

# Clipped Methods for Stochastic Optimization with Heavy-Tailed Noise

TES Conference on Mathematical Optimization for  
Machine Learning, Berlin

---

**Eduard Gorbunov**

September 15, 2023

Mohamed bin Zayed University of Artificial Intelligence

1. Clipping and Heavy-Tailed Noise
2. In-Expectation Guarantees vs High-Probability Convergence
3. Main Results

# The Talk is Based on Four Papers

- Gorbunov, E., Danilova, M., & Gasnikov, A. (2020).  
 . NeurIPS 2020
- Gorbunov, E., Danilova, M., Shibaev, I., Dvurechensky, P., & Gasnikov, A. (2021). <sup>1</sup>  
arXiv:2106.05958
- Gorbunov, E., Danilova, M., Dobre, D., Dvurechenskii, P., Gasnikov, A., & Gidel, G. (2022). ;  
 . NeurIPS 2022.
- Sadiev, A., Danilova, M., Gorbunov, E., Horváth, S., Gidel, G., Dvurechensky, P., Gasnikov, A., & Richtárik, P. (2023). }  
 . ICML 2023.

# Clipping and Heavy-Tailed Noise

---

$$\theta^{+1} = \theta - \eta r(\theta; \mathcal{D}) \quad (1)$$

- $F$  – the function to be minimized
- $r(\theta; \mathcal{D})$  – stochastic gradient, i.e.,  $\mathbb{E}[r(\theta; \mathcal{D})] = \nabla F(\theta)$  estimate of  $\nabla F(\theta)$ :

# Clipped Stochastic Gradient Descent ( $V_{\text{c}}$ XWF: 7)

$$V_{\text{c}}^{+1} = V_{\text{c}}(r(\cdot; \cdot)); \quad (2)$$

- $V_{\text{c}}(\cdot; \cdot) = \min_{\|g\| \leq c} \langle g, \cdot \rangle$
- $V_{\text{c}}(r(\cdot; \cdot); \cdot)$  - estimate of  $r(\cdot)$ :  
 $\mathbb{E}[V_{\text{c}}(r(\cdot; \cdot); \cdot)] \approx r(\cdot)$

# Origin of Clipping

- Gradient clipping was proposed in (Pascanu et al., 2013). Originally it was used to handle exploding and vanishing gradients in RNNs.

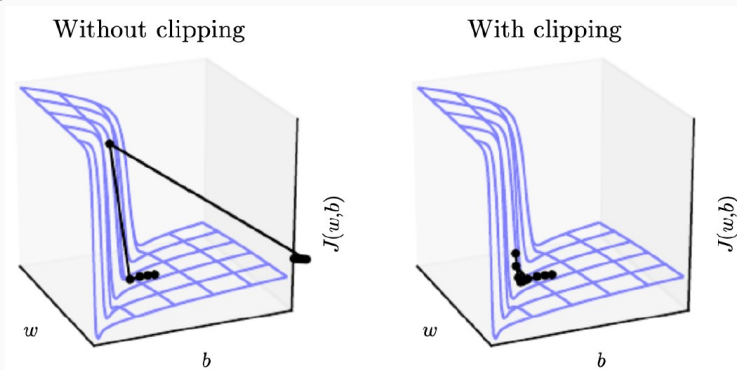


Figure 1: from (Goodfellow et al., 2016)

## Few Years Later in NLP...

- Merity et al. (2017) use gradient clipping for LSTM
- Peters et al. (2017) trained their deep bidirectional language model with  $\text{4W}^{\sim}$  + clipping
- Mosbach et al. (2020) fine-tune BERT using  $\text{4W}^{\sim} J$  + clipping



## Few Years Later in NLP...

- Merity et al. (2017) use gradient clipping for LSTM
- Peters et al. (2017) trained their deep bidirectional language model with  $\text{clip}(\nabla_{\theta} J)$  + clipping
- Mosbach et al. (2020) fine-tune BERT using  $\text{clip}(\nabla_{\theta} J)$  + clipping

It seems that gradient clipping is an important component in training these models. Why?

