
Gradient Clipping Helps in Non-Smooth Stochastic Optimization with Heavy-Tailed Noise

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Thanks to their practical efficiency and random nature of the data, stochastic
2 first-order methods are standard for training large-scale machine learning models.
3 Random behavior may cause a particular run of an algorithm to result in a highly
4 suboptimal objective value, whereas theoretical guarantees are usually proved
5 for the expectation of the objective value. Thus, it is essential to theoretically
6 guarantee that algorithms provide small objective residual with high probability.
7 Existing methods for non-smooth stochastic convex optimization have complexity
8 bounds with the dependence on the confidence level that is either negative-power or
9 logarithmic but under an additional assumption of sub-Gaussian (light-tailed) noise
10 distribution that may not hold in practice, e.g., in several NLP tasks. In our paper,
11 we resolve this issue and derive the first high-probability convergence results with
12 logarithmic dependence on the confidence level for non-smooth convex stochastic
13 optimization problems with non-sub-Gaussian (heavy-tailed) noise. To derive our
14 results, we propose novel stepsize rules for two stochastic methods with gradient
15 clipping. Moreover, our analysis works for generalized smooth objectives with
16 Hölder-continuous gradients, and for both methods, we provide an extension for
17 strongly convex problems. Finally, our results imply that the first (accelerated)
18 method we consider also has optimal iteration and oracle complexity in all the
19 regimes, and the second one is optimal in the non-smooth setting.

20 1 Introduction

21 Stochastic first-order optimization methods like SGD [32], Adam [20], and their various modifi-
22 cations are extremely popular in solving a number of different optimization problems, especially
23 those appearing in statistics [36], machine learning, and deep learning [13]. The success of these
24 methods in real-world applications motivates the researchers to investigate theoretical properties
25 for the methods and to develop new ones with better convergence guarantees. Typically, stochastic
26 methods are analyzed in terms of the convergence in expectation (see [12, 24, 15] and references
27 therein), whereas high-probability complexity results are established much rarely. However, as
28 illustrated in [14], guarantees in terms of the convergence in expectation have much worse correlation
29 with the real behavior of the methods than high-probability convergence guarantees when the noise
30 in the stochastic gradients has *heavy-tailed distribution*.

31 Recent studies [35, 34, 41] show that in several popular problems such as training BERT [37] on
32 Wikipedia dataset the noise in the stochastic gradients is heavy-tailed. Moreover, in [41], the authors
33 justify empirically that in such cases SGD works significantly worse than clipped-SGD [30] and
34 Adam. Therefore, it is important to theoretically study the methods' convergence when the noise is
35 heavy-tailed.

36 For convex and strongly convex problems with Lipschitz continuous gradient, i.e., smooth convex and
37 strongly convex problems, this question was properly addressed in [25, 3, 14] where the first high-
38 probability complexity bounds with logarithmic dependence on the confidence level were derived
39 for the stochastic problems with heavy-tailed noise. However, a number of practically important
40 problems are non-smooth *on the whole space* [40, 22]. For example, in deep neural network training,
41 the loss function often grows polynomially fast when the norm of the network’s weights goes to
42 infinity. Moreover, non-smoothness of the activation functions such as ReLU or loss functions such
43 as hinge loss implies the non-smoothness of the whole problem. While being well-motivated by
44 practical applications, the existing high-probability convergence guarantees for stochastic first-order
45 methods applied to solve non-smooth convex optimization problems with heavy-tailed noise depend
46 on the negative power of the confidence level that dramatically increases the number of iterations
47 required to obtain high accuracy of the solution with probability close to one. Such a discrepancy in
48 the theory between algorithms for stochastic smooth and non-smooth problems leads us to the natural
49 question: *is it possible to obtain high-probability complexity bounds with logarithmic dependence*
50 *on the confidence level for **non-smooth** convex stochastic problems with heavy-tailed noise?* In this
51 paper, we give a positive answer to this question. To achieve this we focus on gradient clipping
52 methods [30, 10, 23, 22, 40, 41].

53 1.1 Preliminaries

54 Before we describe our contributions in detail, we formally state the considered setup.

55 **Stochastic optimization.** We focus on the following problem

$$\min_{x \in \mathbb{R}^n} f(x), \quad f(x) = \mathbb{E}_\xi [f(x, \xi)], \quad (1)$$

56 where $f(x)$ is a convex but possibly non-smooth function. Next, we assume that at each point $x \in \mathbb{R}^n$
57 we have an access to the unbiased estimator $\nabla f(x, \xi)$ of $\nabla f(x)$ with uniformly bounded variance

$$\mathbb{E}_\xi [\nabla f(x, \xi)] = \nabla f(x), \quad \mathbb{E}_\xi \left[\|\nabla f(x, \xi) - \nabla f(x)\|_2^2 \right] \leq \sigma^2, \quad \sigma > 0. \quad (2)$$

58 This assumption on the stochastic oracle is widely used in stochastic optimization literature [11,
59 12, 19, 21, 26]. We emphasize that we do not assume that the stochastic gradients have so-called
60 “light tails” [21], i.e., sub-Gaussian noise distribution meaning that $\mathbb{P}\{\|\nabla f(x, \xi) - \nabla f(x)\|_2 > b\} \leq$
61 $2 \exp(-b^2/(2\sigma^2))$ for all $b > 0$.

62 **Level of smoothness.** Finally, we assume that function f has (ν, M_ν) -Hölder continuous gradients
63 on a compact set $Q \subseteq \mathbb{R}^n$ for some $\nu \in [0, 1]$, $M_\nu > 0$ meaning that

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq M_\nu \|x - y\|_2^\nu \quad \forall x, y \in Q. \quad (3)$$

64 When $\nu = 1$ inequality (3) implies M_1 -smoothness of f , and when $\nu = 0$ we have that $\nabla f(x)$
65 has bounded variation which is equivalent to being uniformly bounded. Moreover, when $\nu = 0$
66 differentiability of f is not needed, and one can assume uniform boundedness of the subgradients of
67 f . Linear regression in the case when the noise has generalized Gaussian distribution (Example 4.4
68 from [2]) serves as a natural example of the situation with $\nu \in (0, 1)$. Moreover, when (3) holds for
69 $\nu = 0$ and $\nu = 1$ simultaneously then it holds for all $\nu \in [0, 1]$ with $M_\nu \leq M_0^{1-\nu} M_1^\nu$ [28]. As we
70 show in our results, the set Q should contain the ball centered at the solution x^* of (1) with radius
71 $2R_0 = 2\|x^0 - x^*\|_2$, where x^0 is a starting point of the method, i.e., our analysis does not require (3)
72 to hold on \mathbb{R}^n .

