^{j+1} j t^{j+1}} + \exp \frac{t^2}{64^2} \frac{M_j}{2^{2n^2}}^k :$$

Proof. To prove Lemma A.5, it is enough to show that

$$\int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}; \frac{nt}{2k}]^k} du_1 \dots du_k \int_0^1 dy \stackrel{(2k)}{=} t \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y) \\ \geq \frac{p}{6} \frac{(2k-1)!}{p^{\frac{2k-1}{2}}} \frac{2^{3+2j} B_j p}{n \prod_{i=1}^k j! t^{j+1}} + \exp \left(- \frac{t^2}{64} \frac{M_j}{2^2 n^2} \right) :$$

Our argument is quite technical, so we divide the proof into several parts.

Step 1: a bound on the $2k$ -th derivative. First, let us consider the $2k$ -th derivative of $\xi_j^{(n-k)}(y)$. Note that

$$\frac{d^{2k}}{dy^{2k}} \xi_j^{(n-k)}(y) = \frac{1}{2k} \frac{d^{2k-1}}{dy^{2k-1}} \xi_j^{(n-k)}(y) = \frac{1}{2k} \frac{p}{2} H_{2k-1}(w) \exp \left(- \frac{w^2}{2} \right) ;$$

where H_{2k-1} is the $(2k-1)$ -th probabilist's Hermite polynomial. We provide a brief information about Hermite polynomials in Appendix D. In particular, we refer to the result of Indritz (1961), which implies that

$$\max_{w \in \mathbb{R}} H_{2k-1}(w) \exp \left(- \frac{w^2}{4} \right) \leq \frac{p}{6} \frac{(2k-1)!}{2} \quad \text{for all } k \geq 2.$$

Thus, it holds that

$$\frac{d^{2k}}{dy^{2k}} \xi_j^{(n-k)}(y) \leq \frac{p}{6} \frac{(2k-1)!}{2k} \exp \left(- \frac{w^2}{4} \right) \quad \text{for all } w \in \mathbb{R} \text{ and all } k, \tag{35}$$

and we obtain the inequality

$$\int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}; \frac{nt}{2k}]^k} du_1 \dots du_k \int_0^1 dy \stackrel{(2k)}{=} t \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y) \\ \geq \frac{p}{6} \frac{(2k-1)!}{2k} \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}; \frac{nt}{2k}]^k} du_1 \dots du_k \int_0^1 dy \exp \left(- \frac{1}{4} \frac{t^2}{n^2} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \right) \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y) ; \tag{36}$$

Step 2: a bound on the convolution. Our next goal is to bound the convolution

$$\int_0^1 \exp \left(- \frac{(w-y)^2}{4} \right) t \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y) dy$$

using the properties of $\xi_j^{(n-k)}$ from Assumption 2.1. Let us fix an arbitrary $w \in \mathbb{R}$ and consider

$$\int_0^1 \exp \left(- \frac{(w-y)^2}{4} \right) \xi_j^{(n-k)}(y) dy :$$

Since $\exp \left(- \frac{(w-y)^2}{4} \right) \geq \frac{1}{6}$, it holds that

$$\int_0^1 \exp \left(- \frac{(w-y)^2}{4} \right) \xi_j^{(n-k)}(y) dy \geq \frac{1}{6} \int_0^1 \xi_j^{(n-k)}(y) dy \geq \frac{1}{6} ; \tag{37}$$

On the other hand, we have

$$\int_1^{\bar{z}^1} \exp\left(-\frac{(w-y)^2}{4}\right) \zeta_j^{(n-k)}(y) dy = n \int_1^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) \zeta_j^{(n-k)}(ny+nw) dy$$

$$= \int_{w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) \frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw+yj^{1+j}}} dy$$

$$+ n \int_{|y|>w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) \zeta_j^{(n-k)}(ny+nw) dy:$$

If $|y| \geq w=2$, then

$$\frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw+yj^{1+j}}} \geq \frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw=2j^{1+j}}}$$

and it holds that

$$\int_{w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) \frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw+yj^{1+j}}} dy$$

$$\geq \frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw=2j^{1+j}}} \int_{w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) dy \tag{38}$$

$$\geq \frac{B_j n(n-k)}{(n-k)^{(j+1)=j+n^{1+j}jw=2j^{1+j}}} \frac{1}{4} \int_{w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) dy:$$

Otherwise,

$$n \int_{|y|>w=2}^{\bar{z}^1} \exp\left(-\frac{y^2}{4}\right) \zeta_j^{(n-k)}(ny+nw) dy \geq \exp\left(-\frac{w^2}{16}\right) \int_1^{\bar{z}^1} \zeta_j^{(n-k)}(ny+nw) ndy$$

$$= \exp\left(-\frac{w^2}{16}\right) : \tag{39}$$

Taking into account (37), (38), and (39), we obtain that

$$\int_1^{\bar{z}^1} \exp\left(-\frac{(w-y)^2}{4}\right) \zeta_j^{(n-k)}(y) dy \geq \frac{1}{n} \left[\frac{2^{2+j} B_j}{n^{j+1} j w^{j+1}} + \exp\left(-\frac{w^2}{16}\right) \right]$$

and, hence,

$$\int_1^{\bar{z}^1} \exp\left(-\frac{1}{4} t \sum_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n}\right) \zeta_j^{(n-k)}(y) dy$$

$$\geq \min\left\{1, \frac{2^{2+j} B_j}{n^{j+1}} t \sum_{i=1}^k \frac{v_i u_i}{n} + \exp\left(-\frac{1}{16} t \sum_{i=1}^k \frac{v_i u_i}{n}\right)\right\} \tag{40}$$

Step 3: nal bound. The inequalities (36) and (40) yield that

$$\int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}, \frac{nt}{2k}]^k} du_1 \dots du_k \int_1^y dy \stackrel{(2k)}{t} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \mathfrak{S}_j^{(n-k)}(y)$$

$$\leq \frac{p}{2k} \frac{(2k-1)!}{2^p} \int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}, \frac{nt}{2k}]^k} du_1 \dots du_k \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2}$$

$$\min_{j \in \mathcal{I}} \left\{ 1; \frac{2^{2+j} B_j}{n^{j-1}} t \prod_{i=1}^k \frac{v_i u_i}{n} + \exp \left\{ -\frac{1}{16} t \prod_{i=1}^k \frac{v_i u_i}{n} \right\} \right\} \frac{1}{2^{9+1}} \frac{1}{A_j} :$$

On the set $[\frac{nt}{2k}; \frac{nt}{2k}]^k$, we have

$$t \prod_{i=1}^k \frac{v_i u_i}{n} > \frac{t}{2};$$

and, hence,

$$\int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}, \frac{nt}{2k}]^k} du_1 \dots du_k \int_1^y dy \stackrel{(2k)}{t} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \mathfrak{S}_j^{(n-k)}(y)$$

$$\leq \frac{p}{2k} \frac{(2k-1)!}{2^p} \frac{2^{3+2+j} B_j}{n^{j-1} j!^{1+j}} + \exp \left\{ -\frac{t^2}{64} \int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{[\frac{nt}{2k}, \frac{nt}{2k}]^k} du_1 \dots du_k \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \right\} :$$

Due to Assumption 2.1, the right-hand side does not exceed

$$\frac{p}{2k} \frac{(2k-1)!}{2^p} \frac{2^{3+2+j} B_j}{n^{j-1} j!^{1+j}} + \exp \left\{ -\frac{t^2}{64} \int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 du_1 \dots du_k \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \right\}$$

$$= \frac{p}{2} \frac{(2k-1)!}{2^p} \frac{2^{3+2+j} B_j}{n^{j-1} j!^{1+j}} + \exp \left\{ -\frac{t^2}{64} \frac{M_j}{2^{2n^2}} \right\} :$$

The proof is finished.

Lemma A.6. Let Assumption 2.1 be fulfilled. Then, for any $k > 2$, it holds that

$$\int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{R^k} du_1 \dots du_k \int_1^y dy \stackrel{(2k)}{t} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \mathfrak{S}_j^{(n-k)}(y)$$

$$\leq \frac{p}{2} \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \frac{4k^2-2}{t^2} + \frac{1}{2} \frac{2k}{t} :$$

Proof. Let $\mathcal{A}_i : 1 \leq i \leq k$ be a collection of sets in R^k , such that

$\mathcal{A}_1, \dots, \mathcal{A}_k$ form a partition of $R^k \setminus [\frac{nt}{2k}, \frac{nt}{2k}]^k$, that is,

$$R^k \setminus [\frac{nt}{2k}, \frac{nt}{2k}]^k = \bigcup_{i=1}^k \mathcal{A}_i \quad \text{and} \quad \mathcal{A}_i \cap \mathcal{A}_\ell = \emptyset; \quad \text{for all } i, \ell \in \mathcal{I};$$

for all $i \in \{1, \dots, k\}$, it holds that

$$A_i \quad u_i : |u_i| > \frac{nt}{2k} :$$

Let us fix any $i \in \{1, \dots, k\}$ and consider the integral

$$\int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{A_i} du_1 \dots du_k \int_0^1 dy \quad (2k) \quad \int_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \int_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y):$$

Applying the Newton-Leibnitz formula, we obtain that

$$\begin{aligned} \frac{u_i^2}{n^2} \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{A_i} du_1 \dots du_k \int_0^1 dy \quad (2k) \quad \int_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} v_i dv_i &= \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 \int_{A_i} \frac{v_i u_i}{n} \frac{y}{n} \frac{C}{A} \\ &+ \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 \int_{A_i} \frac{v_i u_i}{n} \frac{y}{n} \frac{C}{A} \frac{u_i}{n} : \end{aligned}$$

This implies that

$$\begin{aligned} \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{A_i} du_1 \dots du_k \int_0^1 dy \quad (2k) \quad \int_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \int_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \xi_j^{(n-k)}(y) \\ \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{A_i} du_1 \dots du_k \int_0^1 dy \quad (2k-2) \quad \int_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \frac{C}{A} \\ \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 \int_{A_i} \frac{u_i^2 r_j(u_i)}{n^2} \frac{C}{A} r_j(u_i) \xi_j^{(n-k)}(y) \\ + \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 v_1 dv_1 \dots v_k dv_k \int_{A_i} du_1 \dots du_k \int_0^1 dy \quad (2k-1) \quad \int_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \frac{C}{A} \\ \int_0^1 \int_0^1 \dots \int_0^1 \int_0^1 \int_{A_i} \frac{u_i^2 r_j(u_i)}{n^2} \frac{C}{A} \frac{u_i r_j(u_i)}{n} \xi_j^{(n-k)}(y) \end{aligned}$$

Let us apply the inequality (35) we derived in the proof of Lemma A.5 to $(2k-2)$ and $(2k-1)$: for all $w \in \mathbb{R}$, it holds that

$$0 \leq (w) \leq 1; \quad (2k-2)(w) \leq \frac{p \frac{(2k-3)!}{2k-2} \exp \frac{w^2}{4}}{2} \leq \frac{p \frac{(2k-2)!}{2k-2}}{2}; \quad \text{for all } k > 2.$$

and

$$(2k-1)(w) \leq \frac{p \frac{(2k-2)!}{2k-1} \exp \frac{w^2}{4}}{2} \leq \frac{p \frac{(2k-2)!}{2k-1}}{2}.$$

Then, for any $k \geq 2$, N ,

$$\int_0^1 \int_0^1 \dots \int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 du_1 \dots du_k \int_0^1 dy \int_0^1 \frac{v_i u_i}{n} \frac{y}{n} \dots$$

$$\int_0^1 \frac{u_i^2 r_j(u_i)}{n^2} r_j(u) \xi^{(n-k)}(y)$$

$$\int_0^1 \frac{(2k-2)!}{2^{2k-2}} \int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 \frac{u_i^2 r_j(u_i)}{n^2} r_j(u) du_1 \dots du_k$$

Due to Assumption 2.1,

$$\int_0^1 \frac{(2k-2)!}{2^{2k-2}} \int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 \frac{u_i^2 r_j(u_i)}{n^2} r_j(u) du_1 \dots du_k$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \int_0^1 r_j(u) du$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \frac{2k^2}{n^2 t^2} \int_0^1 u^2 r_j(u) du$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \frac{2k^2 M_j}{n^2 t^2}$$

Similarly, it holds that

$$\int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 du_1 \dots du_k \int_0^1 dy \int_0^1 \frac{v_i u_i}{n} \frac{y}{n} \dots$$

$$\int_0^1 \frac{u_i^2 r_j(u_i)}{n^2} \frac{u_i r_j(u)}{n} \xi^{(n-k)}(y)$$

$$\int_0^1 \frac{(2k-2)!}{2^{2k-1} 2} \int_0^1 v_1 dv_1 \dots v_k dv_k \int_0^1 \frac{u_i^2 r_j(u_i)}{n^2} \frac{u_i r_j(u)}{n} du_1 \dots du_k$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \int_0^1 \frac{u_i r_j(u)}{n} du$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \frac{k}{n^2 t} \int_0^1 u^2 r_j(u) du$$

$$\int_0^1 \frac{(2k-2)!}{2^{2n^2}} \frac{M_j}{2^{2n^2}} \frac{k M_j}{n^2 t}$$

Thus, we obtain that

$$\int_0^1 \int_0^1 \dots \int_0^1 \int_{\mathbb{A}} \int_{\mathbb{R}^k} \int_{\mathbb{R}^k} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \frac{M_j}{2 \cdot 2n^2} \frac{4k^2 - 2}{t^2} + \frac{1}{2} \frac{(2k-2)!}{2} \frac{M_j}{2 \cdot 2n^2} \frac{2k}{t} :$$

Hence, due to the triangle inequality,

$$\begin{aligned} & \int_0^1 \int_0^1 \dots \int_0^1 \int_{\mathbb{A}} \int_{\mathbb{R}^k} \int_{\mathbb{R}^k} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \\ & \int_0^1 \int_0^1 \dots \int_0^1 \int_{\mathbb{A}} \int_{\mathbb{R}^k} \int_{\mathbb{R}^k} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \\ & \int_0^1 \int_0^1 \dots \int_0^1 \int_{\mathbb{A}} \int_{\mathbb{R}^k} \int_{\mathbb{R}^k} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \\ & \int_0^1 \int_0^1 \dots \int_0^1 \int_{\mathbb{A}} \int_{\mathbb{R}^k} \int_{\mathbb{R}^k} \prod_{i=1}^k \frac{v_i u_i}{n} \frac{y}{n} \prod_{i=1}^k \frac{u_i^2 r_j(u_i)}{n^2} \zeta_j^{(n-k)}(y) \end{aligned}$$

Lemma A.7. Let $n \geq N$ and $a > 0$ be such that $na \geq 1$. Then it holds that

$$\sum_{k=1}^n \frac{n!}{k!} \frac{a^k}{2} \leq 2na \quad \text{and} \quad \sum_{k=1}^n \frac{n!}{k!} \frac{a^k}{2} \leq 4na:$$

Proof. First, note that, for any positive integer k , we have

$$(2k)! = \prod_{j=1}^k (2j) \cdot \prod_{j=1}^k (2j-1) \leq \prod_{j=1}^k (2j) \cdot \prod_{j=1}^k (2j) \leq 4^k k! k!$$

This implies that

$$\sum_{k=1}^n \frac{n!}{k!} \frac{a^k}{2} \leq \sum_{k=1}^n \frac{n!}{(n-k)!} a^k \leq \sum_{k=1}^n (na)^k \leq (na)^n = \frac{na}{1-na} \leq 2na:$$

Similarly, it holds that

$$\sum_{k=1}^n \frac{n!}{k!} \frac{a^k}{2} \leq \sum_{k=1}^n \frac{n!}{(n-k)!} a^k \leq \sum_{k=1}^n k(na)^k \leq na \sum_{k=0}^n (na)^k - 1 = \frac{na}{(1-na)^2} \leq 4na:$$

B Proofs for clipped-SGD

B.1 Auxiliary Results

Bernstein inequality. The following lemma (known as Bernstein inequality for martingale differences (Bennett, 1962; Dzhaparidze and Van Zanten, 2001; Freedman et al., 1975)) is essential for deriving high-probability upper bounds in our analysis.