# Heavy-Tailed Noise in Stochastic Gradients

Let us look at the distribution of  $\|g_k\|_2$  in two settings:

- Standard CV task: training ResNet50 on ImageNet dataset
- Standard NLP task: training BERT on Wikipedia+Books dataset

# Heavy-Tailed Noise in Stochastic Gradients

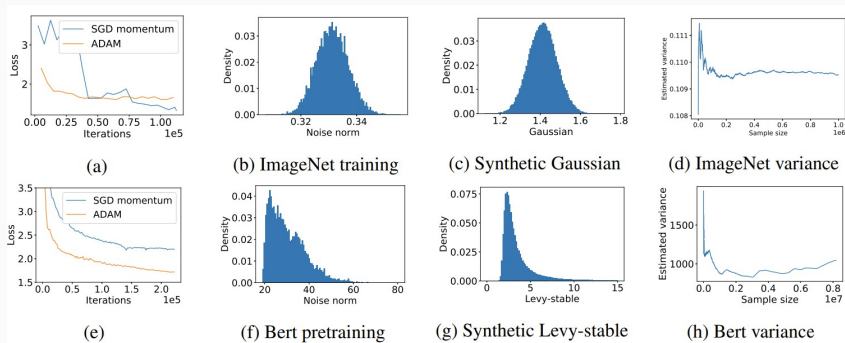


Figure 2: from (Zhang et al., 2020)

We see that  $474@$  is much better than  $F: 7$  when the noise in the stochastic gradient is heavy-tailed

# $4W^*$ and $V_{cc}XZF: 7$

- $V_{cc}XZF: 7$ :

$$+1 = V_{cc} r ( ; );$$

- $4W^*$ :

$$= 1 \quad 1 + (1 \quad 1)r ( ; );$$

$$= 2 \quad 1 + (1 \quad 2)(r ( ; ))^2;$$

$$+1 = \frac{P}{+}$$

- When  $\beta_1 = 0$   $4W^*$  ( $E@Fcebc$ ) can be seen as  $V_{cc}XZF: 7$  with “adaptive”  $\lambda$

# Definition of Heavy-Tailed Noise in Stochastic Gradients

- Random vector  $g$  has light tails if

$$P\|g\| \geq k \leq \frac{E\|g\|^2}{2} \exp\left(-\frac{k^2}{2}\right) \quad \delta > 0: \quad (3)$$

The above condition is equivalent (up to the numerical factor in  $\delta$ ) to

$$E \exp\left(-\frac{k}{2} \frac{E\|g\|^2}{2}\right) \leq \exp(-1): \quad (4)$$

# Definition of Heavy-Tailed Noise in Stochastic Gradients

- Random vector  $g$  has light tails if

$$P\|g\| \geq \delta \leq \frac{E\|g\|^2}{2} \exp\left(-\frac{\delta^2}{2}\right) \quad \delta > 0: \quad (3)$$

The above condition is equivalent (up to the numerical factor in (3)) to

$$E \exp\left(\frac{\|g\|^2}{2}\right) \leq \exp(1): \quad (4)$$

- Otherwise we say that  $g$  has heavy tails. However, in this talk, we will assume that it has bounded central  $2$ -th moment for some  $\sigma^2$ :

$$E\|g\|^2 \leq \sigma^2 \quad (5)$$

# In-Expectation Guarantees vs High-Probability Convergence

---

## Problem and Assumptions

$$\min_{z \in \mathbb{R}} f(z) = \mathbb{E} [l(z; \mathcal{D})] \quad (6)$$

- $l: \mathbb{R} \rightarrow \mathbb{R}$  is convex and  $\mathcal{C}^2$ -smooth, i.e.,  $\exists L, \mu > 0$ ;

$$l(x) \leq l(x_0) + h(x - x_0); \quad \forall x, x_0 \in \mathbb{R}; \quad (7)$$

$$|l''(x)| \leq L; \quad \forall x \in \mathbb{R}; \quad (8)$$



## Problem and Assumptions

$$\min_{\mathcal{Z} \subseteq \mathbb{R}^d} f(\theta) = \mathbb{E} [l(\theta; \mathcal{Z})] g \quad (6)$$

- $\mathcal{Z} \subseteq \mathbb{R}^d$  is convex and  $\mathcal{Z}$ -smooth, i.e.,  $\delta \in \mathcal{Z} \subseteq \mathbb{R}^d$

$$f(\theta + h) \leq f(\theta) + h^T r(\theta); \quad \theta \in \mathcal{Z} \quad (7)$$

$$\|kr(\theta) - r(\theta)k \leq \mathcal{K} \|k\| \quad \theta \in \mathcal{Z} \quad (8)$$

- Stochastic gradient  $r(\theta; \mathcal{Z})$  with bounded central  $\ell_2$ -th moment ( $\mathcal{Z} \subseteq \mathbb{R}^d$ ) is available, i.e.,  $\delta \in \mathcal{Z} \subseteq \mathbb{R}^d$

$$\mathbb{E} [r(\theta; \mathcal{Z})] = r(\theta); \quad \mathbb{E} [kr(\theta; \mathcal{Z}) - r(\theta)k^2] \leq \mathcal{K} \quad \theta \in \mathcal{Z} \quad (9)$$