73 **High-probability convergence.** For a given accuracy $\varepsilon > 0$ and confidence level $\beta \in (0, 1)$ we
74 are interested in finding ε -solutions of problem (1) with probability at least $1 - \beta$, i.e., such \hat{x} that
75 $\mathbb{P}\{f(\hat{x}) - f(x^*) \leq \varepsilon\} \geq 1 - \beta$. For brevity, we will call such (in general, random) points \hat{x} as
76 (ε, β) -solution of (1). Moreover, by high-probability complexity of a stochastic method \mathcal{M} we mean
77 the sufficient number of oracle calls, i.e., number of $\nabla f(x, \xi)$ computations, needed to guarantee that
78 the output of \mathcal{M} is an (ε, β) -solution of (1).

Table 1: Summary of known and new high-probability complexity bounds for solving (1) with f being **convex** and having (ν, M_ν) -Hölder continuous gradients. Columns: “Ref.” = reference, “Complexity” = high-probability complexity (ε – accuracy, β – confidence level, numerical constants and logarithmic factors are omitted), “HT” = heavy-tailed noise, “UD” = unbounded domain, “HCC” = Hölder continuity of the gradient is required only on the compact set. The results labeled by \clubsuit are obtained from the convergence guarantees in expectation via Markov’s inequality. Negative-power dependencies on the confidence level β are colored in red.

Method	Ref.	Complexity	ν	HT?	UD?	HCC?
SGD	[26]	$\max \left\{ \frac{M_0^2 R_0^2}{\varepsilon^2}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	0	✗	✓	✗
AC-SA	[11, 21]	$\max \left\{ \sqrt{\frac{M_1 R_0^2}{\varepsilon}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	1	✗	✓	✗
SIGMA	[6]	$\max \left\{ \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	[0, 1]	✗	✓	✗
SGD	[26] \clubsuit	$\max \left\{ \frac{M_0^2 R_0^2}{\beta^2 \varepsilon^2}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2} \right\}$	0	✓	✗	✗
AC-SA	[11, 21] \clubsuit	$\max \left\{ \sqrt{\frac{M_1 R_0^2}{\beta \varepsilon}}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2} \right\}$	1	✓	✓	✗
SIGMA	[6] \clubsuit	$\max \left\{ \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\beta^{\frac{2}{1+3\nu}} \varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2} \right\}$	[0, 1]	✓	✓	✗
clipped-SSTM	[14]	$\max \left\{ \sqrt{\frac{M_1 R_0^2}{\varepsilon}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	1	✓	✓	✗
clipped-SGD	[14]	$\max \left\{ \frac{M_1 R_0^2}{\varepsilon}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	1	✓	✓	✗
clipped-SSTM	Thm. 2.2	$\max \left\{ \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	[0, 1]	✓	✓	✓
clipped-SGD	Thm. 3.1	$\max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\}$	[0, 1]	✓	✓	✓

79 1.2 Contributions

- 80 • We propose novel stepsize rules for **clipped-SSTM** [14] to handle the problems with Hölder
81 continuous gradients and derive high-probability complexity guarantees for convex stochastic
82 optimization problems without using “light tails” assumption, i.e., we prove that our version of
83 **clipped-SSTM**

$$\mathcal{O} \left(\max \left\{ D \ln \frac{2(1+\nu)}{1+3\nu} \frac{D}{\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{D}{\beta} \right\} \right), \quad D = \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}$$

84 high-probability complexity. Unlike all previous high-probability complexity results in this setup
85 with $\nu < 1$ (see Tbl. 1), our result depends only logarithmically on the confidence level β that
86 is highly important when β is small. Moreover, up to the difference in logarithmic factors the
87 derived complexity guarantees meet the known lower bounds [21, 17] obtained for the problems
88 with light-tailed noise. In particular, when $\nu = 1$ we recover accelerated convergence rate [29, 21].
89 That is, neglecting the logarithmic factors our results are unimprovable and, surprisingly coincide
90 with the best-known results in the “light-tailed case”.

- 91 • We derive the first high-probability complexity bounds for **clipped-SGD** when the objective
92 functions is convex with (ν, M_ν) -Hölder continuous gradient and the noise is heavy tailed., i.e., we
93 derive

$$\mathcal{O} \left(\max \left\{ D^2, \max \left\{ D^{1+\nu}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \right\} \ln \frac{D^2 + D^{1+\nu}}{\beta} \right\} \right), \quad D = \frac{M_\nu^{\frac{1}{1+\nu}} R_0}{\varepsilon^{\frac{1}{1+\nu}}}$$

94 high-probability complexity bound. Interestingly, when $\nu = 0$ the derived bound for **clipped-SGD**
95 has better dependence on the logarithms than the corresponding one for **clipped-SSTM**. Moreover,
96 neglecting the dependence on ε under the logarithm, our bound for **clipped-SGD** has the same

Table 2: Summary of known and new high-probability complexity bounds for solving (1) with f being μ -strongly convex and having (ν, M_ν) -Hölder continuous gradients. Columns: “Ref.” = reference, “Complexity” = high-probability complexity (ε – accuracy, β – confidence level, numerical constants and logarithmic factors are omitted), “HT” = heavy-tailed noise, “UD” = unbounded domain, “HCC” = Hölder continuity of the gradient is required only on the compact set. The results labeled by \clubsuit are obtained from the convergence guarantees in expectation via Markov’s inequality. Negative-power dependencies on the confidence level β are colored in red.