Lemma B.1. Let the sequence of random variables $\{X_i\}_{i>1}$ form a martingale difference sequence, i.e. $E[X_i | X_1, \dots, X_{i-1}] = 0$ for all $i > 1$. Assume that conditional variances $\sigma_i^2 \stackrel{\text{def}}{=} E[X_i^2 | X_1, \dots, X_{i-1}]$ exist and are bounded and assume also that there exists deterministic constant $\alpha > 0$ such that $|X_i| \leq \alpha$ almost surely for all $i > 1$. Then for all $b > 0$, $G > 0$ and $n > 1$

$$P \left(\sum_{i=1}^n X_i > b \text{ and } \sum_{i=1}^n \sigma_i^2 \leq G \right) \leq 2 \exp \left(-\frac{b^2}{2G + \frac{2cb^3}{3}} \right) \quad (41)$$

Bias and variance of the clipped stochastic vector. We also rely on the following result from (Gorbunov et al., 2020).

Lemma B.2 (Simplified version of Lemma F.5 from (Gorbunov et al., 2020)). Let X be a random vector in \mathbb{R}^d and $\mathcal{X} = \text{clip}(X; \cdot)$. Then,

$$E[\mathcal{X}] = E[X] \quad (42)$$

Moreover, if for some $\alpha > 0$

$$E[X] = x \in \mathbb{R}^d; \quad E[\|X - x\|^2] \leq \alpha^2 \quad (43)$$

and $\|x\| \leq \alpha$, then

$$E[\mathcal{X}] = x \leq \frac{4}{3} \alpha^2; \quad (44)$$

$$E[\|\mathcal{X} - E[\mathcal{X}]\|^2] \leq 18 \alpha^2; \quad (45)$$

B.2 Quasi-Convex Case

The analysis of clipped-SGD in the quasi-convex case relies on the following lemma from (Sadiev et al., 2023).

Lemma B.3. Let Assumptions 2.2 and 2.4 with $\beta = 0$ hold on $Q = B_{2R}(x^0)$, where $R > \kappa x^0$, and let stepsize η satisfy $\eta \leq \frac{1}{L}$. If $x^k \in Q$ for all $k = 0, 1, \dots, K+1$, $K > 0$, then after K iterations of clipped-SGD we have

$$f(\bar{x}^K) - f(x^0) \leq \frac{\kappa x^0 - \eta \kappa^2 x^{K+1} - \eta \kappa^2}{K+1} + \frac{2}{K+1} \sum_{k=0}^K \eta \kappa^k \|x^k - r f(x^k)\| + \frac{2}{K+1} \sum_{k=0}^K \eta \kappa^k \kappa^2; \quad (46)$$

$$\bar{x}^K = \frac{1}{K+1} \sum_{k=0}^K x^k; \quad (47)$$

$$\kappa \stackrel{\text{def}}{=} \text{clip}(r f(x^k); \kappa) - r f(x^k); \quad (48)$$

Theorem B.4. Let Assumptions 2.2 and 2.4 with $\beta = 0$ hold on $Q = B_{2R}(x^0)$, where $R > \kappa x^0$. Assume

⁴Although Sadiev et al. (2023) claim that they use Assumption 2.3 with $\beta = 0$, their proof relies on Assumption 2.4 with $\beta = 0$ instead, which is strictly weaker.

that $r f_k(x^k)$ satisfies Assumption 5.1 with parameters b_k ; k for $k = 0; 1; \dots; K$, $K > 0$ and

$$\min_k \left\{ \frac{1}{160L \ln^{4(K+1)}}; \frac{qR}{208k \ln^{4(K+1)}}; \frac{R}{160b_k \ln^{4(K+1)}}; \frac{R}{1600b_k(K+1)} \right\}; \quad (49)$$

$$k = \frac{R}{40 \ln^{4(K+1)}}; \quad (50)$$

for some $\alpha \in (0; 1]$. Then, after K iterations of clipped-SGD the iterates with probability at least $1 - \alpha$ satisfy

$$f(x^K) - f(x^*) \leq \frac{2R^2}{(K+1)} \quad \text{and} \quad \|x^K - x^*\| \leq B_{2R} \alpha; \quad (51)$$

In particular, when

$$\min_k \left\{ \frac{1}{160L \ln^{4(K+1)}}; \frac{qR}{208k \ln^{4(K+1)}}; \frac{R}{160b_k \ln^{4(K+1)}}; \frac{R}{1600b_k(K+1)} \right\}; \quad (52)$$

$$\text{where } k = \min_{k=0;1;\dots;K} k; \quad b = \min_{k=0;1;\dots;K} b_k; \quad (53)$$

then the iterates produced by clipped-SGD after K iterations with probability at least $1 - \alpha$ satisfy

$$f(x^K) - f(x^*) = O\left(\max\left\{ \frac{LR^2 \ln K}{K}; \frac{R}{K}; \frac{bR \ln K}{K}; bR \right\}\right); \quad (54)$$

Proof. Our proof follows similar steps to the one given by [Sadiev et al. \(2023\)](#). The main difference comes due to the presence of the bias in $r f_k(x^k)$. Therefore, for completeness, we provide the full proof here.

Let $R_k = kx^k - x^k$ for all $k > 0$. Our next objective is to establish, by induction, that $\|R_k\| \leq 2R$ with a high probability. This will enable us to apply the result from Lemma B.3 and subsequently utilize Bernstein's inequality to estimate the stochastic component of the upper bound. To be more precise, for each $k = 0; \dots; K+1$, we consider the probability event E_k , defined as follows: inequalities

$$\sum_{l=0}^{k-1} \|x^l - x^*\| \leq R; \quad \sum_{l=0}^{k-1} \|k_l\|^2 \leq R^2; \quad (55)$$

$$\|R_t\| \leq 2R \quad (56)$$

hold for all $t = 0; 1; \dots; k$ simultaneously. We aim to demonstrate through induction that $\mathbb{P}(E_k) > 1 - \alpha^{k/(K+1)}$ for all $k = 0; 1; \dots; K+1$. The base case $k = 0$, is trivial. Assuming that the statement holds for some $k = T - 1 \leq K$, specifically, $\mathbb{P}(E_{T-1}) > 1 - \alpha^{(T-1)/(K+1)}$, we need to establish that $\mathbb{P}(E_T) > 1 - \alpha^{T/(K+1)}$.

To begin, we observe that the probability event E_{T-1} implies that $x_t \in B_{2R}(x^*)$ for all $t = 0; 1; \dots; T-1$. Furthermore, E_{T-1} implies that

$$\|x^T - x^*\| = \|x^T - x^{T-1} + x^{T-1} - x^*\| \leq \|x^T - x^{T-1}\| + \|x^{T-1} - x^*\| \leq \frac{R}{2} + 2R; \quad (50)$$

i.e., $x^0; x^1; \dots; x^T \in B_{2R}(x^*)$. Hence, with E_{T-1} implying $\|x^k - x^*\| \leq 2R$, we confirm that the conditions of Lemma B.3 are met, resulting in

$$f(x^{T-1}) - f(x^*) \leq \frac{\sum_{t=0}^{T-1} \|x^t - x^*\|^2}{T} + \frac{2}{T} \sum_{l=0}^{T-1} \|x^l - x^*\| \leq \frac{2R^2 T}{T} + \frac{2}{T} \sum_{l=0}^{T-1} 2R \leq 2R + \frac{4R}{T}; \quad (57)$$

for all $t = 1; \dots; T$ simultaneously and for all $t = 1; \dots; T - 1$ this probability event also implies that

$$f(x^{t-1}) - f(x^*) \leq \frac{1}{t} R^2 + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, x^t - x^* \rangle \mid \mathcal{F}_t \right] + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\|k^l\|^2 \mid \mathcal{F}_t \right] \stackrel{(55)}{\leq} \frac{2R^2}{t}. \quad (58)$$

Considering that $f(x^{T-1}) - f(x^*) > 0$, we can further deduce from (57) that when E_{T-1} holds, the following holds as well:

$$R^2 \leq R^2 + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, x^T - x^* \rangle \mid \mathcal{F}_T \right] + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\|k^l\|^2 \mid \mathcal{F}_T \right]. \quad (59)$$

Next, we define random vectors

$$z_t = \begin{cases} x^t - x^* - r f(x^t); & \text{if } \|x^t - x^* - r f(x^t)\| \leq 2R; \\ 0; & \text{otherwise;} \end{cases}$$

for all $t = 0; 1; \dots; T - 1$. As per their definition, these random vectors are bounded with probability 1

$$\|z_t\| \leq 2R. \quad (60)$$

Moreover, for $t = 0; \dots; T - 1$ event E_{T-1} implies

$$\|k^t f(x^t)\| \stackrel{(3)}{\leq} L \|x^t - x^*\| \stackrel{(56)}{\leq} p \frac{R}{2L} \stackrel{(49);(50)}{\leq} \frac{R}{2}; \quad (61)$$

$$\mathbb{E} \left[\|r f(x^t)\| \mid \mathcal{F}_t \right] \leq \mathbb{E} \left[\|r f(x^t)\| \mid \mathcal{F}_t \right] + \|k^t f(x^t)\| \stackrel{(7);(61)}{\leq} b_t + p \frac{R}{2L} \stackrel{(49);(50)}{\leq} \frac{R}{2}; \quad (62)$$

$$\|x^t - x^* - r f(x^t)\| \leq \|x^t - x^*\| + \|k^t f(x^t)\| \stackrel{(61)}{\leq} p \frac{R}{2L} (1 + L) \stackrel{(49)}{\leq} 2R;$$

The latter inequality means that E_{T-1} implies $z_t = x^t - x^* - r f(x^t)$ for $t = 0; \dots; T - 1$. Next, we define the unbiased part and the bias of z_t as z_t^u and z_t^b , respectively:

$$z_t^u = \text{clip} \left(r f_t(x^t); z_t \right) - \mathbb{E} \left[\text{clip} \left(r f_t(x^t); z_t \right) \mid \mathcal{F}_t \right]; \quad z_t^b = \mathbb{E} \left[\text{clip} \left(r f_t(x^t); z_t \right) \mid \mathcal{F}_t \right] - r f(x^t). \quad (63)$$

We notice that $z_t = z_t^u + z_t^b$. Using new notation, we get that E_{T-1} implies

$$\begin{aligned} R^2 &\leq R^2 + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^u \rangle \mid \mathcal{F}_t \right] + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^b \rangle \mid \mathcal{F}_t \right] + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\|k^l\|^2 \mid \mathcal{F}_t \right] + \mathbb{E} \left[\|z_t^u\|^2 \mid \mathcal{F}_t \right] \\ &\quad + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\|k^l\|^2 \mid \mathcal{F}_t \right] + 2 \sum_{l=0}^{X-1} \mathbb{E} \left[\|z_t^b\|^2 \mid \mathcal{F}_t \right]. \end{aligned} \quad (64)$$

To conclude our inductive proof successfully, we must obtain sufficiently strong upper bounds with high probability for the terms $\sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^u \rangle \mid \mathcal{F}_t \right]$; $\sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^b \rangle \mid \mathcal{F}_t \right]$; $\sum_{l=0}^{X-1} \mathbb{E} \left[\|z_t^u\|^2 \mid \mathcal{F}_t \right]$; $\sum_{l=0}^{X-1} \mathbb{E} \left[\|z_t^b\|^2 \mid \mathcal{F}_t \right]$. In other words, we need to demonstrate that $\sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^u \rangle \mid \mathcal{F}_t \right] + \sum_{l=0}^{X-1} \mathbb{E} \left[\langle h^l, z_t^b \rangle \mid \mathcal{F}_t \right] + \sum_{l=0}^{X-1} \mathbb{E} \left[\|z_t^u\|^2 \mid \mathcal{F}_t \right] + \sum_{l=0}^{X-1} \mathbb{E} \left[\|z_t^b\|^2 \mid \mathcal{F}_t \right] \leq R^2$ with a high probability. In the subsequent stages of the proof, we will rely on the bounds for the norms and second moments of z_t^u and z_t^b . First, as per the definition of the clipping operator, we can assert with probability 1 that

$$\|z_t^u\| \leq 2; \quad (65)$$

Furthermore, given that E_{T-1} implies $\|k^t f(x^t)\| \leq R$ for $t = 0; 1; \dots; T - 1$ (as per (62)), then, according to Lemma B.2, we can deduce that E_{T-1} implies

$$\|z_t^b\| \leq \mathbb{E} \left[\text{clip} \left(r f_t(x^t); z_t \right) \mid \mathcal{F}_t \right] - r f(x^t) \stackrel{(44);(7)}{\leq} \frac{4}{6} + b_t; \quad (66)$$

$$\mathbb{E} \left[\|z_t^u\|^2 \mid \mathcal{F}_t \right] \stackrel{(45)}{\leq} 18. \quad (67)$$