## F: 7 Does Not Converge When $\alpha < 2$

- In-expectation guarantees:  $E[k^2] = \alpha^2 E[k]$ ,  $E[(k-1)^2] = \alpha^2 E[k] - 2\alpha + 1$ ;

## F: 7 Does Not Converge When $\alpha < 2$

- In-expectation guarantees:  $E[k^2] < \infty$ ,  $E[r(\cdot)] < \infty$ ,  
 $E[kr(\cdot)] < \infty$
- Consider the example from (Zhang et al., 2020):  $r(\cdot) = \frac{1}{2}k^2$  and  
 $r(\cdot; \alpha) = \frac{1}{2}k^2 + \epsilon$ , where  $E[\epsilon] = 0$  and  $E[k^2] < \infty$  but  $E[k^2] = \infty$   
(e.g.,  $\epsilon$  can Lévy  $\alpha$ -stable distribution)

## F: 7 Does Not Converge When $\alpha < 2$

- In-expectation guarantees:  $E[k^2] < \infty$ ,  $E[(\cdot)] < \infty$ ,  $E[kr(\cdot)k^2] < \infty$
- Consider the example from (Zhang et al., 2020):  $(\cdot) = \frac{1}{2}k^2$  and  $r(\cdot; \cdot) = \cdot + \cdot$ , where  $E[\cdot] = 0$  and  $E[k^2] < \infty$  but  $E[k^2] = 1$  (e.g., can Lévy  $\alpha$ -stable distribution)
- Then, after one step of **F: 7** we have

$$\begin{aligned}
 E[k^1 k^2] &= E[k^0 r(\cdot; \cdot)k^2] \\
 &= \underbrace{E[k^0 k^2]}_{\text{infinite}} + \underbrace{2 E[kr(\cdot; \cdot)k^2]}_{=1} \\
 &= \infty
 \end{aligned}$$

The method does not converge in expectation (in  $\mathbb{Q}$ ) when  $\alpha < 2$ !  
 What about the case when  $\alpha = 2$  (bounded variance)?

# In-Expectation Guarantees and Trajectories of the Method

Consider  $F: \mathbb{R}^n$  with constant stepsize

$$x^{i+1} = x^i - \alpha r(x^i; \omega^i)$$

applied to a toy stochastic quadratic problem:

$$\min_{x \in \mathbb{R}^n} f(x) = \mathbb{E}[r(x; \omega)]g; \quad r(x; \omega) = \frac{1}{2}k \|x\|^2 + h; \quad \omega^i;$$

where  $\mathbb{E}[\omega] = 0$  and  $\mathbb{E}[k \|\omega\|^2] = \sigma^2$ .

# In-Expectation Guarantees and Trajectories of the Method

Consider  $F: \mathcal{Z}$  with constant stepsize

$$z^{+1} = z - \eta r(z; \theta)$$

applied to a toy stochastic quadratic problem:

$$\min_{z \in \mathbb{R}} f(z) = \mathbb{E}[r(z; \theta)]g; \quad r(z; \theta) = \frac{1}{2}k z^2 + h; \quad \theta = i;$$

where  $\mathbb{E}[\theta] = 0$  and  $\mathbb{E}[k z^2] = \sigma^2$ . We consider three scenarios:

- $\theta$  has Gaussian distribution
- $\theta$  has Weibull distribution (non-sub-Gaussian)
- $\theta$  has Burr Type XII distribution (non-sub-Gaussian)

## In-Expectation Guarantees and Trajectories of the Method

For all of three cases, state-of-the-art theory on  $F$ :  $\mathcal{F}$  (Ghadimi and Lan, 2013) says

$$\mathbb{E} \left( \sum_{i=0}^{h-1} \left\| \nabla F(x_i) \right\|^2 \right) \leq \frac{2}{\eta} \left( F(x_0) - F^* \right) + \frac{\eta}{2} \sum_{i=0}^{h-1} \left\| \nabla^2 F(x_i) \right\| \left( F(x_i) - F^* \right) \quad (10)$$

# In-Expectation Guarantees and Trajectories of the Method

For all of three cases, state-of-the-art theory on  $F$ :  $\nabla$  (Ghadimi and Lan, 2013) says

$$\mathbb{E} \left( \sum_{i=0}^{h-1} \left\| \nabla F(x_i) \right\|^2 \right) \leq \frac{2}{\epsilon} \left( \sum_{i=0}^{h-1} \left\| \nabla F(x_i) \right\|^2 \right) + \frac{2}{\epsilon} \left( \sum_{i=0}^{h-1} \left\| \nabla F(x_i) \right\|^2 \right) \quad (10)$$