Method	Ref.	Complexity	ν	HT?	UD?	HCC?
SGD	[26]	$\max \left\{ \frac{M_0^2}{\mu\varepsilon}, \frac{\sigma^2}{\mu\varepsilon} \right\}$	0	✗	✓	✗
AC-SA	[11, 21]	$\max \left\{ \sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\varepsilon} \right\}$	1	✗	✓	✗
SIGMA	[6]	$\max \left\{ \hat{N}, \frac{\sigma^2}{\mu\varepsilon} \right\},$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}} \right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu} \varepsilon^{1-\nu}} \right)^{\frac{1}{1+3\nu}}$	$[0, 1]$	✗	✓	✗
SGD	[26] \clubsuit	$\max \left\{ \frac{M_0^2}{\mu\beta\varepsilon}, \frac{\sigma^2}{\mu\beta\varepsilon} \right\}$	0	✓	✗	✗
AC-SA	[11, 21] \clubsuit	$\max \left\{ \sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\beta\varepsilon} \right\}$	1	✓	✓	✗
SIGMA	[6] \clubsuit	$\max \left\{ \hat{N}, \frac{\sigma^2}{\mu\hat{\varepsilon}} \right\}, \hat{\varepsilon} = \beta\varepsilon,$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}} \right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu} \hat{\varepsilon}^{1-\nu}} \right)^{\frac{1}{1+3\nu}}$	$[0, 1]$	✓	✓	✗
R-clipped-SSTM	[14]	$\max \left\{ \sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\varepsilon^2} \right\}$	1	✓	✓	✗
R-clipped-SGD	[14]	$\max \left\{ \frac{M_1}{\mu}, \frac{\sigma^2}{\mu\varepsilon^2} \right\}$	1	✓	✓	✗
R-clipped-SSTM	Thm. 2.1	$\max \left\{ \hat{N}, \frac{\sigma^2}{\mu\varepsilon} \right\},$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}} \right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu} \varepsilon^{1-\nu}} \right)^{\frac{1}{1+3\nu}}$	$[0, 1]$	✓	✓	✓
R-clipped-SGD	Thm. 3.2	$\max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}}}{\mu^{\frac{2}{1+\nu}} R_0^{\frac{2(1-\nu)}{1+\nu}}}, \frac{M_\nu^{\frac{2}{1+\nu}}}{\mu\varepsilon^{\frac{1-\nu}{1+\nu}}}, \frac{\sigma^2}{\mu\varepsilon} \right\}$	$[0, 1]$	✓	✓	✓

97 dependence on the confidence level as the tightest known result in this case under the “light tails”
98 assumption [16].

99 • Using restarts technique we extend the obtained results for clipped-SSTM and clipped-SGD to
100 the strongly convex case (see Tbl. 2). As in the convex case, the obtained results are superior to all
101 previous known results in the general setup we consider.

102 • As one of the key contributions of this work, we emphasize that in our theoretical results it is
103 sufficient to assume Hölder continuity of the gradients of f only on the ball with radius $2R_0 =$
104 $2\|x^0 - x^*\|_2$ and centered at a solution of the problem. This makes our results applicable to much
105 larger class of problems than functions with Hölder continuous gradients on \mathbb{R}^n , e.g., our analysis
106 works even for polynomially growing objectives.

107 • To test the performance of the considered methods we conduct several numerical experiments
108 on image classification and NLP tasks, and observe that 1) clipped-SSTM and clipped-SGD
109 show a comparable performance with SGD on the image classification task, when the noise
110 distribution is almost sub-Gaussian, 2) converge much faster than SGD on the NLP task, when the
111 noise distribution is heavy-tailed, and 3) clipped-SSTM achieves a comparable performance with
112 Adam on the NLP task enjoying both the best known theoretical guarantees and good practical
113 performance.

114 1.3 Related work

115 **Light-tailed noise.** The theory of high-probability complexity bounds for convex stochastic op-
116 timization with light-tailed noise is well-developed. Lower bounds and optimal methods for the
117 problems with (ν, M_ν) -Hölder continuous gradients are obtained in [26] for $\nu = 0$, and in [11] for
118 $\nu = 1$. Up to the logarithmic dependencies these high-probability convergence bounds coincide with

119 the corresponding results for the convergence in expectation (see first two rows of Tbl. 1) While not
 120 being directly derived in the literature, the lower bound for the case when $\nu \in (0, 1)$ can be obtained
 121 as a combination of lower bounds in the deterministic [27, 17] and smooth stochastic settings [11].
 122 The corresponding optimal methods are analyzed in [4, 6] through the lens of inexact oracle.

123 **Heavy-tailed noise.** Unlike in the “light-tailed” case, the first theoretical guarantees with reasonable
 124 dependence on both the accuracy ε and the confidence level β appeared just recently. In [25], the
 125 first such results without acceleration [29] were derived for Mirror Descent with special truncation
 126 technique for smooth ($\nu = 1$) convex problems on a bounded domain, and then were accelerated and
 127 extended in [14]. For the strongly convex problems the first accelerated high-probability convergence
 128 guarantees were obtained in [3] for the special method called proxBOOST requiring solving auxiliary
 129 nontrivial problems at each iteration. These bounds were tightened in [14].

130 In contrast, for the case when $\nu < 1$ and, in particular, when $\nu = 0$ the best-known high-probability
 131 complexity bounds suffer from the negative-power dependence on the confidence level β , i.e., have
 132 a factor $1/\beta^\alpha$ for some $\alpha > 0$, that affects the convergence rate dramatically for small enough
 133 β . Without additional assumptions on the tails these results are obtained via Markov’s inequality
 134 $\mathbb{P}\{f(x) - f(x^*) > \varepsilon\} < \mathbb{E}[f(x) - f(x^*)]/\varepsilon$ from the guarantees for the convergence in expectation to
 135 the accuracy $\varepsilon\beta$, see the results labeled by \clubsuit in Tbl. 1. Under an additional assumption on noise
 136 tails that $\mathbb{P}\{\|\nabla f(x, \xi) - \nabla f(x)\|_2^2 > s\sigma^2\} = O(s^{-\alpha})$ for $\alpha > 2$ these results can be tightened [9]
 137 when $\nu = 0$ as $O\left(M_0^2 R_0^2 \max\left\{\ln(\beta^{-1})/\varepsilon^2, (1/\beta\varepsilon^\alpha)^{2/(3\alpha-2)}\right\}\right)$ without removing the negative-power
 138 dependence on the confidence level β . Different stepsize policies allow to change the last term in
 139 max to $\beta^{-\frac{1}{2\alpha-1}} \varepsilon^{-\frac{2\alpha}{2\alpha-1}}$ without removing the negative-power dependence on β .