Upper bound for -- . By definition of h_t^u , we readily observe that $E_t[h_t^u] = 0$ and

$$E_t[2h_t^u; i] = 0:$$

Next, the sum -- contains only the terms that are bounded with probability 1:

$$|2h_t^u; i| \leq 2k_t^u k_t k_t \stackrel{(60);(65)}{\leq} \frac{R^2}{6} \stackrel{(50)}{\leq} \frac{R^2}{5 \ln^{4(K+1)}} \stackrel{\text{def}}{=} c: \quad (68)$$

The conditional variances $e_t^2 \stackrel{\text{def}}{=} E_t[4h_t^u; i^2]$ of the summands are bounded:

$$e_t^2 \leq 4 \cdot 2k_t^u k_t^2 \stackrel{(60)}{\leq} \frac{16}{6} R^2 E_t[k_t^u k_t^2]: \quad (69)$$

To summarize, we have demonstrated that $2h_t^u; i g_{t=0}^T$ is a bounded martingale difference sequence with bounded conditional variances $e_t^2 g_{t=0}^T$. Therefore, one can apply Bernstein's inequality (Lemma B.1) with $X_t = 2h_t^u; i$, parameter c as in (68), $b = \frac{R^2}{5}$, $G = \frac{R^4}{150 \ln^{4(K+1)}}$ and get

$$\left(P_{j \rightarrow j} > \frac{R^2}{5} \text{ and } \sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln^{4(K+1)}} \right) \leq 2 \exp \left(-\frac{b^2}{2G + 2cb} \right) = \frac{b^2}{2(K+1)};$$

which is equivalent to

$$P\{E_{-g} > 1\} \leq \frac{b^2}{2(K+1)}; \text{ for } E_{-g} = \left(\sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln^{4(K+1)}} \text{ or } P_{j \rightarrow j} > \frac{R^2}{5} \right): \quad (70)$$

Additionally, event E_{T-1} implies that

$$\sum_{t=0}^{T-1} e_t^2 \stackrel{(69)}{\leq} \frac{16}{6} R^2 \sum_{t=0}^{T-1} E_t[k_t^u k_t^2] \stackrel{(67)}{\leq} \frac{16}{6} R^2 \cdot 2T \stackrel{(49)}{\leq} \frac{R^4}{150 \ln^{4(K+1)}}: \quad (71)$$

Upper bound for -- . From E_{T-1} it follows that

$$\begin{aligned} \text{--} &= \sum_{t=0}^{T-1} 2h_t^b; i \leq 2 \sum_{t=0}^{T-1} k_t^b k_t k_t \stackrel{(60);(66)}{\leq} \frac{16}{6} T R + 4 R b T \\ &\stackrel{(50)}{=} \frac{2}{5} \sum_{t=0}^{T-1} 2T \ln^{4(K+1)} + 4 R b T \stackrel{(49)}{\leq} \frac{R^2}{5}: \end{aligned} \quad (72)$$

Upper bound for ⓐ . By construction, we have

$$E_t[2^2 k_t^u k_t^2] - E_t[k_t^u k_t^2] = 0:$$

Next, the sum ⓐ contains only the terms that are bounded with probability 1:

$$\begin{aligned} 2^2 k_t^u k_t^2 - E_t[k_t^u k_t^2] &\leq 2^2 k_t^u k_t^2 + E_t[k_t^u k_t^2] \\ &\stackrel{(65)}{\leq} \frac{16}{6} \cdot 2 \stackrel{(50)}{\leq} \frac{R^2}{100 \ln^{4(K+1)}} \leq \frac{R^2}{5 \ln^{4(K+1)}} \stackrel{\text{def}}{=} c: \end{aligned} \quad (73)$$

The conditional variances $e_t^2 \stackrel{\text{def}}{=} E_t[4^4 k_t^u k_t^2 - E_t[k_t^u k_t^2]^2]$ of the summands are bounded:

$$e_t^2 \stackrel{(73)}{\leq} \frac{R^2}{5 \ln^{4(K+1)}} E_t[2^2 k_t^u k_t^2 - E_t[k_t^u k_t^2]^2] \leq \frac{4}{5 \ln^{4(K+1)}} R^2 E_t[k_t^u k_t^2]: \quad (74)$$

To summarize, we have demonstrated that $\sum_{t=0}^n 2^{-2t} k_t^u k^2 E_t k_t^u k^2$ is a bounded martingale difference sequence with bounded conditional variances $E_t k_t^u k^2 = 1$. Therefore, one can apply Bernstein's inequality (Lemma B.1) with $X_t = 2^{-2t} k_t^u k^2 E_t k_t^u k^2$, parameter c as in (73), $b = \frac{R^2}{5}$, $G = \frac{R^4}{150 \ln^{4(K+1)}}$ and get

$$P(j \circledast j) > \frac{R^2}{5} \text{ and } \sum_{t=0}^{T-1} e_t^2 \leq \frac{R^4}{150 \ln^{4(K+1)}} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb}\right) = \frac{b^2}{2(K+1)};$$

which is equivalent to

$$P(E_{\circledast} > 1) \leq \frac{b^2}{2(K+1)}; \text{ for } E_{\circledast} = \left(\sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln^{4(K+1)}} \text{ or } j \circledast j \leq \frac{R^2}{5} \right); \quad (75)$$

Additionally, event E_{T-1} implies that

$$\sum_{t=0}^{T-1} e_t^2 \stackrel{(74)}{\leq} \frac{4 \cdot 2R^2}{5 \ln^{4(K+1)}} \sum_{t=0}^{T-1} E_t k_t^u k^2 \stackrel{(67)}{\leq} \frac{72 \cdot 2R^2 \cdot 2T}{5 \ln^{4(K+1)}} \stackrel{(49)}{\leq} \frac{R^4}{150 \ln^{4(K+1)}}; \quad (76)$$

Upper bound for \ominus . From E_{T-1} it follows that

$$\ominus = \sum_{t=0}^{T-1} 2^{-2t} \sum_{i=0}^h E_t k_t^u k^2 \stackrel{(67)}{\leq} \frac{36 \cdot 2^2 T}{6} \stackrel{(49)}{\leq} \frac{R^2}{5}; \quad (77)$$

Upper bound for \circ . From E_{T-1} it follows that

$$\circ = \sum_{t=0}^{T-1} 2^{-2t} \sum_{i=0}^h b^2 \stackrel{(66)}{\leq} \frac{32 \cdot 4T^2}{2} + 2 \cdot 2b^2 T \stackrel{(50)}{=} 51200 \frac{4T^2 \ln^{4(K+1)}}{R^2} + 2 \cdot 2b^2 T \stackrel{(49)}{\leq} \frac{R^2}{5}; \quad (78)$$

That is, we derived the upper bounds for \ominus ; $\omin�$; \circledast ; $\omin�$; \circ . More specifically, the probability event E_{T-1} implies:

$$\begin{aligned} R_T^2 &\stackrel{(64)}{\leq} R^2 + \omin� + \omin� + \circledast + \omin� + \circ; \\ \omin� &\stackrel{(72)}{\leq} \frac{R^2}{5}; \quad \omin� &\stackrel{(77)}{\leq} \frac{R^2}{5}; \quad \circ &\stackrel{(78)}{\leq} \frac{R^2}{5}; \\ \sum_{t=0}^{T-1} e_t^2 &\stackrel{(71)}{\leq} \frac{R^4}{150 \ln^{4(K+1)}}; \quad \sum_{t=0}^{T-1} e_t^2 &\stackrel{(76)}{\leq} \frac{R^4}{150 \ln^{4(K+1)}}; \end{aligned}$$

In addition, we also have (see (70), (75) and our induction assumption)

$$P(E_{T-1} > 1) \leq \frac{(T-1)}{K+1}; \quad P(E_{\omin�} > 1) \leq \frac{1}{2(K+1)}; \quad P(E_{\circledast} > 1) \leq \frac{1}{2(K+1)};$$

where

$$\begin{aligned} E_{\omin�} &= \left(\sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln^{4(K+1)}} \text{ or } j \omin� j \leq \frac{R^2}{5} \right); \\ E_{\circledast} &= \left(\sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln^{4(K+1)}} \text{ or } j \circledast j \leq \frac{R^2}{5} \right); \end{aligned}$$

Therefore, probability event $E_{T-1} \setminus E_{\omin�} \setminus E_{\circledast}$ implies

$$R_T^2 \leq R^2 + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} = 2R^2;$$

which is equivalent to (55) and (56) for $t = T$, and

$$\mathbb{P}\{E_T g > 1\} \leq \mathbb{P}\{E_{T-1} \setminus E_{\otimes} \setminus E_{\otimes} g > 1\} \leq \mathbb{P}\{\bar{E}_{T-1} \setminus [\bar{E}_{-} \setminus \bar{E}_{\otimes}] g > 1\} \leq \mathbb{P}\{\bar{E}_{T-1} g > 1\} \leq \mathbb{P}\{\bar{E}_{-} g > 1\} \leq \mathbb{P}\{\bar{E}_{\otimes} g > 1\} \leq \frac{T}{K+1}.$$

We have now completed the inductive part of our proof. That is, for all $k = 0; 1; \dots; K+1$, we have $\mathbb{P}\{E_k g > 1\} \leq \frac{T}{K+1}$. Notably, when $k = K+1$, we can conclude that with a probability of at least $1 - \frac{T}{K+1}$:

$$f(\bar{x}^K) - f(x^*) \leq \frac{2R^2}{6(K+1)} \tag{58}$$

and $f(x^k) \leq g_{k=0}^K$ Q , which follows from (56).

Finally, if

$$= \min \left\{ \frac{1}{160L \ln^{4(K+1)}}, \frac{R}{208K \ln^{4(K+1)}}, \frac{R}{160b \ln^{4(K+1)}}, \frac{R}{1600b(K+1)} \right\};$$

then with probability at least $1 - \frac{T}{K+1}$

$$\begin{aligned} f(\bar{x}^K) - f(x^*) &\leq \frac{2R^2}{6(K+1)} \\ &\leq \max \left\{ \frac{320LR^2 \ln^{4(K+1)}}{K+1}, \frac{416R}{K+1}, \frac{320bR \ln^{4(K+1)}}{K+1}, 3200bR \right\}; \\ &= O \left(\max \left\{ \frac{LR^2 \ln^4 K}{K}, \frac{R}{K}, \frac{bR \ln^4 K}{K}, bR \right\} \right); \end{aligned}$$

This concludes the proof.

B.3 Quasi-Strongly Convex Case

Lemma B.5. Let Assumptions 2.2 and 2.4 with $\mu = 0$ hold on $Q = B_{2R}(x^*)$, where $R > \frac{1}{\mu}$, and let stepsize η_k satisfy $\eta_k \leq \frac{1}{L}$. If $x^k \in Q$ for all $k = 0; 1; \dots; K$, $K > 0$, then after K iterations of clipped-SGD we have

$$\begin{aligned} \|x^{K+1} - x^*\| &\leq \frac{1}{6} \exp(-\mu(K+1)) \|x^0 - x^*\| + \frac{1}{6} \sum_{k=0}^K \exp(-\mu(K-k)) \eta_k \|g^k - \nabla f(x^k)\| \\ &\quad + \frac{1}{6} \sum_{k=0}^K \exp(-\mu(K-k)) \eta_k^2 L^2 \|x^k - x^*\|^2; \end{aligned} \tag{79}$$

where η_k is defined in (48).

Proof. Using the update rule of clipped-SGD we obtain

$$\begin{aligned}
 \|x^{k+1} - x^k\|^2 &= \|x^k - x^k - 2\eta x^k + \eta \text{clip}(r f_k(x^k); \kappa) + 2\eta \text{clip}(r f_k(x^k); \kappa)\|^2 \\
 &= \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k) + 2\eta r f(x^k) - 2\eta r f(x^k) + 2\eta \kappa\|^2 \\
 &= \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k) + 2\eta r f(x^k) - 2\eta r f(x^k) + 2\eta \kappa\|^2 \\
 &\stackrel{(6);(4)}{\leq} (1 - 2\eta) \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k)\|^2 + 2L^2 \eta^2 \|r f(x^k) - r f(x^k)\|^2 + 2\eta^2 \kappa^2 \\
 &\stackrel{6.1=}{\leq} (1 - 2\eta) \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k) + 2\eta \kappa\|^2 + 2\eta^2 \kappa^2 \\
 &\stackrel{6.1=}{\leq} (1 - 2\eta) \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k) + 2\eta \kappa\|^2 + 2\eta^2 \kappa^2 \\
 &\leq \exp(-2\eta) \|x^k - x^k - 2\eta x^k + \eta r f(x^k) - 2\eta x^k + \eta r f(x^k) + 2\eta \kappa\|^2 + 2\eta^2 \kappa^2.
 \end{aligned}$$

Unrolling the recurrence, we obtain (79).