However, the behavior in practice does depend on the distribution:

**Figure 3:** from (Gorbunov et al., 2020)



# In-Expectation Guarantees vs High-Probability Guarantees

- In-expectation guarantees:  $E[k^2] \leq \sigma^2$ ,  $E[\text{error}] \leq \sigma^2$ ;  $E[kr^2] \leq \sigma^2$ 
  - Typically, depend only on some moments of stochastic gradient, e.g., variance

# In-Expectation Guarantees vs High-Probability Guarantees

- In-expectation guarantees:  $E[k^2] \leq \sigma^2$ ,  $E[f(\theta) - f(\theta^*)] \leq \sigma^2$ ;  $E[kr(\theta)k^2] \leq \sigma^2$ 
  - Typically, depend only on some moments of stochastic gradient, e.g., variance
- High-probability guarantees:  $P\{k^2 \leq \sigma^2/g\} \geq 1 - \delta$ ,  $P\{f(\theta) - f(\theta^*) \leq \sigma^2/g\} \geq 1 - \delta$ ,  $P\{kr(\theta)k^2 \leq \sigma^2/g\} \geq 1 - \delta$ 
  - Sensitive to the distribution of the stochastic gradient noise

# High-Probability Convergence of $F$ : 7 under Bounded Variance Assumption

Natural idea: apply Markov's inequality:

$$P(f^{(\wedge)}(\cdot) > \epsilon) < \frac{E[f^{(\wedge)}(\cdot)]}{\epsilon}.$$

# High-Probability Convergence of $F: \mathcal{Z}$ under Bounded Variance Assumption

Natural idea: apply Markov's inequality:

$$P(f(\hat{\theta}) > \epsilon) < \frac{E[f(\hat{\theta})]}{\epsilon}.$$

Taking  $\epsilon = \epsilon$  of  $F: \mathcal{Z}$ , we can guarantee  $E[f(\hat{\theta})] < \epsilon$  that implies  $P(f(\hat{\theta}) > \epsilon) < \frac{\epsilon}{\epsilon}$  or, equivalently,  $P(f(\hat{\theta}) > \epsilon) < 1$ .

# High-Probability Convergence of $F: \mathcal{Z}$ under Bounded Variance Assumption

Natural idea: apply Markov's inequality:

$$P[f(\hat{\theta}) - g < -\epsilon] \leq \frac{E[f(\hat{\theta}) - g]}{\epsilon}.$$

Taking  $\epsilon = \frac{\sigma_0^2}{2n}$  of  $F: \mathcal{Z}$ , we can guarantee  $E[f(\hat{\theta}) - g] \leq \frac{\sigma_0^2}{2n}$  that implies  $P[f(\hat{\theta}) - g < -\frac{\sigma_0^2}{2n}] \leq \frac{\sigma_0^2}{2n \cdot \frac{\sigma_0^2}{2n}} = \frac{1}{2}$ .  
 or, equivalently,  $P[f(\hat{\theta}) - g < -\frac{\sigma_0^2}{2n}] \leq \frac{1}{2}$ .

**Bad news:** to ensure  $E[f(\hat{\theta}) - g] \leq \frac{\sigma_0^2}{2n}$   $F: \mathcal{Z}$  needs

$$O \left( \max \left\{ \frac{\sigma_0^2}{n}, \frac{\sigma_0^2}{n^2} \right\} \right) \text{ oracle calls}$$

Negative-power dependence on  $n$ :

# High-Probability Convergence of $F: \mathcal{Z}$ under Bounded Variance Assumption

Natural idea: apply Markov's inequality:

$$P[f(\hat{\theta}) - g < -\epsilon] \leq \frac{E[f(\hat{\theta}) - g]}{\epsilon}.$$

Taking  $\epsilon = \epsilon_0$  of  $F: \mathcal{Z}$ , we can guarantee  $E[f(\hat{\theta}) - g] \leq \epsilon_0$  that implies  $P[f(\hat{\theta}) - g < -\epsilon_0] \leq \frac{\epsilon_0}{\epsilon_0} = 1$ . or, equivalently,  $P[f(\hat{\theta}) - g < -\epsilon_0] \leq \frac{\epsilon_0}{\epsilon_0} = 1$ .