140 **Gradient clipping.** The methods based on gradient clipping [30] and normalization [18] are popular
 141 in different machine learning and deep learning tasks due to their robustness in practice to the noise
 142 in the stochastic gradients and rapid changes of the objective function [13]. In [40, 22], clipped-GD
 143 and clipped-SGD are theoretically studied in applications to non-smooth problems that can grow
 144 polynomially fast when $\|x - x^*\|_2 \rightarrow \infty$ showing the superiority of gradient clipping methods
 145 to the methods without clipping. The results from [40] are obtained for non-convex problems
 146 with almost surely bounded noise, and in [22], the authors derive the stability and expectation
 147 convergence guarantees for strongly convex under assumption that the central p -th moment of the
 148 stochastic gradient is bounded for $p \geq 2$. Since the authors of [22] do not provide convergence
 149 guarantees with explicit dependencies on all important parameters of the problem it complicates direct
 150 comparison with our results. Nevertheless, convergence guarantees from [22] are sub-linear and are
 151 given for the convergence in expectation, and, as a consequence, the corresponding high-probability
 152 convergence results obtained via Markov’s inequality also suffer from negative-power dependence on
 153 the confidence level. Next, in [41], the authors establish several expectation convergence guarantees
 154 for clipped-SGD and prove their optimality in the non-convex case under assumption that the central
 155 α -moment of the stochastic gradient is uniformly bounded, where $\alpha \in (1, 2]$. It turns out that
 156 clipped-SGD is able to converge even when $\alpha < 2$, whereas vanilla SGD can diverge in this setting.

157 2 Clipped Stochastic Similar Triangles Method

158 In this section, we propose a novel variation of Clipped Stochastic Similar Triangles Method [14]
 159 adjusted to the class of objectives with Hölder continuous gradients (clipped-SSTM, see Alg. 1).

160 The method is based on the clipping of the stochastic gradients:

$$\text{clip}(\nabla f(x, \xi), \lambda) = \min\left\{1, \frac{\lambda}{\|\nabla f(x, \xi)\|_2}\right\} \nabla f(x, \xi) \quad (4)$$

161 where $\nabla f(x, \xi) = \frac{1}{m} \sum_{i=1}^m \nabla f(x, \xi_i)$ is a mini-batched stochastic gradient. Gradient clipping
 162 ensures that the resulting vector has a norm bounded by the clipping level λ . Since the clipped
 163 stochastic gradient cannot have arbitrary large norm, the clipping helps to avoid unstable behavior of
 164 the method when the noise is heavy-tailed and the clipping level λ is properly adjusted.

165 However, unlike the stochastic gradient, clipped stochastic gradient is a *biased* estimate of $\nabla f(x)$:
 166 the smaller the clipping level the larger the bias. The biasedness of the clipped stochastic gradient

Algorithm 1 Clipped Stochastic Similar Triangles Method (clipped-SSTM): case $\nu \in [0, 1]$

Input: starting point x^0 , number of iterations N , batchsizes $\{m_k\}_{k=1}^N$, stepsize parameter α , clipping parameter B , Hölder exponent $\nu \in [0, 1]$.

- 1: Set $A_0 = \alpha_0 = 0, y^0 = z^0 = x^0$
 - 2: **for** $k = 0, \dots, N - 1$ **do**
 - 3: Set $\alpha_{k+1} = \alpha(k+1)^{\frac{2\nu}{1+\nu}}, A_{k+1} = A_k + \alpha_{k+1}, \lambda_{k+1} = \frac{B}{\alpha_{k+1}}$
 - 4: $x^{k+1} = (A_k y^k + \alpha_{k+1} z^k) / A_{k+1}$
 - 5: Draw mini-batch m_k of fresh i.i.d. samples $\xi_1^k, \dots, \xi_{m_k}^k$ and compute $\nabla f(x^{k+1}, \xi^k) = \frac{1}{m_k} \sum_{i=1}^{m_k} \nabla f(x^{k+1}, \xi_i^k)$
 - 6: Compute $\tilde{\nabla} f(x^{k+1}, \xi^k) = \text{clip}(\nabla f(x^{k+1}, \xi^k), \lambda_{k+1})$ using (4)
 - 7: $z^{k+1} = z^k - \alpha_{k+1} \tilde{\nabla} f(x^{k+1}, \xi^k)$
 - 8: $y^{k+1} = (A_k y^k + \alpha_{k+1} z^{k+1}) / A_{k+1}$
 - 9: **end for**
- Output:** y^N
-

167 complicates the analysis of the method. On the other hand, to circumvent the negative effect of
 168 the heavy-tailed noise on the high-probability convergence one should choose λ to be not too large.
 169 Therefore, the question on the appropriate choice of the clipping level is highly non-trivial.

170 Fortunately, there exists a simple but insightful observation that helps us to obtain the right formula
 171 for the clipping level λ_k in clipped-SSTM: if λ_k is chosen in such a way that $\|\nabla f(x^k)\|_2 \leq \lambda_k/2$
 172 with high probability, then for the realizations $\nabla f(x^{k+1}, \xi^k)$ of the mini-batched stochastic gradient
 173 such that $\|\nabla f(x^{k+1}, \xi^k) - \nabla f(x^{k+1})\|_2 \leq \lambda_k/2$ the clipping is an identity operator. Next, if the
 174 probability mass of such realizations is big enough then the bias of the clipped stochastic gradient is
 175 properly bounded that helps to derive needed convergence guarantees. It turns out that the choice
 176 $\lambda_k \sim 1/\alpha_k$ ensures the method convergence with needed rate and high enough probability.

