

Sign-SGD with Heavy Tails and Differential Privacy

Speaker: [Alexey Kravatsky](#), 3rd student, MIPT

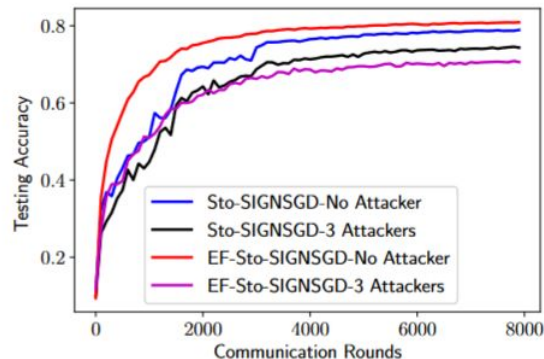
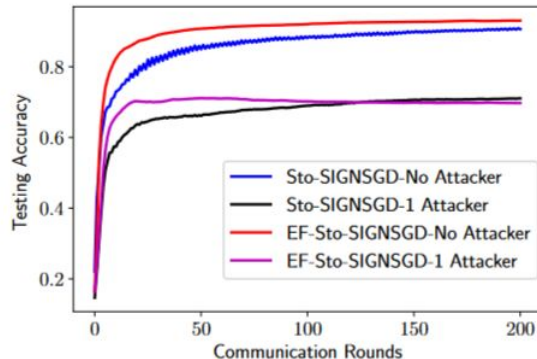
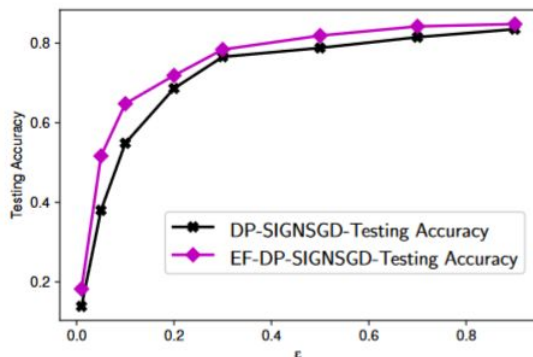
Advisor: [Savelii Chezhegov](#)

03/10/2025, MIPT

Meld the approaches of two papers

- Paper 1: relatively general proofs for vanilla Sign-SGD convergence
- Paper 2: weaker proofs, but theoretical and experimental analysis of the application to the federated learning

We want to mix it to create an algorithm to securely train LLMs on user data.
However, we will start from MNIST to test our ideas.



Algorithm 1 SignSGD

Input: Starting point $x^1 \in \mathbb{R}^d$, number of iterations T ,
stepsizes $\{\gamma_k\}_{k=1}^T$.

1: **for** $k = 1, \dots, T$ **do**

2: Sample ξ^k and compute estimate $g^k = \nabla f(x^k, \xi^k)$;

3: Set $x^{k+1} = x^k - \gamma_k \cdot \text{sign}(g^k)$;

4: **end for**

Output: uniformly random point from $\{x^1, \dots, x^T\}$.

Algorithm 1 Stochastic-Sign SGD with majority vote

Input: learning rate η , current hypothesis vector $w^{(t)}$, M workers each with an independent gradient $\mathbf{g}_m^{(t)}$, the 1-bit compressor $q(\cdot)$.

on server:

pull $q(\mathbf{g}_m^{(t)})$ from worker m .

push $\tilde{\mathbf{g}}^{(t)} = \text{sign}(\frac{1}{M} \sum_{m=1}^M q(\mathbf{g}_m^{(t)}))$ to all the workers.

on each worker:

update $w^{(t+1)} = w^{(t)} - \eta \tilde{\mathbf{g}}^{(t)}$.

Definition 2. For any given gradient $\mathbf{g}_m^{(t)}$, the compressor $dp\text{-sign}$ outputs $dp\text{-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)$. The i -th entry of $dp\text{-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)$ is given by

$$dp\text{-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)_i = \begin{cases} 1, & \text{with probability } \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \\ -1, & \text{with probability } 1 - \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \end{cases} \quad (15)$$

Assumption 1 (Lower bound). *The objective function f is lower bounded by $f^* > -\infty$, i.e., $f(x) \geq f^*, \forall x \in \mathbb{R}^d$.*

Assumption 2 (Smoothness). *The objective function f is differentiable and L -smooth, i.e., for the positive constant L*

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2, \quad \forall x, y \in \mathbb{R}^d.$$

Assumption 3 (Heavy-tailed noise in gradient estimates). *The unbiased estimate $\nabla f(x, \xi)$ has bounded κ -th moment $\kappa \in (1, 2]$ for each coordinate, i.e., $\forall x \in \mathbb{R}^d$:*

- $\mathbb{E}_\xi[\nabla f(x, \xi)] = \nabla f(x),$
- $\mathbb{E}_\xi[|\nabla f(x, \xi)_i - \nabla f(x)_i|^\kappa] \leq \sigma_i^\kappa, i \in \overline{1, d},$

where $\vec{\sigma} = [\sigma_1, \dots, \sigma_d]$ are non-negative constants. If $\kappa = 2$, then the noise is called a bounded variance.