Theorem B.6. Let Assumptions 2.2 and 2.4 with $\eta > 0$ hold on $Q = B_{2R}(x^0)$, where $R > \kappa^0 / \eta$. Assume that $r f_k(x^k)$ satisfies Assumption 5.1 with parameters b_k, κ_k for $k = 0; 1; \dots; K$, $K > 0$ and

$$0 < \eta \leq \min \left\{ \frac{1}{40L \ln \frac{4(K+1)}{b_0}}, \frac{\ln(B_K)}{(K+1)}, \frac{\ln(C_K)}{(1+\kappa^2)}, \frac{2 \ln(D)}{(K+1)} \right\}; \tag{80}$$

$$B_K = \max \left\{ 2, \frac{(K+1)^2 R^2}{5400 \ln \frac{4(K+1)}{b_0} \ln^2(B_K)} \right\} = O \left(\max \left\{ 2, \frac{K^2 R^2}{2 \ln \frac{K}{b_0} \ln^2 \max \left\{ 2, \frac{K^2 R^2}{2 \ln \frac{K}{b_0}} \right\}} \right\} \right); \tag{81}$$

$$C_K = \max \left\{ 2, \frac{(\frac{K}{2} + 1) R}{480 \ln \frac{4(K+1)}{b_0} \ln(C_K)} \right\} = O \left(\max \left\{ 2, \frac{K R}{b \ln \frac{K}{b_0} \ln \max \left\{ 2, \frac{K R}{b \ln \frac{K}{b_0}} \right\}} \right\} \right); \tag{82}$$

$$D = \max \left\{ 2, \frac{R}{80 \ln(D)} \right\} = O \left(\max \left\{ 2, \frac{R}{b \ln \max \left\{ 2, \frac{R}{b} \right\}} \right\} \right); \tag{83}$$

$$\kappa_k = \frac{\exp(-(1+\kappa^2)k) R}{120 \ln \frac{4(K+1)}{b_0}}; \tag{84}$$

for some $\eta \in (0; 1]$ and $b = \max_{k=0; 1; \dots; K} b_k, \kappa = \max_{k=0; 1; \dots; K} \kappa_k$. Then, after K iterations the iterates produced by clipped-SGD with probability at least $1 - \delta$ satisfy

$$\|x^{K+1} - x^k\|^2 \leq 2 \exp(-\eta K) R^2; \tag{85}$$

In particular, when η equals the minimum from (80), then the iterates produced by clipped-SGD after K iterations with probability at least $1 - \delta$ satisfy

$$\|x^K - x^k\|^2 = O \left(\max \left\{ R^2 \exp \left(-\frac{K}{L \ln \frac{K}{b_0}} \right); \frac{2 \ln \frac{K}{b_0} \ln^2(B_K)}{K^2}; \frac{b R \ln \frac{K}{b_0} \ln(C_K)}{K}; \frac{b R \ln(D)}{K} \right\} \right); \tag{86}$$

Proof. Our proof follows similar steps to the one given by [Sadiev et al. \(2023\)](#). The main difference comes due to the presence of the bias in $r f_k(x^k)$. Therefore, for completeness, we provide the full proof here.

Let $R_k = \|x^k - x^k\|$ for all $k > 0$. As in the previous results, the main part of the proof is inductive. More precisely, for each $k = 0; 1; \dots; K + 1$ we consider probability event E_k as follows: inequalities

$$R_k^2 \leq 2 \exp(-\eta k) R^2 \tag{87}$$

hold for $t = 0; 1; \dots; k$ simultaneously. We aim to demonstrate through induction that $\text{Pf} E_k g > 1 - \epsilon^{k+1}$ for all $k = 0; 1; \dots; K + 1$. The base case $k = 0$, is trivial. Assuming that the statement holds for some $k = T - 1 \in K$, specially, $\text{Pf} E_{T-1} g > 1 - \epsilon^{(T-1)+1}$, we need to establish that $\text{Pf} E_T g > 1 - \epsilon^{T+1}$. Since $R_t^2 \leq 2 \exp(-t) R^2 \leq 2R^2$, we have $x^t \in B_{2R}(x)$, where function f is L -smooth. Thus, E_{T-1} implies

$$\|k r f(x^t)\| \leq L \|k x^t - x\| \stackrel{(87)}{\leq} \rho \frac{1}{2L} \exp(-t) R; \quad (88)$$

$$E_t[r f_t(x^t)] \stackrel{(80);(82);(84)}{\leq} E_t[r f_t(x^t)] - r f(x^t) + \|k r f(x^t)\| \stackrel{(7);(88)}{\leq} \rho \frac{1}{2L} \exp(-t) R + \frac{t}{2} \quad (89)$$

and

$$\|k_t\|^2 \leq 2k \rho f(x^t) k^2 + 2k r f(x^t) k^2 \stackrel{(88)}{\leq} \frac{5}{2} \rho^2 \frac{(84)}{6} \frac{\exp(-t) R^2}{4^2} \quad (90)$$

for all $t = 0; 1; \dots; T - 1$, where we use that $a + bk^2 \leq 2ak^2 + 2bk^2$ holding for all $a; b \in \mathbb{R}^d$.

Using Lemma B.5, we obtain that E_{T-1} implies

$$R_T^2 \leq \exp(-T) R^2 + 2 \sum_{t=0}^{T-1} \exp(-(T-1-t)) \|h_t - r f(x^t)\|_t + 2 \sum_{t=0}^{T-1} \exp(-(T-1-t)) \|k_t\|^2.$$

Next, we define random vectors

$$h_t = \begin{cases} x^t - x - r f(x^t); & \text{if } \|k x^t - x\| \leq \rho \frac{1}{2L} \exp(-t) R; \\ 0; & \text{otherwise;} \end{cases} \quad (91)$$

for $t = 0; 1; \dots; T - 1$. As per their definition, these random vectors are bounded with probability 1

$$\|k_t\| \leq \rho \frac{1}{2L} \exp(-t) R \quad (92)$$

for all $t = 0; 1; \dots; T - 1$. Moreover, for $t = 0; \dots; T - 1$ event E_{T-1} implies $\|k r f(x^t)\| \leq \rho \frac{1}{2L} \exp(-t) R$ (due to (88)) and

$$\|k x^t - x - r f(x^t)\| \leq \|k x^t - x\| + \|k r f(x^t)\| \stackrel{(88)}{\leq} \rho \frac{1}{2L} \exp(-t) R$$

for $t = 0; 1; \dots; T - 1$. The latter inequality means that E_{T-1} implies $h_t = x^t - x - r f(x^t)$ for all $t = 0; 1; \dots; T - 1$, meaning that from E_{T-1} it follows that

$$R_T^2 \leq \exp(-T) R^2 + 2 \sum_{t=0}^{T-1} \exp(-(T-1-t)) \|h_t\|_t + 2 \sum_{t=0}^{T-1} \exp(-(T-1-t)) \|k_t\|^2.$$

Next, we define the unbiased part and the bias of h_t as h_t^u and h_t^b , respectively:

$$h_t^u \stackrel{\text{def}}{=} \text{clip}(r f_t(x^t); h_t) - E_t[\text{clip}(r f_t(x^t); h_t)]; \quad h_t^b \stackrel{\text{def}}{=} E_t[\text{clip}(r f_t(x^t); h_t)] - r f(x^t); \quad (93)$$

for all $t = 0; \dots; T - 1$. We notice that $\tilde{u}_t = u_t + b_t$. Using new notation, we get that E_{T-1} implies

$$\begin{aligned}
 R_T^2 &\leq \exp(-T) R^2 \int_{t=0}^{T-1} \exp(-(T-1-t)) h_t; u_t \\
 &\quad \int_{t=0}^{T-1} \exp(-(T-1-t)) h_t; b_t + 2 \int_{t=0}^{T-1} \exp(-(T-1-t)) E k_t^u k^2 \\
 &\quad + 2 \int_{t=0}^{T-1} \exp(-(T-1-t)) k_t^u k^2 E k_t^u k^2 + 2 \int_{t=0}^{T-1} \exp(-(T-1-t)) k_t^b k^2: \quad (94)
 \end{aligned}$$

where we also apply inequality $ka + bk^2 \leq 2ka^2 + 2bk^2$ holding for all $a; b \in \mathbb{R}^d$ to upper bound $k_t^u k^2$. To conclude our inductive proof successfully, we must obtain sufficiently strong upper bounds with high probability for the terms $\tilde{-}; \tilde{-}; \tilde{0}; \tilde{-}; \tilde{0}$. In other words, we need to demonstrate that $\tilde{-} + \tilde{-} + \tilde{0} + \tilde{-} + \tilde{0} \leq \exp(-T) R^2$ with high probability. In the subsequent stages of the proof, we will rely on the bounds for the norms and second moments of u_t and b_t . First, as per the definition of the clipping operator, we can assert with probability 1 that

$$k_t^u k \leq 2 \tilde{t}: \quad (95)$$

Furthermore, given that E_{T-1} implies $E \|r f_t(x^t)\| \leq \tilde{c}$ for $t = 0; 1; \dots; T - 1$ (as per (89)), then, according to Lemma B.2, we can deduce that E_{T-1} implies

$$\begin{aligned}
 b_t &\leq E_t \text{clip}(r f_t(x^t); \tilde{t}) - E_t r f_t(x^t) + E_t r f_t(x^t) - r f_t(x^t) \\
 &\leq \frac{4}{\tilde{t}} + b; \quad (96)
 \end{aligned}$$

$$E_t k_t^u k^2 \leq 18 \tilde{c}^2; \quad (97)$$

for all $t = 0; 1; \dots; T - 1$.

Upper bound for $\tilde{-}$. By definition of u_t , we readily observe that $E_t[u_t] = 0$ and

$$E_t \left[\int_{t=0}^{T-1} \exp(-(T-1-t)) h_t; u_t \right] = 0:$$

Next, the sum $\tilde{-}$ contains only the terms that are bounded with probability 1:

$$\begin{aligned}
 \int_{t=0}^{T-1} \exp(-(T-1-t)) h_t; u_t &\leq \int_{t=0}^{T-1} \exp(-(T-1-t)) k_t k k_t^u k \\
 &\stackrel{(92);(95)}{\leq} \frac{P}{4} \frac{1}{2} (1 + L) \exp(-(T-1-t)) R_t \\
 &\stackrel{(80);(84)}{\leq} \frac{\exp(-T) R^2}{5 \ln \frac{4(K+1)}{}} \stackrel{\text{def}}{=} c: \quad (98)
 \end{aligned}$$

The conditional variances $\tilde{v}_t \stackrel{\text{def}}{=} E_t \left[\int_{t=0}^{T-1} \exp(-2(T-1-t)) h_t; u_t \right]^2$ of the summands are bounded:

$$\begin{aligned}
 \tilde{v}_t &\leq E_t \left[\int_{t=0}^{T-1} \exp(-2(T-1-t)) k_t k^2 k_t^u k^2 \right] \\
 &\stackrel{(92)}{\leq} 8 \frac{P^2}{4} (1 + L)^2 \exp(-2(T-1-t)) R^2 E_t k_t^u k^2 \\
 &\stackrel{(80)}{\leq} 10 \frac{P^2}{4} \exp(-2(T-1-t)) R^2 E_t k_t^u k^2: \quad (99)
 \end{aligned}$$

To summarize, we have demonstrated that $2 \left(\int_{t=0}^{T-1} \exp(-2(T-1-t)) h_t; u_t \right)^T$ is a bounded martingale difference sequence with bounded conditional variances \tilde{v}_t . Therefore, one can apply Bernstein's inequality (Lemma B.1)

with $X_t = 2(1 - \frac{1}{5})^{T-1} h_t; u_t^i$, parameter c as in (98), $b = \frac{1}{5} \exp(-T) R^2$, $G = \frac{\exp(-2T) R^4}{150 \ln \frac{4(K+1)}{1}}$ and get

$$P_{j \rightarrow j} > \frac{1}{5} \exp(-T) R^2 \text{ and } \sum_{t=0}^{K-1} \frac{\exp(-2T) R^4}{150 \ln \frac{4(K+1)}{1}} \leq 2 \exp\left(\frac{b^2}{2F + 2cb-3}\right) = \frac{1}{2(K+1)};$$

which is equivalent to

$$P_{E_{-g} > 1} \frac{1}{2(K+1)}; \text{ for } E_{-g} = \text{either } \sum_{t=0}^{K-1} \frac{\exp(-2T) R^4}{150 \ln \frac{4(K+1)}{1}} > \frac{1}{5} \exp(-T) R^2 \text{ or } j \rightarrow j \leq \frac{1}{5} \exp(-T) R^2 : (100)$$

Additionally, event E_{T-1} implies that

$$\begin{aligned} \sum_{t=0}^{K-1} \frac{\exp(-2T) R^4}{150 \ln \frac{4(K+1)}{1}} &\stackrel{(99)}{\leq} 10^2 \exp(-2T) R^2 \sum_{t=0}^{K-1} \frac{E_t k_t^u k^2}{\exp(-t)} \\ &\stackrel{(97); T \leq K+1}{\leq} 180^2 \exp(-2T) R^2 \sum_{t=0}^{K-1} \frac{1}{\exp(-t)} \\ &\stackrel{(84)}{\leq} 180^2 \exp(-2T) R^2 2(K+1) \exp(-K) \\ &\stackrel{(80); (81)}{\leq} \frac{\exp(-2T) R^4}{150 \ln \frac{4(K+1)}{1}} : \end{aligned} \quad (101)$$

Upper bound for ① . From E_{T-1} it follows that

$$\begin{aligned} \text{①} &\leq \sum_{t=0}^{K-1} 2 \exp(-T) \sum_{k=1}^K \frac{k_t k k^b k}{\exp(-t)} \\ &\stackrel{(92); (96)}{\leq} \frac{1}{2} (1+L) \exp(-T) R \sum_{t=0}^{K-1} \exp(-t) \leq \frac{4}{t} + b \\ &\stackrel{(80); (84)}{\leq} 3840^2 \exp(-T) 2(K+1) \exp(-T) \ln \frac{4(K+1)}{1} \\ &\quad + 2 \exp(-T) R (K+1) \exp(-T) b \\ &\stackrel{(80); (81); (83)}{\leq} \frac{1}{5} \exp(-T) R^2 : \end{aligned} \quad (102)$$

Upper bound for ② . From E_{T-1} it follows that

$$\begin{aligned} \text{②} &= 2^2 \exp(-T) \sum_{t=0}^{K-1} \frac{E_t k_t^u k^2}{\exp(-t)} \\ &\stackrel{(97)}{\leq} 144^2 \exp(-T) \sum_{t=0}^{K-1} \frac{1}{\exp(-t)} \\ &\stackrel{(84)}{\leq} 144^2 \exp(-T) 2(K+1) \exp(-K) \\ &\stackrel{(80)}{\leq} \frac{1}{5} \exp(-T) R^2 : \end{aligned} \quad (103)$$

Upper bound for ③ . By construction, we have

$$2^2 \exp(-T) \sum_{t=0}^{K-1} E_t k_t^u k^2 - E_t k_t^u k^2 = 0:$$

Next, the sum $\bar{\cdot}$ contains only the terms that are bounded with probability 1:

$$\begin{aligned}
 2^{-2} \exp(-\sum_{t=0}^{T-1} k_t^u k^2) E_t k_t^u k^2 &\stackrel{(95)}{=} \frac{16^{-2} \exp(-\sum_{t=0}^{T-1} k_t^u k^2)}{\exp(-\sum_{t=0}^{T-1} k_t^u k^2)} \\
 &\stackrel{(84)}{=} \frac{\exp(-\sum_{t=0}^{T-1} k_t^u k^2) R^2}{5 \ln^{4(K+1)}} \\
 &\stackrel{\text{def}}{=} c: \tag{104}
 \end{aligned}$$