177 Guided by this observation we derive the precise expressions for all the parameters of clipped-SSTM
 178 and derive high-probability complexity bounds for the method. Below we provide a simplified version
 179 of the main result for clipped-SSTM in the convex case. The complete formulation and the full proof
 180 of the theorem are deferred to Appendix B.1 (see Thm. B.1).

181 **Theorem 2.1.** Assume that function f is convex and its gradient satisfy (3) with $\nu \in [0, 1], M_\nu > 0$
 182 on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist such
 183 a choice of parameters that clipped-SSTM achieves $f(y^N) - f(x^*) \leq \varepsilon$ with probability at least

184 $1 - \beta$ after $\mathcal{O}\left(D \ln \frac{2(1+\nu)}{1+3\nu} \frac{D}{\beta}\right)$ iterations with $D = \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}$ and requires

$$\mathcal{O}\left(\max\left\{D \ln \frac{2(1+\nu)}{1+3\nu} \frac{D}{\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{D}{\beta}\right\}\right) \text{ oracle calls.} \quad (5)$$

185 The obtained result has only logarithmic dependence on the confidence level β and optimal depen-
 186 dence on the accuracy ε up to logarithmic factors [21, 17] for all $\nu \in [0, 1]$. Moreover, we emphasize
 187 that our result does not require f to have (ν, M_ν) -Hölder continuous gradient on the whole space.
 188 This is because we prove that for the proposed choice of parameters the iterates of clipped-SSTM
 189 stay inside the ball $B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$ with probability at least $1 - \beta$, and,
 190 as a consequence, Hölder continuity of the gradient is required only inside this ball. In particular,
 191 this means that the better starting point leads not only to the reduction of R_0 , but also it can reduce
 192 M_ν . Moreover, our result is applicable to much wider class of functions than the standard class of
 193 functions with Hölder continuous gradients in \mathbb{R}^n , e.g., to the problems with polynomial growth.

194 For the strongly convex problems, we consider restarted version of Alg. 1 (R-clipped-SSTM, see
 195 Alg. 2) and derive high-probability complexity result for this version. Below we provide a simplified
 196 version of the result. The complete formulation and the full proof of the theorem are deferred to
 197 Appendix B.2 (see Thm. B.2).

198 **Theorem 2.2.** Assume that function f is μ -strongly convex and its gradient satisfy (3) with $\nu \in [0, 1]$,
 199 $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist

Algorithm 2 Restarted clipped-SSTM (R-clipped-SSTM): case $\nu \in [0, 1]$

Input: starting point x^0 , number of restarts τ , number of steps of clipped-SSTM in restarts $\{N_t\}_{t=1}^\tau$, batchsizes $\{m_k^1\}_{k=1}^{N_1-1}, \{m_k^2\}_{k=1}^{N_2-1}, \dots, \{m_k^\tau\}_{k=1}^{N_\tau-1}$, stepsize parameters $\{\alpha^t\}_{t=1}^\tau$, clipping parameters $\{B_t\}_{t=1}^\tau$, Hölder exponent $\nu \in [0, 1]$.

- 1: $\hat{x}^0 = x^0$
- 2: **for** $t = 1, \dots, \tau$ **do**
- 3: Run clipped-SSTM (Alg. 1) for N_t iterations with batchsizes $\{m_k^t\}_{k=1}^{N_t-1}$, stepsize parameter α_t , clipping parameter B_t , and starting point \hat{x}^{t-1} . Define the output of clipped-SSTM by \hat{x}^t .
- 4: **end for**

Output: \hat{x}^τ

200 *such a choice of parameters that R-clipped-SSTM achieves $f(\hat{x}^\tau) - f(x^*) \leq \varepsilon$ with probability at*
 201 *least $1 - \beta$ after*

$$\hat{N} = O\left(D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}\right), \quad D = \max\left\{\left(\frac{M_\nu}{\mu R_0^{1-\nu}}\right)^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, \left(\frac{M_\nu^2}{\mu^{1+\nu} \varepsilon^{1-\nu}}\right)^{\frac{1}{1+3\nu}}\right\} \quad (6)$$

202 *iterations of Alg. 1 in total and requires*

$$O\left(\max\left\{D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}, \frac{\sigma^2}{\mu \varepsilon} \ln \frac{D}{\beta}\right\}\right) \text{ oracle calls.} \quad (7)$$

203 Again, the obtained result has only logarithmic dependence on the confidence level β and, as our
 204 result in the convex case, it has optimal dependence on the accuracy ε up to logarithmic factors
 205 depending on β [21, 17] for all $\nu \in [0, 1]$.

206 3 SGD with clipping

207 In this section, we present a new variant of clipped-SGD [30] properly adjusted to the class of
 208 objectives with (ν, M_ν) -Hölder continuous gradients (see Alg. 3).

Algorithm 3 Clipped Stochastic Gradient Descent (clipped-SGD): case $\nu \in [0, 1]$

Input: starting point x^0 , number of iterations N , batchsize m , stepsize γ , clipping parameter $B > 0$.

- 1: **for** $k = 0, \dots, N - 1$ **do**
- 2: Draw mini-batch of m fresh i.i.d. samples ξ_1^k, \dots, ξ_m^k and compute $\nabla f(x^{k+1}, \xi^k) = \frac{1}{m} \sum_{i=1}^m \nabla f(x^{k+1}, \xi_i^k)$
- 3: Compute $\tilde{\nabla} f(x^k, \xi^k) = \text{clip}(\nabla f(x^k, \xi^k), \lambda)$ using (4) with $\lambda = B/\gamma$
- 4: $x^{k+1} = x^k - \gamma \tilde{\nabla} f(x^k, \xi^k)$
- 5: **end for**

Output: $\bar{x}^N = \frac{1}{N} \sum_{k=0}^{N-1} x^k$

209 We emphasize that as for clipped-SSTM we use clipping level λ inversely proportional to the stepsize
 210 γ . Below we provide a simplified version of the main result for clipped-SGD in the convex case. The
 211 complete formulation and the full proof of the theorem are deferred to Appendix C.1 (see Thm. C.1).