The idea for our proofs is to be borrowed from Kornilov et al.

Lemma 1 (SignSGD Convergence Lemma). *Consider lower-bounded L -smooth function f (As. 1, 2) and HT gradient estimates (As. 4). Then Alg. 1 after T iterations with constant stepsizes $\gamma_k \equiv \gamma$ achieves with probability at least $1 - \delta$ starting with $\Delta_1 = f(x^1) - f^*$:*

$$\begin{aligned} \frac{1}{T} \sum_{k=1}^T \|\nabla f(x^k)\|_1 &\leq \frac{2\Delta_1}{T\gamma} + 16Ld\gamma \log(1/\delta) + 4\|\vec{\sigma}\|_1 \\ &\quad + 12 \frac{d\|\nabla f(x^1)\|_1}{T} \log(1/\delta). \end{aligned} \quad (2)$$

From China with love

Definition 3. Given a set of local datasets \mathcal{D} provided with a notion of neighboring local datasets $\mathcal{N}_{\mathcal{D}} \subset \mathcal{D} \times \mathcal{D}$ that differ in only one data point. For a query function $f : \mathcal{D} \rightarrow \mathcal{X}$, a mechanism $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{O}$ to release the answer of the query is defined to be (ϵ, δ) -locally differentially private if for any measurable subset $\mathcal{S} \subseteq \mathcal{O}$ and two neighboring local datasets $(D_1, D_2) \in \mathcal{N}_{\mathcal{D}}$,

$$P(\mathcal{M}(f(D_1)) \in \mathcal{S}) \leq e^\epsilon P(\mathcal{M}(f(D_2)) \in \mathcal{S}) + \delta. \quad (21)$$

A key quantity in characterizing local differential privacy for many mechanisms is the sensitivity of the query f in a given norm l_r , which is defined as

$$\Delta_r = \max_{(D_1, D_2) \in \mathcal{N}_{\mathcal{D}}} \|f(D_1) - f(D_2)\|_r. \quad (22)$$

Theorem 6. Let u_1, u_2, \dots, u_M be M known and fixed real numbers. Further define random variables $\hat{u}_i = \text{dp-sign}(u_i, \epsilon, \delta), \forall 1 \leq i \leq M$. Then there always exist a constant σ_0 such that when $\sigma \geq \sigma_0$, $P(\text{sign}(\frac{1}{M} \sum_{m=1}^M \hat{u}_i) \neq \text{sign}(\frac{1}{M} \sum_{m=1}^M u_i)) < [(1 - x^2)]^{\frac{M}{2}}$, where $x = \frac{|\sum_{m=1}^M u_m|}{2\sigma M}$.

This is impractical. Our goal is to condense it.

Theorem 2. Suppose Assumptions 1, 2 and 4 are satisfied, and the learning rate is set as $\eta = \frac{1}{\sqrt{Td}}$. Then by running Algorithm 1 with $q(\mathbf{g}_m^{(t)}) = \text{sto-sign}(\nabla f_m(w^{(t)}), \mathbf{b})$ (termed as *Sto-SIGNSGD*) for T iterations, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \|\nabla F(w^{(t)})\|_1 &\leq \frac{1}{c} \left[\frac{\mathbb{E}[F(w^{(0)}) - F(w^{(T+1)})] \sqrt{d}}{\sqrt{T}} + \frac{L\sqrt{d}}{2\sqrt{T}} + \frac{2}{T} \sum_{t=1}^T \sum_{i=1}^d |\nabla F(w^{(t)})_i| \mathbf{1}_{p_i^{(t)} > \frac{1-c}{2}} \right] \\ &\leq \frac{1}{c} \left[\frac{(F(w^{(0)}) - F^*) \sqrt{d}}{\sqrt{T}} + \frac{L\sqrt{d}}{2\sqrt{T}} + 2 \sum_{i=1}^d b_i \Delta(M) \right], \end{aligned} \tag{10}$$

where $0 < c < 1$ is some positive constant, $p_i^{(t)}$ is the probability that the aggregation on the i -coordinate of the gradient is wrong during the t -th communication round, and $\Delta(M)$ is the solution to $(1 - x^2)^{\frac{M}{2}} = \frac{1-c}{2}$. The second inequality is due to the fact that $p_i^{(t)} > \frac{1-c}{2}$ only if $\frac{|\nabla F(w^{(t)})_i|}{b_i} \leq \Delta(M)$.

Bibliography

- Jin et al., 2020: Chinese [article](#) with compressors and guarantees of differential privacy
- Kornilov et al., 2025: Russian [preprint](#) with proofs of high-probability convergence in the case of heavy-tailed noise
- These articles are key to our research. Others will be used only as a reference.