The conditional variances

$$e_t^2 \stackrel{\text{def}}{=} E_t \left[4^{-4} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) k_t^u k^2 \right] E_t k_t^u k^2^{-2}$$

of the summands are bounded:

$$\begin{aligned}
 e_t^2 &\stackrel{(104)}{\leq} \frac{2^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2)}{5 \exp(-\sum_{t=0}^{T-1} k_t^u k^2) \ln^{4(K+1)}} E_t k_t^u k^2 E_t k_t^u k^2 \\
 &\leq \frac{4^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2)}{5 \exp(-\sum_{t=0}^{T-1} k_t^u k^2) \ln^{4(K+1)}} E_t k_t^u k^2 : \tag{105}
 \end{aligned}$$

To summarize, we have demonstrated that $2^{-2} (1 - \sum_{t=0}^{T-1} k_t^u k^2) E_t k_t^u k^2$ is a bounded martingale difference sequence with bounded conditional variances $e_t^2 \leq \frac{4^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2)}{5 \exp(-\sum_{t=0}^{T-1} k_t^u k^2) \ln^{4(K+1)}} E_t k_t^u k^2$. Therefore, one can apply Bernstein's inequality (Lemma B.1) with $X_t = 2^{-2} (1 - \sum_{t=0}^{T-1} k_t^u k^2) E_t k_t^u k^2$, parameter c as in (104), $b = \frac{1}{5} \exp(-\sum_{t=0}^{T-1} k_t^u k^2) R^2$, $G = \frac{\exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^4}{150 \ln^{4(K+1)}}$ and get

$$P \left[\sum_{j=1}^n |j| > \frac{1}{5} \exp(-\sum_{t=0}^{T-1} k_t^u k^2) R^2 \right] \leq \sum_{l=0}^{K-1} e_t^2 \leq \frac{\exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^4}{150 \ln^{4(K+1)}} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb^{-3}}\right) = \frac{1}{2(K+1)};$$

which is equivalent to

$$P \left[\sum_{j=1}^n |j| > \frac{1}{5} \exp(-\sum_{t=0}^{T-1} k_t^u k^2) R^2 \right] \leq \frac{1}{2(K+1)}; \text{ for } E^- = \left(\sum_{t=0}^{K-1} e_t^2 > \frac{\exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^4}{150 \ln^{4(K+1)}} \text{ or } \sum_{j=1}^n |j| > \frac{1}{5} \exp(-\sum_{t=0}^{T-1} k_t^u k^2) R^2 \right) : \tag{106}$$

Additionally, event E_{T-1} implies that

$$\begin{aligned}
 \sum_{l=0}^{K-1} e_t^2 &\stackrel{(105)}{\leq} \frac{4^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^2}{5 \ln^{4(K+1)}} \sum_{t=0}^{K-1} E_t k_t^u k^2 \\
 &\stackrel{(97); T \leq 6K+1}{\leq} \frac{72^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^2}{5 \ln^{4(K+1)}} \sum_{t=0}^{K-1} \frac{1}{\exp(-\sum_{t=0}^{T-1} k_t^u k^2)} \\
 &\stackrel{(84)}{\leq} \frac{72^{-2} \exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^2}{5 \ln^{4(K+1)}} (K+1) \exp(-\sum_{t=0}^{K-1} k_t^u k^2) \\
 &\stackrel{(80)}{\leq} \frac{\exp(-2 \sum_{t=0}^{T-1} k_t^u k^2) R^4}{150 \ln^{4(K+1)}} : \tag{107}
 \end{aligned}$$

Upper bound for \circ . From E_{T-1} it follows that

$$\begin{aligned}
 \circ &= \sum_{l=0}^{K-1} 2^{-2l} \exp(-l(T-1)) k_l^2 \\
 &\stackrel{(96)}{\leq} \sum_{t=0}^{K-1} 2^{-2t} \exp(-t(T-1)) \exp(t) \left(\frac{64}{t^2} + 2b^2 \right) \\
 &\stackrel{(84); T \leq K+1}{\leq} \frac{921600^4 \exp(-t(T-1)) \ln^{4(K+1)} \sum_{t=0}^{K-1} \exp(2t) \left(1 + \frac{t}{2}\right) \exp(t)}{R^2} \\
 &\quad + 2 \sum_{t=0}^{K-1} \exp(-t(T-1)) b^2 \sum_{t=0}^{K-1} \exp(t) \\
 &\stackrel{(80); (81); (83)}{\leq} \frac{921600^4 \exp(-t(T-3)) \ln^{4(K+1)} (K+1) \exp(2K)}{R^2} \\
 &\quad + 2 \sum_{t=0}^{K-1} \exp(-t(T-1)) b^2 \exp(K) (K+1) \\
 &\stackrel{(80); (81); (83)}{\leq} \frac{1}{5} \exp(-T) R^2. \tag{108}
 \end{aligned}$$

That is, we derived the upper bounds for \neg ; $-$; \circledast ; $\bar{-}$; \circ . More specifically, the probability event E_{T-1} implies:

$$\begin{aligned}
 R_T^2 &\stackrel{(94)}{\leq} \frac{1}{6} \exp(-T) R^2 + \neg + - + \circledast + \bar{-} + \circ; \\
 - &\stackrel{(102)}{\leq} \frac{1}{6} \exp(-T) R^2; \quad \circledast \stackrel{(103)}{\leq} \frac{1}{6} \exp(-T) R^2; \\
 \circ &\stackrel{(108)}{\leq} \frac{1}{6} \exp(-T) R^2 \\
 \sum_{t=0}^{K-1} 2^{-2t} &\stackrel{(101)}{\leq} \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}}; \quad \sum_{t=0}^{K-1} e^{2t} \stackrel{(107)}{\leq} \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}}.
 \end{aligned}$$

Moreover, we also have (see (100), (106) and our induction assumption)

$$\begin{aligned}
 \text{Pf } E_{T-1} g > 1 &\leq \frac{(T-1)}{K+1}; \\
 \text{Pf } E_{-} g > 1 &\leq \frac{1}{2(K+1)}; \quad \text{Pf } E_{\bar{-}} g > 1 \leq \frac{1}{2(K+1)};
 \end{aligned}$$

where

$$\begin{aligned}
 E_{-} &= \left(\text{either } \sum_{t=0}^{K-1} 2^{-2t} > \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}} \text{ or } \sum_{t=0}^{K-1} e^{2t} > \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}} \right) \\
 E_{\bar{-}} &= \left(\text{either } \sum_{t=0}^{K-1} 2^{-2t} > \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}} \text{ or } \sum_{t=0}^{K-1} e^{2t} > \frac{\exp(-2T) R^4}{150 \ln^{4(K+1)}} \right)
 \end{aligned}$$

Therefore, probability event $E_{T-1} \setminus E_{-} \setminus E_{\bar{-}}$ implies

$$R_T^2 \stackrel{(94)}{\leq} \frac{1}{6} \exp(-T) R^2 + \neg + - + \circledast + \bar{-} + \circ$$

which is equivalent to (87) for $t = T$, and

$$\text{Pf } E_T g > \text{Pf } E_{T-1} \setminus E_{-} \setminus E_{\bar{-}} g = 1 - \text{Pf } \bar{E}_{T-1} [\bar{E}_{-} \cap \bar{E}_{\bar{-}} g > 1] \leq \frac{T}{K+1};$$

We have now completed the inductive part of our proof. That is, for all $k = 0; 1; \dots; K + 1$, we have $\Pr E_k > 1 - \epsilon^{k+1}$. Notably, when $k = K + 1$, we can conclude that with a probability of at least $1 - \epsilon$:

$$kx^{K+1} \leq k^2 \cdot 6 \cdot 2 \exp(-((K+1))R^2):$$

Finally, if

$$\begin{aligned} &= \min \left(\frac{1}{400L \ln^{4(K+1)}}; \frac{\ln(B_K)}{(K+1)}; \frac{\ln(C_K)}{(1+K^2)}; \frac{2 \ln(D)}{(K+1)} \right); \\ B_K &= \max \left(2; \frac{(K+1)^2 R^2}{5400 \cdot 2 \ln^{4(K+1)} \ln^2(B_K)} \right); \\ C_K &= \max \left(2; \frac{(\frac{K}{2} + 1) R}{480 \ln^{4(K+1)} \ln(C_K)} \right); \\ D &= \max \left(2; \frac{R}{80 \ln(D)} \right) \end{aligned}$$

then with probability at least $1 - \epsilon$

$$\begin{aligned} &kx^{K+1} \leq k^2 \cdot 6 \cdot 2 \exp(-((K+1))R^2) \\ &\leq 2R^2 \max \left(\frac{(K+1)}{400L \ln^{4(K+1)}}; \frac{1}{B_K}; \frac{1}{C_K}; \frac{1}{D} \right) \\ &= O \left(\max \left(R^2 \exp \left(\frac{K}{L \ln K} \right); \frac{2 \ln K}{K^2}; \frac{bR \ln K}{K}; \frac{bR \ln(D)}{K} \right) \right) \end{aligned}$$

This concludes the proof.

C Proofs for clipped-SSTM

C.1 Convex Case

The analysis of clipped-SSTM in the convex case relies on the following lemma from (Sadiev et al., 2023).

Lemma C.1 (Lemma F.1 from (Sadiev et al., 2023)) Let Assumptions 2.2 and 2.3 with $\epsilon = 0$ hold on $Q = B_{3R}(x)$, where $R > kx^0$, and let stepsize parameter a satisfy $a > 1$. If $x^k, y^k, z^k \in B_{3R}(x)$ for all $k = 0, 1, \dots, N$, $N > 0$, then after N iterations of clipped-SSTM for all $z \in B_{3R}(x)$ we have

$$\begin{aligned} A_N f(y^N) - f(z) &\leq \frac{1}{2}kz^0 - zk^2 - \frac{1}{2}kz^N - zk^2 + \sum_{k=0}^{N-1} \frac{1}{k+1} \frac{1}{k+1}; z = z^k + \frac{1}{k+1} r f(x^{k+1}) \\ &\quad + \sum_{k=0}^{N-1} \frac{1}{k+1} \frac{1}{k+1} k^2; \end{aligned} \quad (109)$$

$$\frac{1}{k+1} \stackrel{\text{def}}{=} \text{clip}(r f(x^{k+1}); k) = r f(x^{k+1}); \quad (110)$$

Next, we also use the following technical result from (Gorbunov et al., 2020).

Lemma C.2 (Lemma E.1 from (Gorbunov et al., 2020)). Let sequences $\{g_k\}_{k>0}$ and $\{A_k\}_{k>0}$ satisfy

$$A_0 = 0; \quad A_{k+1} = A_k + g_{k+1}; \quad g_{k+1} = \frac{k+2}{2aL} g_k > 0; \quad (111)$$

where $a > 0, L > 0$. Then for all $k > 0$

$$A_{k+1} = \frac{(k+1)(k+4)}{4aL}; \quad (112)$$

$$A_{k+1} > aL \frac{1}{k+1}; \quad (113)$$

Theorem C.3. Let Assumptions 2.2 and 2.3 with $\epsilon = 0$ hold on $Q = B_{3R}(x)$, where $R > kx^0$. Assume that $r f(x^{k+1})$ satisfies Assumption 5.1 with parameters b_k, g_k for $k = 0, 1, \dots, K$, $K > 0$ and

$$a > \max \left\{ \frac{8}{97200} \ln^2 \frac{4K}{LR}; \frac{1800(K+1)^p K^q \ln \frac{4K}{LR}}{LR}; \frac{4b(K+2)^2}{15LR}; \frac{60b(K+2) \ln \frac{4K}{LR}}{LR}; \right\}; \quad (114)$$

$$k = \frac{R}{30 \frac{1}{k+1} \ln \frac{4K}{LR}}; \quad (115)$$

for some $\beta \in (0, 1]$ and $b = \max_{k=0,1,\dots,K} b_k, g = \max_{k=0,1,\dots,K} g_k$. Then, after K iterations of clipped-SSTM the iterates with probability at least $1 - \beta$ satisfy

$$f(y^K) - f(x) \leq \frac{6aLR^2}{K(K+3)} \quad \text{and} \quad f(x^k) g_{k=0}^K; f(z^k) g_{k=0}^K; f(y^k) g_{k=0}^K \in B_{2R}(x); \quad (116)$$

In particular, when parameter a equals the maximum from (114), then the iterates produced by clipped-SSTM after K iterations with probability at least $1 - \beta$ satisfy

$$f(y^K) - f(x) = O \left(\max \left\{ \frac{8}{K^2}; \frac{R \ln \frac{4K}{LR}}{K^2}; \frac{bR \ln \frac{4K}{LR}}{K}; bR; \frac{1}{K}; \right\} \right); \quad (117)$$

Proof. Our proof follows similar steps to the one given by Sadiev et al. (2023). The main difference comes due to the presence of the bias $r f(x^k)$. Therefore, for completeness, we provide the full proof here.

Let $R_k = kz^k - x^k$, $R_0 = R_0$, and $R_{k+1} = \max\{R_k, R_{k+1} g\}$ for all $k > 0$. We will initially demonstrate through induction that for all $k > 0$, the iterates x^{k+1}, z^k, y^k belong to $B_{R_k}(x)$. The base of the induction is straightforward because $y^0 = z^0, R_0 = R_0$, and $x^1 = \frac{A_0 y^0 + z^0}{A_1} = z^0$. Now, assume that for some $k > 1$,

$x^l; z^{l-1}; y^{l-1} \in B_{R_{l-1}}(x)$. By the definitions of R_l and R_{l-1} , we have $z^l \in B_{R_l}(x) \cap B_{R_{l-1}}(x)$. As y^l is a convex combination of $y^{l-1} \in B_{R_{l-1}}(x) \cap B_{R_l}(x)$, it follows that $z^l \in B_{R_l}(x)$ and, given the convex nature of $B_{R_l}(x)$, we can conclude that $y^l \in B_{R_l}(x)$. Finally, as x^{l+1} is a convex combination of y^l and z^l , it is evident that x^{l+1} also lies in $B_{R_l}(x)$.