212 **Theorem 3.1.** *Assume that function f is convex and its gradient satisfy (3) with $\nu \in [0, 1]$, $M_\nu > 0$
 213 on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist such a
 214 choice of parameters that clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$
 215 after*

$$O\left(\max\left\{D^2, D^{1+\nu} \ln \frac{D^2 + D^{1+\nu}}{\beta}\right\}\right), \quad D = \frac{M_\nu^{\frac{1}{1+\nu}} R_0}{\varepsilon^{\frac{1}{1+\nu}}} \quad (8)$$

216 *iterations and requires*

$$O\left(\max\left\{D^2, \max\left\{D^{1+\nu}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\} \ln \frac{D^2 + D^{1+\nu}}{\beta}\right\}\right) \text{ oracle calls.} \quad (9)$$

217 As all our results in the paper, this result for clipped-SGD has two important features: 1) the
 218 dependence on the confidence level β is logarithmic and 2) Hölder continuity is required only on
 219 the ball B_{2R_0} centered at the solution. Moreover, up to the difference in the expressions under
 220 the logarithm the dependence on ε in the result for clipped-SGD is the same as in the tightest
 221 known results for non-accelerated SGD-type methods [4, 16]. Finally, we emphasize that for $\nu < 1$
 222 the logarithmic factors appearing in the complexity bound for clipped-SSTM are worse than the
 223 corresponding factor in the complexity bound for clipped-SGD. Therefore, clipped-SGD has the
 224 best known high-probability complexity results in the case when $\nu = 0$ and f is convex.

225 For the strongly convex problems, we consider restarted version of Alg. 3 (R-clipped-SGD, see
 Alg. 4) and derive high-probability complexity result for this version. Below we provide a simplified

Algorithm 4 Restarted clipped-SGD (R-clipped-SGD): case $\nu \in [0, 1]$

Input: starting point x^0 , number of restarts τ , number of steps of clipped-SGD in restarts $\{N_t\}_{t=1}^\tau$,
 batchsizes $\{m_t\}_{k=1}^\tau$, stepsizes $\{\gamma_t\}_{t=1}^\tau$, clipping parameters $\{B_t\}_{t=1}^\tau$
 1: $\hat{x}^0 = x^0$
 2: **for** $t = 1, \dots, \tau$ **do**
 3: Run clipped-SGD (Alg. 3) for N_t iterations with batchsize m_t , stepsize γ_t , clipping parameter
 B_t , and starting point \hat{x}^{t-1} . Define the output of clipped-SGD by \hat{x}^t .
 4: **end for**
Output: \hat{x}^τ

226 version of the result. The complete formulation and the full proof of the theorem are deferred to
 227 Appendix C.2 (see Thm. C.2).

229 **Theorem 3.2.** *Assume that function f is μ -strongly convex and its gradient satisfy (3) with $\nu \in [0, 1]$,
 230 $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist
 231 such a choice of parameters that R-clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at
 232 least $1 - \beta$ after*

$$\mathcal{O} \left(\max \left\{ D_1^{\frac{2}{1+\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max \left\{ D_1 \ln \frac{\mu R_0^2}{\varepsilon}, D_2 \right\} \ln \frac{D}{\beta} \right\} \right)$$

233 iterations of Alg. 3 in total and requires

$$\mathcal{O} \left(\max \left\{ D_1^{\frac{2}{1+\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max \left\{ D_1 \ln \frac{\mu R_0^2}{\varepsilon}, D_2, \frac{\sigma^2}{\mu \varepsilon} \right\} \ln \frac{D}{\beta} \right\} \right) \text{ oracle calls, where}$$

$$234 \quad D_1 = \frac{M_\nu}{\mu R_0^{1-\nu}}, \quad D_2 = \frac{M_\nu}{\mu^{\frac{1+\nu}{2}} \varepsilon^{\frac{1-\nu}{2}}}, \quad D = (D_1^{\frac{2}{1+\nu}} + D_1) \ln \frac{\mu R_0^2}{\varepsilon} + D_2 + D_2^{\frac{2}{1+\nu}}.$$

235 As in the convex case, for $\nu < 1$ the log factors appearing in the complexity bound for R-clipped-
 236 SSTM are worse than the corresponding factor in the bound for R-clipped-SGD. Thus, R-clipped-
 237 SGD has the best known high-probability complexity results for strongly convex f and $\nu = 0$.

238 4 Numerical experiments

239 We tested the performance of the methods on the following problems:

- 240 • BERT fine-tuning on CoLA dataset [38]. We use pretrained BERT from Transformers library [39]
 241 (bert-base-uncased) and freeze all layers except the last two linear ones.
- 242 • ResNet-18 training on ImageNet-100 (first 100 classes of ImageNet [33]).

243 First, we study the noise distribution for both problem as follows: at the starting point we sample
 244 large enough number of batched stochastic gradients $\nabla f(x^0, \xi_1), \dots, \nabla f(x^0, \xi_K)$ with batchsize
 245 32 and plot the histograms for $\|\nabla f(x^0, \xi_1) - \nabla f(x^0)\|_2, \dots, \|\nabla f(x^0, \xi_K) - \nabla f(x^0)\|_2$, see Fig. 1.
 246 As one can see, the noise distribution for BERT + CoLA is substantially non-sub-Gaussian, whereas
 247 the distribution for ResNet-18 + Imagenet-100 is almost Gaussian.

248 Next, we compared 4 different optimizers on these problems: Adam, SGD (with Momentum),
 249 clipped-SGD (with Momentum and coordinate-wise clipping) and clipped-SSTM (with norm-
 250 clipping and $\nu = 1$). The results are presented in Fig. 2. We observed that the noise distributions do

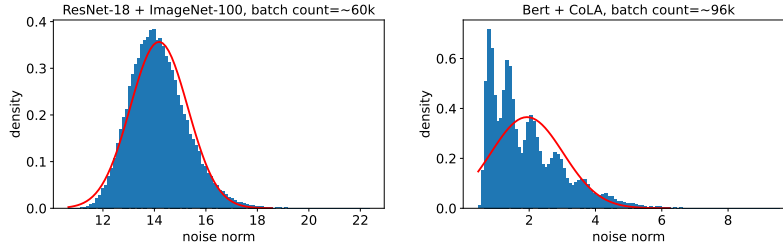


Figure 1: Noise distribution of the stochastic gradients for ResNet-18 on ImageNet-100 and BERT fine-tuning on the CoLA dataset before the training. Red lines: probability density functions with means and variances empirically estimated by the samples. Batch count is the total number of samples used to build a histogram.