Bad news: to ensure  $E[f(\hat{\theta}) - g] \leq \epsilon_0$   $F: \mathcal{Z}$  needs

$$O \left( \max \left\{ \frac{\sigma_0^2}{\epsilon_0^2}, \frac{\sigma_0^2}{\epsilon_0^2} \right\} \right) \text{ oracle calls}$$

Negative-power dependence on  $\epsilon_0$ :

**Natural question:** can we analyze high-probability convergence of  $F: \mathcal{Z}$  better?

# High-Probability Convergence of $F: \mathbb{R}^n$ under Bounded Variance Assumption

## Failure of $F: \mathbb{R}^n$

For any  $\epsilon > 0$ ,  $\delta \in (0, 1)$ , and  $F: \mathbb{R}^n$  parameterized by the number of steps  $N$  and stepsize  $\eta$ , there exists  $\epsilon$ -strongly convex  $L$ -smooth problem and stochastic oracle with noise having bounded  $L$ -th moment with  $L = 2$ ,  $0 < \epsilon < \frac{1}{2L}$  such that for the iterates produced by  $F: \mathbb{R}^n$  with any stepsize  $0 < \eta < \frac{1}{2L}$

$$\mathbb{P} \left( \|x_k - x^*\| \leq \frac{\epsilon}{k^2} \right) \geq 1 - \delta \quad \Rightarrow \quad \mathbb{P} \left( \sum_{k=1}^N \|x_k - x^*\|^2 \leq \frac{\epsilon^2}{\delta} \right) \geq 1 - \delta \quad (11)$$

This illustrates the necessity of modifying the method, e.g., one can use gradient clipping

## Main Results

---



# Key Challenge in the Analysis of $V_{cc} \setminus \text{CCXVF}$ : 7

$$+1 = \frac{V_{cc} \setminus r(\cdot; \cdot)}{\tilde{r}(\cdot; \xi) \{Z\}}$$

- Key challenge:  $E^h(\cdot; \cdot) \neq E^i(\cdot; \cdot)$

# Analysis of $V_{cc} \setminus CCXVF$ : 7: Key Idea

- We start the proof classically:

$$\begin{aligned} k^{+1} \quad k^2 &= k \quad k^2 \quad 2 \quad h \quad ; \mathcal{P} ( ; ) i \\ &+ {}^2 k \mathcal{P} ( ; ) k^2 \\ &\vdots \end{aligned}$$

# Analysis of $V_{cc} \setminus \text{CCXVF}$ : 7: Key Idea

- We start the proof classically:

$$\begin{aligned}
 k^{+1} \quad k^2 &= k \quad k^2 \quad 2 \quad h \quad ; \mathcal{F} ( ; ) i \\
 &+ \quad k^2 \mathcal{F} ( ; ) k^2 \\
 &\vdots
 \end{aligned}$$

- Using convexity and smoothness of  $\mathcal{F}$  and simple rearrangements, we eventually get for  $\mathcal{F} = ( ; ) \quad ( ; )$ ,
 
$$\mathcal{F} = k \quad k, \quad \mathcal{F} = \mathcal{F} ( ; ) \quad r ( ; )$$

$$\begin{aligned}
 \frac{2 (1 \quad 2 \quad \mathcal{Q}) \mathcal{X}^1}{1} &= 0 \quad \frac{1}{1} \quad \mathcal{F}_0^2 \quad \mathcal{F}_i^2 \\
 &+ \frac{2}{1} \quad \mathcal{X}^1 \quad h \quad ; \quad i + \frac{2}{1} \quad \mathcal{X}^1 \quad k \quad k^2 \\
 &= 0 \quad = 0
 \end{aligned}$$

How to upper bound the sums in red?

# Bernstein Inequality for Martingale Differences

**Lemma 1 (Bennett, 1962; Dzhaparidze and Van Zanten, 2001; Freedman et al., 1975)**

Let the sequence of random variables  $f_j$  form a martingale difference sequence, i.e.  $\mathbb{E}[f_j | \mathcal{F}_{j-1}] = 0$  for all  $j \geq 1$ . Assume that conditional variances  $\sigma_j^2 \stackrel{\text{def}}{=} \mathbb{E}[f_j^2 | \mathcal{F}_{j-1}]$  exist and are bounded and assume also that there exists deterministic constant  $\beta > 0$  such that  $\sigma_j \geq \beta$  almost surely for all  $j \geq 1$ .