Our next objective is to establish, by induction, that $R_t \leq 3R$ with high probability. This will enable us to apply the result from Lemma B.3 and subsequently utilize Bernstein's inequality to estimate the stochastic component of the upper bound. To be more precise, for each $k = 0; \dots; K+1$, we consider the probability event E_k , defined as follows: inequalities

$$\sum_{l=0}^{X-1} \sum_{i=1}^{l+1} \sum_{j=1}^{i+1} x^j \cdot z^l + \sum_{l=1}^{i+1} r f_l(x^{l+1}) + \sum_{l=0}^{X-1} \sum_{i=1}^{l+1} k_{i+1} k^2 \leq 6R^2; \tag{118}$$

$$R_t \leq 2R \tag{119}$$

hold for all $t = 0; 1; \dots; k$ simultaneously. We aim to demonstrate through induction that $\Pr[E_k] > 1 - \epsilon^{k/(K+1)}$ for all $k = 0; 1; \dots; K+1$. The base case $k = 0$, is trivial: the left-hand side of (118) equals zero and $R > R_0$ by definition. Assuming that the statement holds for some $k = T-1 \leq K$, specifically, $\Pr[E_{T-1}] > 1 - \epsilon^{(T-1)/(K+1)}$, we need to establish that $\Pr[E_T] > 1 - \epsilon^{T/(K+1)}$.

To begin, we observe that the probability event E_{T-1} implies that $R_t \leq 2R$ for all $t = 0; 1; \dots; T-1$. Moreover, it implies that

$$\|kz^T - x\|_k \leq \|kz^T - x\|_k + \sum_{t=0}^{T-1} k \Pr[E_{T-1}(x^T)] \leq 2R + \sum_{t=0}^{T-1} \epsilon^{(T-1)/(K+1)} \tag{115}$$

Hence, with E_{T-1} implying $\sum_{k=0}^T g_k^T \leq Q$, we confirm that the conditions of Lemma C.1 are met, resulting in

$$A_t f(y^t) - f(x) \leq \frac{1}{2} R_0^2 + \frac{1}{2} R_t^2 + \sum_{l=0}^{X-1} \sum_{i=1}^{l+1} \sum_{j=1}^{i+1} x^j \cdot z^l + \sum_{l=1}^{i+1} r f_l(x^{l+1}) + \sum_{l=0}^{X-1} \sum_{i=1}^{l+1} k_{i+1} k^2 \tag{120}$$

for all $t = 0; 1; \dots; T$ simultaneously and for all $t = 1; \dots; T-1$ this probability event also implies that

$$f(y^t) - f(x) \stackrel{(118);(120)}{\leq} \frac{1}{6} R_0^2 + \frac{\frac{1}{2} R_t^2 + R^2}{A_t} \leq \frac{3R^2}{2A_t} = \frac{6aLR^2}{t(t+3)} \tag{121}$$

Considering that $f(y^T) - f(x) > 0$, we can further deduce from (120) that when E_{T-1} holds, the following holds as well:

$$R_T^2 \leq R_0^2 + 2 \sum_{t=0}^{X-1} \sum_{i=1}^{t+1} \sum_{j=1}^{i+1} x^j \cdot z^t + \sum_{t=1}^{i+1} r f_t(x^{t+1}) + 2 \sum_{t=0}^{X-1} \sum_{i=1}^{t+1} k_{i+1} k^2$$

$$\leq R^2 + 2B_T \tag{122}$$

Prior to our estimation of B_T , we need to establish several helpful inequalities. We start with showing that E_{T-1} implies $\|k f(x^{l+1})\|_k \leq L k x^0$ for all $t = 0; 1; \dots; T-1$. For $t = 0$ we have $x^1 = x^0$ and

$$\|k f(x^1)\|_k = \|k f(x^0)\|_k \stackrel{(3)}{\leq} L k x^0 = k \leq \frac{R}{a-1} = \frac{0}{2} + \frac{60 \ln \frac{4K}{a}}{a} \stackrel{(114)}{\leq} \frac{0}{4} \tag{123}$$

Next, for $t = 1; \dots; T - 1$ we have $\|x^{t+1} - z^t\| = A_t \|y^t - x^{t+1}\|$ and event E_{T-1} implies

$$\begin{aligned}
 \|r f(x^{t+1}) - r f(y^t)\| &\stackrel{(3);(4)}{\leq} L \|x^{t+1} - y^t\| + \frac{\rho}{2L} \frac{\|f(y^t) - f(x^t)\|^2}{S} \\
 &\stackrel{(121)}{\leq} \frac{L_{t+1}}{A_t} \|x^{t+1} - z^t\| + \frac{12aL^2R^2}{t(t+3)} \\
 &\leq \frac{4LR_{t+1}}{A_t} + \frac{12aL^2R^2}{t(t+3)} \\
 &= \frac{R}{60} \frac{240L_{t+1}^2 \ln^{4K}}{t(t+3)} + 60 \frac{12aL^2 \frac{t+2}{2aL} \ln^{4K}}{t(t+3)} \\
 &\stackrel{(112);(115)}{\leq} \frac{R}{2} \frac{240L_{t+1}^2 \ln^{4K}}{t(t+3)} + 60 \frac{12aL^2 \frac{t+2}{2aL} \ln^{4K}}{t(t+3)} \\
 &= \frac{R}{2} \frac{240(t+2)^2 \ln^{4K}}{t(t+3)a} + 60 \frac{3(t+2)^2 \ln^{4K}}{t(t+3)a} \\
 &\leq \frac{R}{2} \frac{540 \ln^{4K}}{a} + \frac{90 \rho \ln^{4K}}{a} \stackrel{(114)}{\leq} \frac{R}{4}; \tag{124}
 \end{aligned}$$

where in the last row we use $\frac{(t+2)^2}{t(t+3)} \leq \frac{9}{4}$ for all $t > 1$. Therefore, probability event E_{T-1} implies that

$$\|E_t[r f_t(x^{t+1})] - E_t[r f_t(x^{t+1})] - r f(x^{t+1})\| \leq \frac{R}{4} \leq \frac{t}{2} \tag{125}$$

and

$$\|x^t - z^t + \sum_{i=0}^{t-1} r f(x^{i+1})\| \leq \|x^t - z^t\| + \sum_{i=0}^{t-1} \|r f(x^{i+1})\| \stackrel{(119);(123);(124)}{\leq} 2R + \frac{R}{60 \ln^{4K}} \leq 3R \tag{126}$$

for all $t = 0; 1; \dots; T - 1$. Next, we define random vectors

$$\xi_t = \begin{cases} x^t - z^t + \sum_{i=0}^{t-1} r f(x^{i+1}); & \text{if } \|x^t - z^t + \sum_{i=0}^{t-1} r f(x^{i+1})\| \leq 3R; \\ 0; & \text{otherwise;} \end{cases}$$

for all $t = 0; 1; \dots; T - 1$. As per their definition, these random vectors are bounded with probability 1

$$\|\xi_t\| \leq 3R; \tag{127}$$

This means that E_{T-1} implies $\xi_t = x^t - z^t + \sum_{i=0}^{t-1} r f(x^{i+1})$ for all $t = 0; 1; \dots; T - 1$. Then, from E_{T-1} it follows that

$$B_T = \sum_{t=0}^{T-1} \xi_{t+1} h_{t+1}; \quad \sum_{t=0}^{T-1} \|\xi_{t+1}\|^2 \leq \sum_{t=0}^{T-1} 9R^2;$$

Next, we define the unbiased part and the bias of ξ_t as ξ_t^u and ξ_t^b , respectively:

$$\xi_t^u = \text{clip}(r f_t(x^{t+1}); t) - E_t[\text{clip}(r f_t(x^{t+1}); t)]; \quad \xi_t^b = E_t[\text{clip}(r f_t(x^{t+1}); t)] - r f(x^{t+1}); \tag{128}$$

We notice that $\tilde{z}_t = \tilde{u}_t + \tilde{b}_t$. Using new notation, we get that E_{T-1} implies

$$\begin{aligned}
 B_T = & \sum_{t=0}^{T-1} \mathbb{E}_{t+1} \left[\sum_{i=0}^h \left(\tilde{z}_{t+1}^{(i)} \right)^2 \right] \\
 & + 2 \sum_{t=0}^{T-1} \mathbb{E}_t \left[\sum_{i=0}^h \tilde{z}_{t+1}^{(i)} \right] \\
 & + 2 \sum_{t=0}^{T-1} \mathbb{E}_t \left[\sum_{i=0}^h \left(\tilde{z}_{t+1}^{(i)} \right)^2 \right] \\
 & + 2 \sum_{t=0}^{T-1} \mathbb{E}_t \left[\sum_{i=0}^h \left(\tilde{z}_{t+1}^{(i)} \right)^2 \right]
 \end{aligned} \tag{129}$$

To conclude our inductive proof successfully, we must obtain sufficiently strong upper bounds with high probability for the terms $\tilde{z}_t^{(i)}$; $\tilde{z}_t^{(i)}$; $\tilde{z}_t^{(i)}$; $\tilde{z}_t^{(i)}$. In other words, we need to demonstrate that $\tilde{z}_t^{(i)} \leq R^2$ with a high probability. In the subsequent stages of the proof, we will rely on the bounds for the norms and second moments of \tilde{u}_t and \tilde{b}_t . First, as per the definition of the clipping operator, we can assert with probability 1 that

$$\|\tilde{u}_{t+1}\| \leq 2 \|\tilde{u}_t\| \tag{130}$$

Moreover, since E_{T-1} implies that $\mathbb{E}_t \|\tilde{r}_t(x^{t+1})\| \leq 2$ for $t = 0, 1, \dots, T-1$ (see (125)), then, according to Lemma B.2, we can deduce that E_{T-1} implies

$$\|\tilde{b}_{t+1}\| \leq \mathbb{E}_t \text{clip}(\tilde{r}_t(x^{t+1}); \tilde{u}_t) \leq \mathbb{E}_t \|\tilde{r}_t(x^{t+1})\| + \mathbb{E}_t \|\tilde{r}_t(x^{t+1}) - \tilde{r}_t(x^{t+1})\| \leq \frac{4}{t} + b; \tag{131}$$

$$\mathbb{E}_t \|\tilde{u}_{t+1}\|^2 \leq 18 \|\tilde{u}_t\|^2; \tag{132}$$

Upper bound for \tilde{z}_t . By definition of \tilde{u}_t , we readily observe that $\mathbb{E}_t[\tilde{u}_t] = 0$ and

$$\mathbb{E}_t \|\tilde{u}_{t+1}\| = 0;$$

Next, the sum \tilde{z}_t contains only the terms that are bounded with probability 1:

$$\|\tilde{z}_{t+1}^{(i)}\| \leq \|\tilde{u}_{t+1}^{(i)}\| \leq \|\tilde{u}_{t+1}\| \leq \|\tilde{u}_t\| \leq \frac{R^2}{6} \leq \frac{R^2}{6} \leq \frac{R^2}{6} \stackrel{(115)}{=} \frac{R^2}{5 \ln 4K} \stackrel{\text{def}}{=} c; \tag{133}$$

The conditional variances $\sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{z}_{t+1}^{(i)} \right)^2 \right]$ of the summands are bounded:

$$\sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{z}_{t+1}^{(i)} \right)^2 \right] \leq \sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{u}_{t+1}^{(i)} \right)^2 \right] \leq \sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{u}_{t+1} \right)^2 \right] \leq \frac{9}{6} \sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{u}_{t+1} \right)^2 \right] \leq \frac{9}{6} \sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{u}_{t+1} \right)^2 \right]; \tag{134}$$

To summarize, we have demonstrated that $\sum_{i=0}^h \tilde{z}_{t+1}^{(i)} \mathbf{g}_{t=0}^T$ is a bounded martingale difference sequence with bounded conditional variances $\sum_{i=0}^h \mathbb{E}_t \left[\left(\tilde{z}_{t+1}^{(i)} \right)^2 \right]$. Therefore, one can apply Bernstein's inequality (Lemma B.1) with $X_t = \sum_{i=0}^h \tilde{z}_{t+1}^{(i)} \mathbf{g}_{t=0}^T$, parameter c as in (133), $b = \frac{R^2}{5}$, $G = \frac{R^4}{150 \ln 4K}$ and get

$$\left(\sum_{j=1}^T \|\tilde{z}_j\| > \frac{R^2}{5} \text{ and } \sum_{t=0}^{T-1} \mathbb{E}_t \left[\sum_{i=0}^h \left(\tilde{z}_{t+1}^{(i)} \right)^2 \right] \leq 2 \exp \left(- \frac{b^2}{2G + 2cb_3} \right) = \frac{1}{2K}; \right.$$

which is equivalent to

$$\text{Pf } E_{\tilde{z}} \left[\sum_{j=1}^T \|\tilde{z}_j\| > \frac{R^2}{2K}; \text{ for } E_{\tilde{z}} = \left(\text{either } \sum_{t=0}^{T-1} \mathbb{E}_t \left[\sum_{i=0}^h \left(\tilde{z}_{t+1}^{(i)} \right)^2 \right] > \frac{R^4}{150 \ln 4K} \text{ or } \sum_{j=1}^T \|\tilde{z}_j\| \geq \frac{R^2}{5} \right); \tag{135}$$

In addition, E_{T-1} implies that

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbb{E} \left[\frac{1}{6} \left(9R^2 \sum_{t=0}^{T-1} \mathbb{E} \left[k_{t+1}^u k^2 \right] - 162 \sum_{t=0}^{T-1} R^2 \right) \right] \\ = \frac{162 \sum_{t=0}^{T-1} R^2}{4a^2 L^2} \mathbb{E} \left[(t+2)^2 \right] \\ = \frac{1}{a^2} \frac{81 \sum_{t=0}^{T-1} R^2 T(T+1)^2}{2a^2 L^2} \stackrel{(114)}{=} \frac{R^4}{150 \ln \frac{4K}{\epsilon}}. \end{aligned} \quad (136)$$