251 not change significantly along the trajectories of the considered methods, see Appendix D. During
 252 the hyper-parameters search we compared different batchsizes, emulated via gradient accumulation
 253 (thus we compare methods with different batchsizes by the number of base batches used). The base
 254 batchsize was 32 for both problems, stepsizes and clipping levels were tuned. One can find additional
 255 details regarding our experiments in Appendix D.

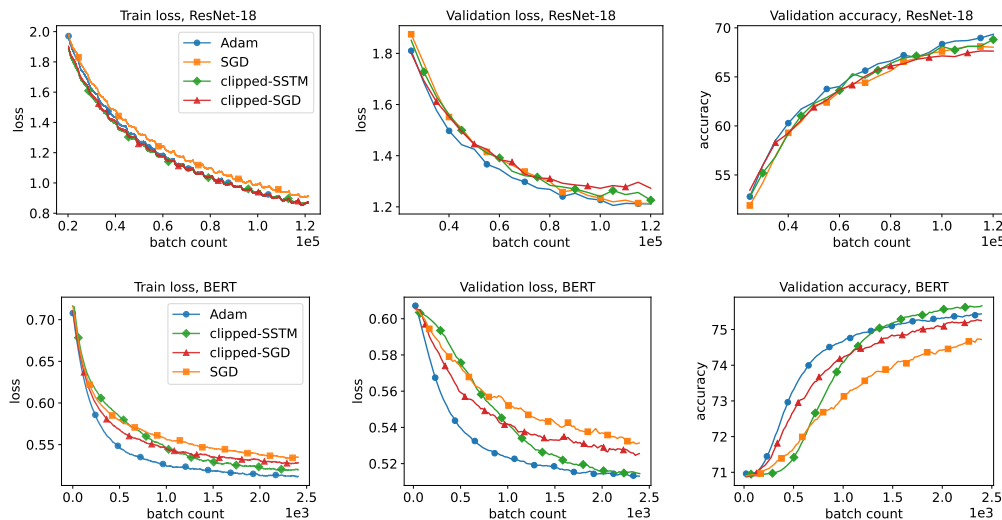


Figure 2: Train and validation loss + accuracy for different optimizers on both problems. Here, “batch count” denotes the total number of used stochastic gradients.

256 **Image classification.** On ResNet-18 + ImageNet-100 task, SGD performs relatively well, and
 257 even ties with Adam (with batchsize of 4×32) in validation loss. clipped-SSTM (with batchsize of
 258 2×32) also ties with Adam and clipped-SGD is not far from them. The results were averaged from
 259 5 different launches (with different starting points/weight initializations). Since the noise distribution
 260 is almost Gaussian even vanilla SGD performs well, i.e., gradient clipping is not required. At the
 261 same time, the clipping does not slow down the convergence significantly.

262 **Text classification.** On BERT + CoLA task, when the noise distribution is heavy-tailed, the methods
 263 with clipping outperform SGD by a large margin. This result is in good correspondence with the
 264 derived high-probability complexity bounds for clipped-SGD, clipped-SSTM and the best-known
 265 ones for SGD. Moreover, clipped-SSTM (with batchsize of 8×32) achieves the same loss on
 266 validation as Adam, and has better accuracy. These results were averaged from 5 different train-val
 267 splits and 20 launches (with different starting points/weight initializations) for each of the splits, 100
 268 launches in total.

269 **References**

- 270 [1] George Bennett. Probability inequalities for the sum of independent random variables. *Journal*
271 *of the American Statistical Association*, 57(297):33–45, 1962.
- 272 [2] Caroline Chaux, Patrick L Combettes, Jean-Christophe Pesquet, and Valérie R Wajs. A varia-
273 tional formulation for frame-based inverse problems. *Inverse Problems*, 23(4):1495–1518, jun
274 2007.
- 275 [3] Damek Davis, Dmitriy Drusvyatskiy, Lin Xiao, and Junyu Zhang. From low probability to
276 high confidence in stochastic convex optimization. *Journal of Machine Learning Research*,
277 22(49):1–38, 2021.
- 278 [4] Olivier Devolder. *Exactness, inexactness and stochasticity in first-order methods for large-scale*
279 *convex optimization*. PhD thesis, PhD thesis, 2013.
- 280 [5] Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex
281 optimization with inexact oracle. *Mathematical Programming*, 146(1):37–75, 2014.
- 282 [6] Pavel Dvurechensky and Alexander Gasnikov. Stochastic intermediate gradient method for con-
283 vex problems with stochastic inexact oracle. *Journal of Optimization Theory and Applications*,
284 171(1):121–145, 2016.
- 285 [7] Kacha Dzhaparidze and JH Van Zanten. On bernstein-type inequalities for martingales. *Stochas-
286 tic processes and their applications*, 93(1):109–117, 2001.
- 287 [8] David A Freedman et al. On tail probabilities for martingales. *the Annals of Probability*,
288 3(1):100–118, 1975.
- 289 [9] Alexander Vladimirovich Gasnikov, Yu E Nesterov, and Vladimir Grigor’evich Spokoiny.
290 On the efficiency of a randomized mirror descent algorithm in online optimization problems.
291 *Computational Mathematics and Mathematical Physics*, 55(4):580–596, 2015.
- 292 [10] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional
293 sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine*
294 *Learning-Volume 70*, pages 1243–1252. JMLR. org, 2017.
- 295 [11] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly
296 convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal*
297 *on Optimization*, 22(4):1469–1492, 2012.
- 298 [12] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex
299 stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 300 [13] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
301 <http://www.deeplearningbook.org>.
- 302 [14] Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with
303 heavy-tailed noise via accelerated gradient clipping. In H. Larochelle, M. Ranzato, R. Had-
304 sell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,
305 volume 33, pages 15042–15053. Curran Associates, Inc., 2020.
- 306 [15] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter
307 Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine*
308 *Learning*, pages 5200–5209, 2019.
- 309 [16] Vincent Guigues, Anatoli Juditsky, and Arkadi Nemirovski. Non-asymptotic confidence bounds
310 for the optimal value of a stochastic program. *Optimization Methods and Software*, 32(5):1033–
311 1058, 2017.
- 312 [17] Cristóbal Guzmán and Arkadi Nemirovski. On lower complexity bounds for large-scale smooth
313 convex optimization. *Journal of Complexity*, 31(1):1–14, 2015.
- 314 [18] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex
315 optimization. In *Advances in Neural Information Processing Systems*, pages 1594–1602, 2015.