# Bernstein Inequality for Martingale Differences

**Lemma 1 (Bennett, 1962; Dzhaparidze and Van Zanten, 2001; Freedman et al., 1975)**

Let the sequence of random variables  $f_j$  form a martingale difference sequence, i.e.  $E[f_j | \mathcal{F}_{j-1}] = 0$  for all  $j \geq 1$ . Assume that conditional variances  $\sigma_j^2 \stackrel{\text{def}}{=} E[f_j^2 | \mathcal{F}_{j-1}]$  exist and are bounded and assume also that there exists deterministic constant  $r > 0$  such that  $\sigma_j \leq r$  almost surely for all  $j \geq 1$ . Then for all  $\epsilon > 0, r > 0$  and  $t \geq 1$

$$P\left(\sum_{j=1}^t X_j > \epsilon \text{ and } \sum_{j=1}^t \sigma_j^2 \leq r t\right) \leq 2 \exp\left(-\frac{\epsilon^2}{2r + \epsilon}\right)$$

# Bernstein Inequality for Martingale Differences

**Lemma 1 (Bennett, 1962; Dzhaparidze and Van Zanten, 2001; Freedman et al., 1975)**

Let the sequence of random variables  $f_j, g_j$  form a martingale difference sequence, i.e.  $\mathbb{E}[f_j | \mathcal{F}_{j-1}] = 0$  for all  $j \geq 1$ . Assume that conditional variances  $\sigma_j^2 \stackrel{\text{def}}{=} \mathbb{E}[f_j^2 | \mathcal{F}_{j-1}]$  exist and are bounded and assume also that there exists deterministic constant  $r > 0$  such that  $\sigma_j \leq r$  almost surely for all  $j \geq 1$ . Then for all  $\epsilon > 0, r > 0$  and  $t \geq 1$

$$\mathbb{P}\left(\sum_{j=1}^t X_j > \epsilon\right) \leq \exp\left(-\frac{\epsilon^2}{2r + \frac{\epsilon^2}{3}}\right)$$

To bound  $\sum_{j=0}^t h_j$ ;  $\sum_{j=0}^t k_j$  we need to

- upper bound bias, variance, and distortion of
- have high-prob. upper bounds for  $k_j$  and  $k_j$

## Lemma 2

Let  $\mathbf{e}$  be a random vector in  $\mathbb{R}^2$  and  $\mathbf{e} = \mathbf{V} \mathbf{c}(\mathbf{e}; \mathbf{c})$ . Then,  $\|\mathbf{e}\|^2 = \mathbb{E}[\|\mathbf{e}\|^2]$ . Moreover, if for some  $\alpha > 0$  and  $\beta \in (1, 2]$  we have  $\mathbb{E}[\|\mathbf{e}\|^\beta] = \beta \mathbb{R}$ ,  $\mathbb{E}[\|\mathbf{e}\|^\alpha] = \alpha \mathbb{R}$ , and  $\|\mathbf{e}\|^\beta \leq \alpha \|\mathbf{e}\|^\alpha$ , then

$$\mathbb{E}[\|\mathbf{e}\|] \leq \frac{\beta}{\alpha - \beta} \mathbb{R}; \quad (12)$$

$$\mathbb{E}[\|\mathbf{e}\|^2] \leq \frac{\beta}{\alpha - \beta} \mathbb{R}^2; \quad (13)$$

$$\mathbb{E}[\|\mathbf{e}\|^2] \leq \frac{\beta}{\alpha - \beta} \mathbb{E}[\|\mathbf{e}\|^\alpha]; \quad (14)$$

# Bound on the Distance to the Solution

Inequality

$$\frac{2}{1} \frac{(1 - 2^i) \epsilon^i}{2^i} \leq \frac{1}{2^i} \epsilon_0^2 + \epsilon_i^2 + \frac{2}{2^i} \epsilon^i h \quad ; \quad i + \frac{2}{2^i} \epsilon^i k \leq k^2$$

implies

$$\epsilon_i^2 \leq \epsilon_0^2 + 2 \frac{\epsilon^i}{2^i} h \quad ; \quad i + 2 \frac{\epsilon^i}{2^i} k \leq k^2$$



# Bound on the Distance to the Solution

Inequality

$$\frac{2}{1} (1 - 2^i) \leq \frac{2^i}{1} \leq \frac{2^i}{1} + \frac{2^i}{1} h \leq \frac{2^i}{1} + \frac{2^i}{1} k^2$$

implies

$$\frac{2^i}{1} \leq \frac{2^i}{1} + \frac{2^i}{1} h \leq \frac{2^i}{1} + \frac{2^i}{1} k^2$$

Key idea: prove  $\frac{2^i}{1} \leq \frac{2^i}{1} + \frac{2^i}{1} h$  ;  $\frac{2^i}{1} \leq \frac{2^i}{1} + \frac{2^i}{1} k^2$  with high probability for some numerical constant ; using the induction!