Upper bound for ② . From E_{T-1} it follows that

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbb{E} \left[\frac{1}{6} \sum_{t=1}^T k_{t+1}^b k_{t+1} k_{t+1} \right] &\stackrel{(127);(131)}{\leq} \sum_{t=0}^{T-1} \mathbb{E} \left[\frac{1}{6} \sum_{t=1}^T 3R k_{t+1} \right] \leq \frac{4}{t} + b \\ &\stackrel{(115)}{\leq} \frac{360 \sum_{t=0}^{T-1} \ln \frac{4K}{\epsilon} \mathbb{E} \left[\frac{1}{6} \sum_{t=1}^T k_{t+1}^2 \right] + 3bR}{4a^2 L^2} \mathbb{E} \left[(t+2)^2 \right] + \frac{3bR}{2aL} \mathbb{E} \left[(t+2) \right] \\ &\leq \frac{360 \sum_{t=0}^{T-1} \ln \frac{4K}{\epsilon} T(T+1)^2}{4a^2 L^2} + \frac{3bRT(T+1)}{2aL} \stackrel{(114)}{\leq} \frac{R^2}{5}. \end{aligned} \quad (137)$$

Upper bound for ③ . First, we have

$$\mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] - \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] = 0:$$

Next, the sum ③ contains only the terms that are bounded with probability 1:

$$\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 + \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] \\ \stackrel{(130)}{\leq} \frac{16 \sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2}{5 \ln \frac{4K}{\epsilon}} \stackrel{\text{def}}{=} c. \quad (138)$$

The conditional variances $e_t^2 \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^T k_{t+1}^u k^2 \right] - \mathbb{E} \left[\sum_{t=1}^T k_{t+1}^u k^2 \right]^2$ of the summands are bounded:

$$e_t^2 \stackrel{(138)}{\leq} \frac{R^2}{5 \ln \frac{4K}{\epsilon}} \mathbb{E} \left[\sum_{t=1}^T k_{t+1}^u k^2 \right] - \mathbb{E} \left[\sum_{t=1}^T k_{t+1}^u k^2 \right]^2 \leq \frac{1}{6} \sum_{t=1}^T R^2 \mathbb{E} \left[k_{t+1}^u k^2 \right]; \quad (139)$$

To summarize, we have demonstrated that $\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 + \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right]$ is a bounded martingale difference sequence with bounded conditional variances $e_t^2 \leq \frac{R^2}{5}$. Therefore, one can apply Bernstein's inequality (Lemma B.1) with $X_t = \sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 + \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right]$, parameter c as in (138), $b = \frac{R^2}{5}$, $G = \frac{R^4}{150 \ln \frac{4K}{\epsilon}}$ and get

$$\mathbb{P} \left(\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 + \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] > \frac{R^2}{5} \right) \leq \frac{R^4}{150 \ln \frac{4K}{\epsilon}} \exp \left(- \frac{b^2}{2G + 2cb} \right) = \frac{R^4}{150 \ln \frac{4K}{\epsilon}};$$

which is equivalent to

$$\mathbb{P} \left(\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 + \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] > \frac{R^2}{5} \right) \leq \frac{R^4}{150 \ln \frac{4K}{\epsilon}}; \quad \text{for } \mathbb{E} \left[\sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] > \frac{R^2}{5} \text{ or } \sum_{t=0}^{T-1} \frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 > \frac{R^2}{5}; \quad (140)$$

In addition, E_{T-1} implies that

$$\sum_{t=0}^{T-1} \mathbb{E} \left[\frac{1}{6} \sum_{t=1}^T k_{t+1}^u k^2 \right] \stackrel{(139)}{\leq} \frac{R^2}{5} + \frac{1}{6} \sum_{t=0}^{T-1} \mathbb{E} \left[\sum_{t=1}^T k_{t+1}^u k^2 \right] \stackrel{(136)}{\leq} \frac{R^4}{150 \ln \frac{4K}{\epsilon}}; \quad (141)$$

Upper bound for $\bar{\tau}$. From E_{T-1} it follows that

$$\bar{\tau} = 2 \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 E_{t+1}^h \bar{u}_{t+1}^2 \leq \frac{1}{R^2} 9R^2 \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 E_{t+1}^h \bar{u}_{t+1}^2 \stackrel{(136)}{\leq} \frac{R^2}{150 \ln 4K} \leq \frac{R^2}{5} \quad (142)$$

Upper bound for $\bar{\sigma}$. From E_{T-1} it follows that

$$\begin{aligned} \bar{\sigma} &= 2 \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 \bar{b}_{t+1}^2 \leq 4 \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 \frac{16}{t} + b^2 \\ &\stackrel{(115)}{=} \frac{57600 \ln^2 4K}{R^2} \sum_{t=0}^{T-1} \bar{X}_{t+1}^4 + 4b^2 \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 = \frac{57600 \ln^2 4K}{16a^4 L^4 R^2} \sum_{t=0}^{T-1} (t+2)^4 + \frac{4b^2}{4a^2 L^2} \sum_{t=0}^{T-1} (t+2)^2 \\ &\leq \frac{57600 \ln^2 4K T(T+1)^4}{16a^4 L^4 R^2} + \frac{4b^2 T(T+1)^2}{4a^2 L^2} \stackrel{(114)}{\leq} \frac{R^2}{5} \end{aligned} \quad (143)$$

That is, we derived the upper bounds for $\bar{\tau}$; $\bar{\sigma}$; $\bar{\tau} + \bar{\sigma}$; $\bar{\tau} + \bar{\sigma} + \bar{\tau} + \bar{\sigma}$. More specifically, the probability event E_{T-1} implies:

$$\begin{aligned} B_T &\stackrel{(129)}{\leq} \frac{R^2}{6} + \bar{\tau} + \bar{\sigma} + \bar{\tau} + \bar{\sigma} + \bar{\tau} + \bar{\sigma}; \\ \bar{\tau} &\stackrel{(137)}{\leq} \frac{R^2}{5}; \quad \bar{\sigma} \stackrel{(142)}{\leq} \frac{R^2}{5}; \quad \bar{\tau} + \bar{\sigma} \stackrel{(143)}{\leq} \frac{R^2}{5}; \\ \bar{X}_{t+1}^2 &\stackrel{(136)}{\leq} \frac{R^4}{150 \ln 4K}; \quad e_t^2 \stackrel{(141)}{\leq} \frac{R^4}{150 \ln 4K}; \end{aligned}$$

In addition, we also have (see (135), (140) and our induction assumption)

$$\text{Pf } E_{T-1} g > 1 \leq \frac{(T-1)}{K}; \quad \text{Pf } E_{-g} > 1 \leq \frac{1}{2K}; \quad \text{Pf } E_{\otimes g} > 1 \leq \frac{1}{2K};$$

where

$$\begin{aligned} E_{-} &= \left(\text{either } \sum_{t=0}^{T-1} \bar{X}_{t+1}^2 > \frac{R^4}{150 \ln 4K} \text{ or } \sum_{j=1}^T \bar{X}_j > \frac{R^2}{5} \right); \\ E_{\otimes} &= \left(\text{either } \sum_{t=0}^{T-1} e_t^2 > \frac{R^4}{150 \ln 4K} \text{ or } \sum_{j=1}^T \bar{X}_j > \frac{R^2}{5} \right); \end{aligned}$$

Therefore, probability event $E_{T-1} \setminus E_{-} \setminus E_{\otimes}$ implies

$$\begin{aligned} B_T &\leq \frac{R^2}{6} + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} + \frac{R^2}{5} = 2R^2; \\ R_T^2 &\stackrel{(122)}{\leq} \frac{R^2}{6} + 2R^2 \leq (2R)^2; \end{aligned}$$

which is equivalent to (118) and (119) for $t = T$, and

$$\text{Pf } E_T g > 1 \leq \text{Pf } E_{T-1} \setminus E_{-} \setminus E_{\otimes} g > 1 \leq \text{P } \bar{E}_{T-1} [\bar{E}_{-} [\bar{E}_{\otimes} > 1] \text{Pf } \bar{E}_{T-1} g] \text{Pf } \bar{E}_{-g} \text{Pf } \bar{E}_{\otimes g} > 1 \leq \frac{T}{K};$$

We have now completed the inductive part of our proof. That is, for all $k = 0; 1; \dots; K$, we have $\text{Pf } E_k g > 1 \leq \frac{k}{K}$. Notably, when $k = K$, we can conclude that with a probability of at least $1 - \frac{K}{K} = 0$:

$$f(y^K) - f(x) \stackrel{(121)}{\leq} \frac{6aLR^2}{K(K+3)}$$

and $f(x) \leq g_{k=0}^{K+1}; f(z) \leq g_{k=0}^K; f(y) \leq g_{k=0}^K \leq B_{2R}(x)$, which follows from (119).

Finally, if

$$a = \max \left\{ 97200 \ln^2 \frac{4K}{\delta}, \frac{1800\sigma(K+1)\sqrt{K}\sqrt{\ln \frac{4K}{\delta}}}{LR}, \frac{4b(K+2)^2}{15LR}, \frac{60b(K+2) \ln \frac{4K}{\delta}}{LR} \right\},$$

then with probability at least $1 - \delta$

$$\begin{aligned} f(y^K) - f(x) &\leq \frac{6aLR^2}{K(K+3)} \\ &= \mathcal{O} \left(\max \left\{ \frac{LR^2 \ln^2 \frac{K}{\delta}}{K^2}, \frac{\sigma R \sqrt{\ln \frac{K}{\delta}}}{\sqrt{K}}, \frac{bR \ln \frac{K}{\delta}}{K}, bR \right\} \right). \end{aligned}$$

□

C.2 Strongly Convex Case

In the strongly convex case, we consider a restarted version of SSTM, see Algorithm 1.

Algorithm 1 Restarted clipped-SSTM (R-clipped-SSTM) (Gorbunov et al., 2020)

Input: starting point x^0 , number of restarts τ , number of steps of clipped-SSTM between restarts $\{K_t\}_{t=1}^\tau$, stepsize parameters $\{a_t\}_{t=1}^\tau$, clipping levels $\{\lambda_k^1\}_{k=0}^{K_1-1}, \{\lambda_k^2\}_{k=0}^{K_2-1}, \dots, \{\lambda_k^\tau\}_{k=0}^{K_\tau-1}$, smoothness constant L .

1: $\hat{x}^0 = x^0$

2: **for** $t = 1, \dots, \tau$ **do**

3: Run clipped-SSTM for K_t iterations with stepsize parameter a_t , clipping levels $\{\lambda_k^t\}_{k=0}^{K_t-1}$, and starting point \hat{x}^{t-1} . Define the output of clipped-SSTM by \hat{x}^t .

4: **end for**

Output: \hat{x}^τ

The main result for R-clipped-SSTM is given below.

Theorem C.4. *Let Assumptions 2.2 and 2.3 with $\mu > 0$ hold on $Q = B_{3R}(x)$, where $R \geq \|x^0 - x\|$ and R-clipped-SSTM runs clipped-SSTM τ times. Assume that estimator $\nabla f_{k,t}(x^{k+1,t})$ used in clipped-SSTM at k -th iteration of stage t satisfies Assumption 5.1 with parameters b_k^t and σ_k^t such that*

$$b_k^t \leq \frac{15\mu R}{24 \cdot 2^{t+1}}. \quad (144)$$

Let

$$K_t = \left\lceil \max \left\{ 2160 \sqrt{\frac{LR_{t-1}^2}{\varepsilon_t}} \ln \frac{4320 \sqrt{LR_{t-1}^2} \tau}{\sqrt{\varepsilon_t} \delta}, 4 \left(\frac{5400 \sigma^t R_{t-1}}{\varepsilon_t} \right)^2 \ln \left(\frac{8\tau}{\delta} \left(\frac{5400 \sigma^t R_{t-1}}{\varepsilon_t} \right)^2 \right) \right\} \right\rceil, \quad (145)$$

$$\varepsilon_t = \frac{\mu R_{t-1}^2}{4}, \quad R_{t-1} = \frac{R}{2^{(t-1)/2}}, \quad \tau = \left\lceil \log_2 \frac{\mu R^2}{2\varepsilon} \right\rceil, \quad (146)$$

$$a_t = \max \left\{ 97200 \ln^2 \frac{4K_t \tau}{\delta}, \frac{1800\sigma(K_t+1)\sqrt{K_t}\sqrt{\ln \frac{4K_t \tau}{\delta}}}{LR_t}, \frac{4b_t(K_t+2)^2}{15LR_t}, \frac{60b_t(K_t+2) \ln \frac{4K_t}{\delta}}{LR_t} \right\}, \quad (147)$$

$$\lambda_k^t = \frac{R_t}{30\alpha_{k+1}^t \ln \frac{4K_t \tau}{\delta}} \quad (148)$$

for $t = 1, \dots, \tau$. Then to guarantee $f(\hat{x}^\tau) - f(x) \leq \varepsilon$ with probability $\geq 1 - \delta$ R-clipped-SSTM requires

$$\mathcal{O} \left(\max \left\{ \sqrt{\frac{L}{\mu}} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \ln \left(\frac{\sqrt{L}}{\sqrt{\mu} \delta} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right), \frac{\sigma^2}{\mu \varepsilon} \ln \left(\frac{\sigma^2}{\mu \varepsilon \delta} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right) \right\} \right) \quad (149)$$

iterations. Moreover, with probability $\geq 1 - \delta$ the iterates of R-clipped-SSTM at stage t stay in the ball $B_{2R_{t-1}}(x)$.