- 316 [19] Anatoli Juditsky, Arkadi Nemirovski, et al. First order methods for nonsmooth convex large-
317 scale optimization, i: general purpose methods. *Optimization for Machine Learning*, pages
318 121–148, 2011.
- 319 [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
320 *arXiv:1412.6980*, 2014.
- 321 [21] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical*
322 *Programming*, 133(1-2):365–397, 2012.
- 323 [22] Vien V Mai and Mikael Johansson. Stability and convergence of stochastic gradient clipping:
324 Beyond lipschitz continuity and smoothness. *arXiv preprint arXiv:2102.06489*, 2021.
- 325 [23] Aditya Krishna Menon, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar. Can gradient
326 clipping mitigate label noise? In *International Conference on Learning Representations*, 2020.
- 327 [24] Eric Moulines and Francis R Bach. Non-asymptotic analysis of stochastic approximation
328 algorithms for machine learning. In *Advances in Neural Information Processing Systems*, pages
329 451–459, 2011.
- 330 [25] Aleksandr Viktorovich Nazin, AS Nemirovsky, Aleksandr Borisovich Tsybakov, and AB Judit-
331 sky. Algorithms of robust stochastic optimization based on mirror descent method. *Automation*
332 *and Remote Control*, 80(9):1607–1627, 2019.
- 333 [26] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic
334 approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–
335 1609, 2009.
- 336 [27] Arkadi Semenovich Nemirovsky and David Borisovich Yudin. Problem complexity and method
337 efficiency in optimization. 1983.
- 338 [28] Yu Nesterov. Universal gradient methods for convex optimization problems. *Mathematical*
339 *Programming*, 152(1-2):381–404, 2015.
- 340 [29] Yurii E Nesterov. A method for solving the convex programming problem with convergence
341 rate $o(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- 342 [30] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent
343 neural networks. In *International conference on machine learning*, pages 1310–1318, 2013.
- 344 [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
345 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
346 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
347 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
348 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
349 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,
350 pages 8024–8035. Curran Associates, Inc., 2019.
- 351 [32] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of*
352 *mathematical statistics*, pages 400–407, 1951.
- 353 [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
354 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
355 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
356 (*IJCV*), 115(3):211–252, 2015.
- 357 [34] Umut Şimşekli, Mert Gürbüzbalaban, Thanh Huy Nguyen, Gaël Richard, and Levent Sagun. On
358 the heavy-tailed theory of stochastic gradient descent for deep neural networks. *arXiv preprint*
359 *arXiv:1912.00018*, 2019.
- 360 [35] Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic
361 gradient noise in deep neural networks. *arXiv preprint arXiv:1901.06053*, 2019.
- 362 [36] Vladimir Spokoiny et al. Parametric estimation. finite sample theory. *The Annals of Statistics*,
363 40(6):2877–2909, 2012.

- 364 [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 365 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
 366 *processing systems*, pages 5998–6008, 2017.
- 367 [38] Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability
 368 judgments. *arXiv preprint arXiv:1805.12471*, 2018.
- 369 [39] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
 370 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,
 371 Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain
 372 Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-
 373 art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods*
 374 *in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.
 375 Association for Computational Linguistics.
- 376 [40] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
 377 training: A theoretical justification for adaptivity. In *International Conference on Learning*
 378 *Representations*, 2020.
- 379 [41] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi,
 380 Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? In
 381 H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural*
 382 *Information Processing Systems*, volume 33, pages 15383–15393. Curran Associates, Inc.,
 383 2020.

384 Checklist

- 385 1. For all authors...
- 386 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 387 contributions and scope? [Yes]
- 388 (b) Did you describe the limitations of your work? [Yes] Section 1.1 describes all assump-
 389 tions that we use
- 390 (c) Did you discuss any potential negative societal impacts of your work? [No] Our results
 391 are primarily theoretical, therefore, such a discussion is not applicable.
- 392 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 393 them? [Yes]
- 394 2. If you are including theoretical results...
- 395 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Section 1.1
 396 describes all assumptions that we use.
- 397 (b) Did you include complete proofs of all theoretical results? [Yes] Appendix B and C
 398 include the complete proofs of all the results we derive.
- 399 3. If you ran experiments...
- 400 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 401 mental results (either in the supplemental material or as a URL)? [Yes] See our code in
 402 the supplementary material.
- 403 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 404 were chosen)? [Yes] See Appendix D.
- 405 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
 406 iments multiple times)? [No] Instead of it, we show the averaged trajectories of the
 407 methods’ convergence.
- 408 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 409 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.
- 410 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 411 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 412 (b) Did you mention the license of the assets? [No] We use only publicly available
 413 resources.

- 414 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
415 (d) Did you discuss whether and how consent was obtained from people whose data you're
416 using/curating? [No] We use only publicly available resources.
417 (e) Did you discuss whether the data you are using/curating contains personally identifiable
418 information or offensive content? [N/A]
- 419 5. If you used crowdsourcing or conducted research with human subjects...
- 420 (a) Did you include the full text of instructions given to participants and screenshots, if
421 applicable? [N/A]
- 422 (b) Did you describe any potential participant risks, with links to Institutional Review
423 Board (IRB) approvals, if applicable? [N/A]
- 424 (c) Did you include the estimated hourly wage paid to participants and the total amount
425 spent on participant compensation? [N/A]