## Theorem 1

Let  $f$  be convex and  $g$  smooth on

$\mathbb{R}^k$ . Let  $\mathcal{C} = \{x \in \mathbb{R}^k : g(x) \leq \tau\}$  and (9) holds on  $\mathcal{C}$ .

# High-Probability Convergence of $V_{\text{CC}}^{\text{XZF}}$ : 7

## Theorem 1

Let  $f$  be convex and  $g$  smooth on  $\mathbb{R}^k$  and (9) holds on  $\mathbb{R}^k$ . Then, for all  $\epsilon \in (0; 1)$ ,  $\delta > 0$  such that  $\ln(\frac{1}{\delta}) \leq \frac{1}{\epsilon}$  there exists a choice of  $\eta$  such that  $V_{\text{CC}}^{\text{XZF}}$  with clipping level  $\eta$  and batchsize  $m = 1$  finds  $x^*$  satisfying  $\|x^* - x^*\| \leq \epsilon$  with probability at least  $1 - \delta$  using

# High-Probability Convergence of $V_{\text{CC}}^{\text{XZF}}$ : 7

## Theorem 1

Let  $f$  be convex and  $g$  smooth on  $\mathcal{C}$ . Let  $\mathcal{C} \subseteq \mathbb{R}^n$  and  $g$  and (9) holds on  $\mathcal{C}$ . Then, for all  $\epsilon \in (0, 1)$ ,  $\delta > 0$  such that  $\ln(\frac{1}{\delta}) \leq \frac{1}{\epsilon}$  there exists a choice of  $\eta$  such that  $V_{\text{CC}}^{\text{XZF}}$ : 7 with clipping level  $\eta$  and batchsize  $m = 1$  finds  $x^*$  satisfying  $\|g(x^*)\| \leq \epsilon$  with probability at least  $1 - \delta$  using

$$O \left( \max \left\{ \frac{1}{\epsilon^2}; \frac{1}{\delta} \ln \frac{1}{\delta} \right\} \right)$$

iterations/oracle calls.

In (Gorbunov et al., 2020, 2021, 2022; Sadiev et al., 2023) we also have

- Accelerated method (Clipped Stochastic Similar Triangles Method)
- Results for the non-convex objectives
- Results for the strongly convex objectives
- Results for the functions with Hölder continuous gradient
- Results for the variational inequalities

# Numerical Experiments: Setup

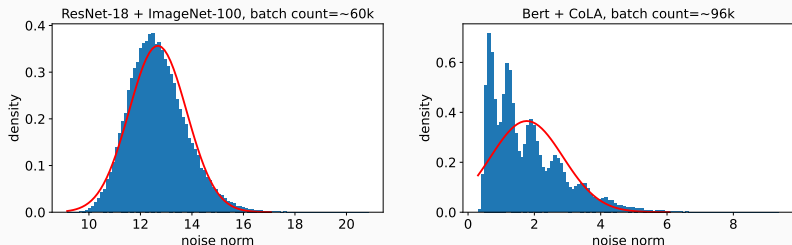
We tested the performance of the methods on the following problems<sup>1</sup>:

- **58EG** ( 0.6 $\mu$  parameters) fine-tuning on **6b?4** dataset. We use pretrained **58EG** and freeze all layers except the last two linear ones. This dataset contains 8551 sentences, and the task is binary classification – to determine if sentence is grammatically correct.
- **EXf AXgž \$+** ( 11.7 $\mu$  parameters) training on **<`TZXAXgž \$##** (first 100 classes of **<`TZXAXg**). It has 134395 images.

---

<sup>1</sup>The code is available at [\[ ggcf- ""Z\g\[ hU! Vb` " 6\\_\ccXWf gbV\[ Tf g\ V@Xg\[ bW" V\\_\ccXWFFG@](https://github.com/robertostromberg/llm-clip)

# Numerical Experiments: Noise Distribution



**Figure 4:** Noise distribution of the stochastic gradients for  $EXAXgžS+$  on  $<`TZXAXgžS##$  and  $58EG$  fine-tuning on the  $6b?4$  dataset before the training. Red lines: probability density functions of normal distributions with means and variances empirically estimated by the samples. Batch count is the total number of samples used to build a histogram.