Proof. The proof of this theorem follows the same steps as the one given for Theorem F.3 from (Sadiev et al., 2023). By induction we derive that for any $t = 1, \dots, \tau$ with probability at least $1 - t\delta/\tau$ inequalities

$$f(\hat{x}^l) - f(x) \leq \varepsilon_l, \quad \|\hat{x}^l - x\|^2 \leq R_l^2 = \frac{R^2}{2^l} \quad (150)$$

hold for $l = 1, \dots, t$ simultaneously. We start with the base of the induction. Theorem C.3 implies that with probability at least $1 - \delta/\tau$

$$\begin{aligned} f(\hat{x}^1) - f(x) &\leq \frac{6a_1LR^2}{K_1(K_1+3)} \\ (147) \quad &\stackrel{=}{=} \max \left\{ \frac{583200LR^2 \ln^2 \frac{4K_1\tau}{\delta}}{K_1(K_1+3)}, \frac{10800\sigma R(K_1+1)\sqrt{K_1 \ln \frac{4K_1\tau}{\delta}}}{K_1(K_1+3)}, \right. \\ &\quad \left. \frac{24b_1R(K_1+2)^2}{15K_1(K_1+3)}, \frac{360b_1R(K_1+2) \ln \frac{4K_1}{\delta}}{K_1(K_1+3)} \right\} \\ &\leq \max \left\{ \frac{583200LR^2 \ln^2 \frac{4K_1\tau}{\delta}}{K_1^2}, \frac{10800\sigma R\sqrt{\ln \frac{4K_1\tau}{\delta}}}{\sqrt{K_1}}, \right. \\ &\quad \left. \frac{24b_1R}{15}, \frac{360b_1R \ln \frac{4K_1}{\delta}}{K_1} \right\} \\ (144),(145) \quad &\stackrel{\leq}{=} \varepsilon_1 = \frac{\mu R^2}{4} \end{aligned}$$

and, due to the strong convexity,

$$\|\hat{x}^1 - x\|^2 \leq \frac{2(f(\hat{x}^1) - f(x))}{\mu} \leq \frac{R^2}{2} = R_1^2.$$

The base of the induction is proven. Now, assume that the statement holds for some $t = T < \tau$, i.e., with probability at least $1 - T\delta/\tau$ inequalities

$$f(\hat{x}^l) - f(x) \leq \varepsilon_l, \quad \|\hat{x}^l - x\|^2 \leq R_l^2 = \frac{R^2}{2^l} \quad (151)$$

hold for $l = 1, \dots, T$ simultaneously. In particular, with probability at least $1 - T\delta/\tau$ we have $\|\hat{x}^T - x\|^2 \leq R_T^2$. Applying Theorem C.3 and using union bound for probability events, we get that with probability at least $1 - (T+1)\delta/\tau$

$$\begin{aligned} f(\hat{x}^{T+1}) - f(x) &\leq \frac{6a_{T+1}LR_T^2}{K_{T+1}(K_{T+1}+3)} \\ (147) \quad &\stackrel{=}{=} \max \left\{ \frac{583200LR_T^2 \ln^2 \frac{4K_{T+1}\tau}{\delta}}{K_{T+1}(K_{T+1}+3)}, \frac{10800\sigma R_T(K_{T+1}+1)\sqrt{K_{T+1} \ln \frac{4K_{T+1}\tau}{\delta}}}{K_{T+1}(K_{T+1}+3)}, \right. \\ &\quad \left. \frac{24b_{T+1}R_T(K_{T+1}+2)^2}{15K_{T+1}(K_{T+1}+3)}, \frac{360b_{T+1}R_T(K_{T+1}+2) \ln \frac{4K_{T+1}}{\delta}}{K_{T+1}(K_{T+1}+3)} \right\} \\ &\leq \max \left\{ \frac{583200LR_T^2 \ln^2 \frac{4K_{T+1}\tau}{\delta}}{K_{T+1}^2}, \frac{10800\sigma R_T\sqrt{\ln \frac{4K_{T+1}\tau}{\delta}}}{\sqrt{K_{T+1}}}, \right. \\ &\quad \left. \frac{24b_{T+1}R_T}{15}, \frac{360b_{T+1}R_T \ln \frac{4K_{T+1}}{\delta}}{K_{T+1}} \right\} \\ (144),(145) \quad &\stackrel{\leq}{=} \varepsilon_{T+1} = \frac{\mu R_T^2}{4} \end{aligned}$$

and, due to the strong convexity,

$$\|\hat{x}^{T+1} - x\|^2 \leq \frac{2(f(\hat{x}^{T+1}) - f(x))}{\mu} \leq \frac{R_T^2}{2} = R_{T+1}^2.$$

Thus, we finished the inductive part of the proof. In particular, with probability at least $1 - \delta$ inequalities

$$f(\hat{x}^l) - f(x) \leq \varepsilon_l, \quad \|\hat{x}^l - x\|^2 \leq R_l^2 = \frac{R^2}{2^l}$$

hold for $l = 1, \dots, \tau$ simultaneously, which gives for $l = \tau$ that with probability at least $1 - \delta$

$$f(\hat{x}^\tau) - f(x) \leq \varepsilon_\tau = \frac{\mu R_{\tau-1}^2}{4} = \frac{\mu R^2}{2^{\tau+1}} \stackrel{(146)}{\leq} \varepsilon.$$

It remains to calculate the overall number of iterations during all runs of clipped-SSTM. We have

$$\begin{aligned} \sum_{t=1}^{\tau} K_t &= \mathcal{O} \left(\sum_{t=1}^{\tau} \max \left\{ \sqrt{\frac{LR_{t-1}^2}{\varepsilon_t}} \ln \left(\frac{\sqrt{LR_{t-1}^2} \tau}{\sqrt{\varepsilon_t} \delta} \right), \left(\frac{\sigma R_{t-1}}{\varepsilon_t} \right)^2 \ln \left(\frac{\tau}{\delta} \left(\frac{\sigma R_{t-1}}{\varepsilon_t} \right)^2 \right) \right\} \right) \\ &= \mathcal{O} \left(\sum_{t=1}^{\tau} \max \left\{ \sqrt{\frac{L}{\mu}} \ln \left(\frac{\sqrt{L} \tau}{\sqrt{\mu} \delta} \right), \left(\frac{\sigma}{\mu R_{t-1}} \right)^2 \ln \left(\frac{\tau}{\delta} \left(\frac{\sigma}{\mu R_{t-1}} \right)^2 \right) \right\} \right) \\ &= \mathcal{O} \left(\max \left\{ \tau \sqrt{\frac{L}{\mu}} \ln \left(\frac{\sqrt{L} \tau}{\sqrt{\mu} \delta} \right), \sum_{t=1}^{\tau} \left(\frac{\sigma \cdot 2^{t/2}}{\mu R} \right)^2 \ln \left(\frac{\tau}{\delta} \left(\frac{\sigma \cdot 2^{t/2}}{\mu R} \right)^2 \right) \right\} \right) \\ &= \mathcal{O} \left(\max \left\{ \sqrt{\frac{L}{\mu}} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \ln \left(\frac{\sqrt{L}}{\sqrt{\mu} \delta} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right), \left(\frac{\sigma}{\mu R} \right)^2 \ln \left(\frac{\tau}{\delta} \left(\frac{\sigma \cdot 2^{\tau/2}}{\mu R} \right)^2 \right) \sum_{t=1}^{\tau} 2^t \right\} \right) \\ &= \mathcal{O} \left(\max \left\{ \sqrt{\frac{L}{\mu}} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \ln \left(\frac{\sqrt{L}}{\sqrt{\mu} \delta} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right), \left(\frac{\sigma}{\mu R} \right)^2 \ln \left(\frac{\tau}{\delta} \left(\frac{\sigma}{\mu R} \right)^2 \cdot 2 \right) 2^\tau \right\} \right) \\ &= \mathcal{O} \left(\max \left\{ \sqrt{\frac{L}{\mu}} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \ln \left(\frac{\sqrt{L}}{\sqrt{\mu} \delta} \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right), \left(\frac{\sigma^2}{\mu \varepsilon} \right) \ln \left(\frac{1}{\delta} \left(\frac{\sigma^2}{\mu \varepsilon} \right) \ln \left(\frac{\mu R^2}{\varepsilon} \right) \right) \right\} \right), \end{aligned}$$

which concludes the proof. \square

D PROPERTIES OF HERMITE POLYNOMIALS

This section collects some properties of Hermite polynomials which are used in the proof of Lemma A.4. First, let us recall the definition. There are two versions of Hermite polynomials, which are referred to as “physicist’s” and “probabilist’s”, given by

$$\mathcal{H}_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} e^{-x^2} \quad \text{and} \quad \mathcal{H}_n(x) = (-1)^n e^{x^2/2} \frac{d^n}{dx^n} e^{-x^2/2} \quad n \in \mathbb{N},$$

respectively. Obviously, for any positive integer n , the following relation holds true:

$$\mathcal{H}_n(x) \equiv 2^{-n/2} \mathcal{H}_n\left(\frac{x}{\sqrt{2}}\right). \quad (152)$$

Hence, “probabilist’s” Hermite polynomials inherit all the properties of “physicist’s” ones. In (Indritz, 1961), the author proved that

$$\max_{x \in \mathbb{R}} \left| \mathcal{H}_n(x) e^{-x^2/2} \right| \leq \sqrt{2^n \cdot n!} \quad \text{for all } n \in \mathbb{N}.$$

In view of (152), this implies that

$$\max_{x \in \mathbb{R}} \left| \mathcal{H}_n(x) e^{-x^2/4} \right| = 2^{-n/2} \max_{x \in \mathbb{R}} \left| \mathcal{H}_n\left(\frac{x}{\sqrt{2}}\right) e^{-x^2/4} \right| \leq \sqrt{n!} \quad \text{for all } n \in \mathbb{N}.$$

E NUMERICAL EXPERIMENTS: ADDITIONAL DETAILS

For every combination of noise distribution and method, we tuned optimal parameters for 70000 steps and ran methods on 95000 steps, where one step is one oracle call.

The optimal values of the learning rate and the clipping parameter were selected via grid search over the sets $\{0.002, 0.004, 0.008, 0.01, 0.02, 0.04\}$ and $\{0.75, 1, 1.5, 2, 4, 8\}$, respectively.

Distribution	Method	Learning Rate	Clipping parameter
Cauchy	clipped-MB-SGD	0.004	4
	Med-MB-SGD	0.002	-
	MB-clipped-SGD	0.01	4
	clipped-Med-MB-SGD	0.002	2
	SMoM-MB-SGD	0.002	-
	clipped-SMoM-MB-SGD	0.008	1
Cauchy + Exponential	clipped-MB-SGD	0.008	1.5
	Med-MB-SGD	0.002	-
	MB-clipped-SGD	0.04	4
	clipped-Med-MB-SGD	0.002	8
	SMoM-MB-SGD	0.002	-
	clipped-SMoM-MB-SGD	0.002	8
Cauchy + Pareto	clipped-MB-SGD	0.008	1.5
	Med-MB-SGD	0.002	-
	MB-clipped-SGD	0.02	0.75
	clipped-Med-MB-SGD	0.002	8
	SMoM-MB-SGD	0.002	-
	clipped-SMoM-MB-SGD	0.008	1

Table 1: Optimal parameters for different distributions and methods.

We also provide plots, reflecting the dependence of the error on the number of iterations, where one iteration is one method’s update, see Figure 2 below. As we can see, clipped-SMoM-MB-SGD converges much faster than the competitors due to the larger batch size.

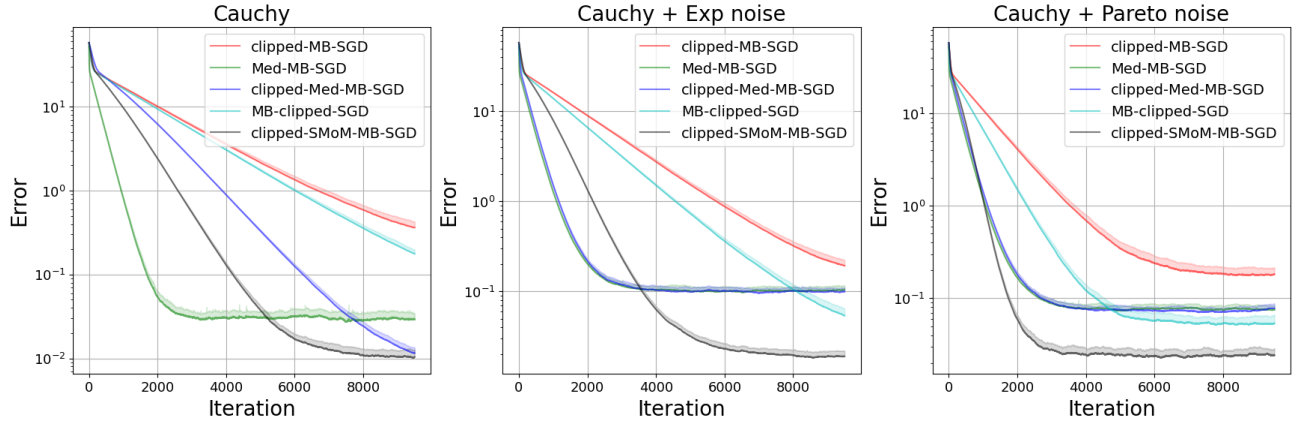


Figure 2: Dependence of the mean error on the number of iterations with a standard deviation upper bound.

Finally, according to our theoretical findings, the error bound for the mini-batched SGD with clipped smoothed median of means grows logarithmically with $1/\delta$. In Figure 3, we plot the dependence of the confidence interval width on the number of iterations to illustrate this point.

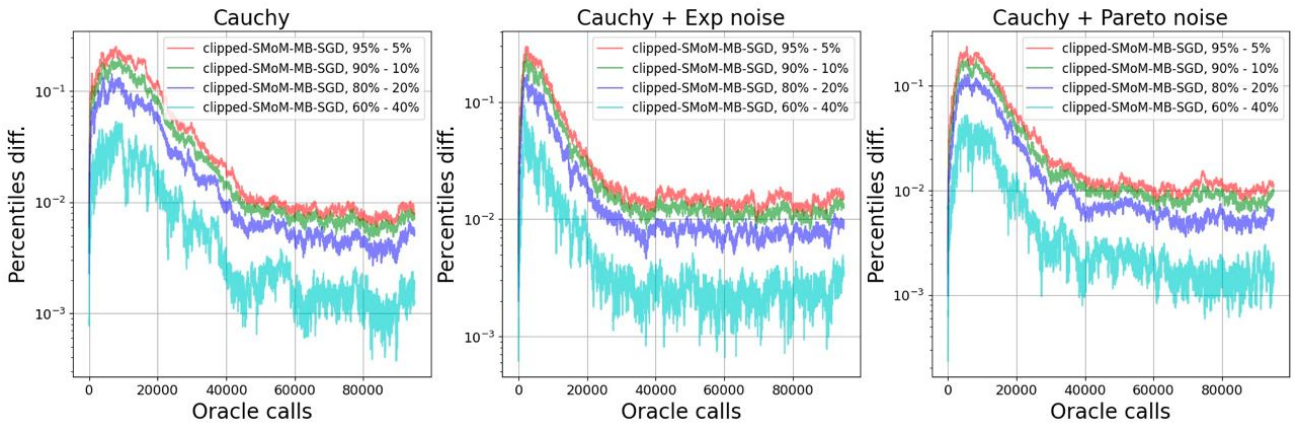


Figure 3: Dependence of the confidence interval width for the error of mini-batched SGD with clipped smoothed median of means on the number of iterations.