# Numerical Results, Image Classification

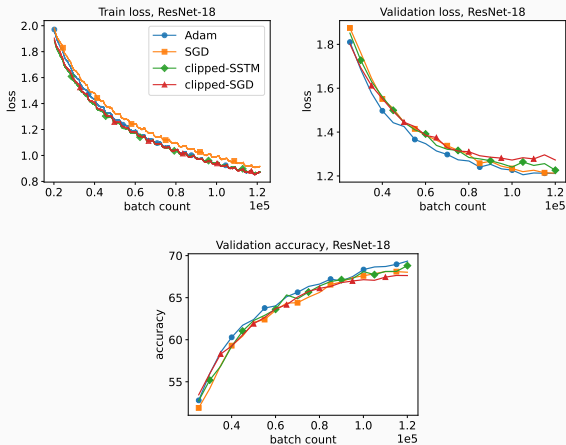


Figure 5: Train and validation loss + accuracy for different optimizers on  $EXFAXgžS+ + <`TZXAXgžS##$  problem. Here, “batch count” denotes the total number of used stochastic gradients. The noise distribution is almost Gaussian, even vanilla  $F: 7$  performs well.



# Numerical Results, Text Classification

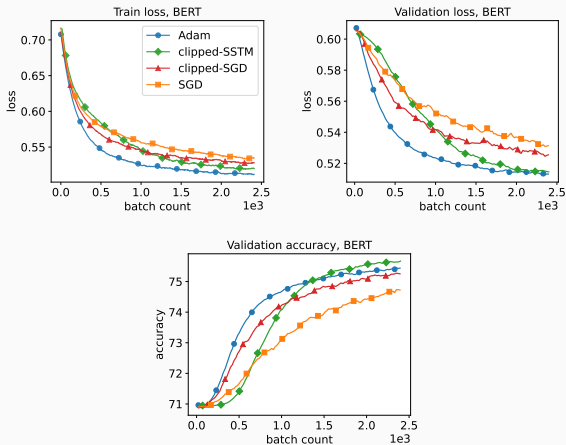


Figure 6: Train and validation loss + accuracy for different optimizers on *58EG*+ *6b?4* problem. The noise distribution is heavy-tailed, the methods with clipping outperform *F: 7* by a large margin.

# Conclusion

- Some popular problems have heavy-tailed noise: in NLP it was observed before, for GANs we demonstrated empirically
- Clipping is a simple way to deal with heavy-tailed noise
- High-probability convergence results for methods with clipping are better than known high-probability convergence results for methods without it
- Partial explanation of the success of adaptive methods like **4W** on GANs and NLP tasks

# About MBZUAI

Established in 2019, located in Masdar City (Abu Dhabi, UAE)

First classes started in January 2021 (because of COVID-19)

Three departments: NLP, CV, and ML

Some numbers: 300 students, 50 faculties, 20th in CSRankings (AI, CV, ML, and NLP)



Figure 7: [ggcf- " "jjj ! TeTUaXj f! Vb` " abWX" S\*% \$\$\$" T` c

## References

---

- Bennett, G. (1962). Probability inequalities for the sum of independent random variables. *J. Amer. Statist. Assoc.*, 57(297):33–45.
- Dzhaparidze, K. and Van Zanten, J. (2001). On Bernstein-type inequalities for martingales. *Stochastics*, 93(1):109–117.
- Freedman, D. A. et al. (1975). On tail probabilities for martingales. *Ann. Probab.*, 3(1):100–118.
- Ghadimi, S. and Lan, G. (2013). Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM J. Optim.*, 23(4):2341–2368.

## References ii

- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *D* © . MIT Press. [ggc- " "jjj! Wxc\_XTea\ aZUbb^! beZ
- Gorbunov, E., Danilova, M., Dobre, D., Dvurechensky, P., Gasnikov, A., and Gidel, G. (2022). Clipped stochastic methods for variational inequalities with heavy-tailed noise.
- Gorbunov, E., Danilova, M., and Gasnikov, A. (2020). Stochastic optimization with heavy-tailed noise via accelerated gradient clipping.  $^1 \quad \ddagger \quad \tilde{n}$  , 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.

- Merity, S., Keskar, N. S., and Socher, R. (2017). Regularizing and optimizing lstm language models.
- Mosbach, M., Andriushchenko, M., and Klakow, D. (2020). On the stability of fine-tuning bert: Misconceptions, explanations, and strong baselines.
- 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.
- Peters, M. E., Ammar, W., Bhagavatula, C., and Power, R. (2017). Semi-supervised sequence tagging with bidirectional language models.

- 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.
- Zhang, J., Karimireddy, S. P., Veit, A., Kim, S., Reddi, S. J., Kumar, S., and Sra, S. (2020). Why are adaptive methods good for attention models? *arXiv preprint arXiv:2006.05964*, 33.