

Sign SGD with Heavy-Tailed Noise and Differential Privacy

Alexey Kravatskiy

`kravtskii.aiu@phystech.edu`

Anton Plusnin

`plusnin.aa@phystech.edu`

Savelii Chezhegov

`chezhegov.sa@phystech.edu`

Alexander Beznosikov

`beznosikov.an@phystech.edu`

April 3, 2025

Abstract

In the era of large models, federated learning has become indispensable like never before. Sound modern federated learning must meet three key requirements. First, the method must process correctly real-world data, which, in case of Large Language Models, means the algorithm must tolerate noise with unbounded variance. Second, to ensure applicability, the algorithm must converge under these conditions with high probability, not only in expectation. Third, the whole procedure must not jeopardize user data. To address these natural requirements, we have constructed a novel modification of Sign version of Stochastic Gradient Descent. In this paper, we demonstrate that it meets all three earlier stated requisites. We start with proving algorithm’s high-probability convergence on data with heavy-tailed noise. Then, we prove its differential privacy. Finally, we show the superior performance of the algorithm in training Large Language Models.

Keywords: Sign SGD, differential privacy, high-probability convergence, federated learning, heavy-tailed noise.

Highlights:

1. Sign Stochastic Gradient Descent can be used to train LLMs on real data.

2. Our modification of Sign Stochastic Gradient Descent keeps user data private.
3. Our modification of Sign Stochastic Gradient Descent does not require tuning.

1 Introduction

Federated Learning is a useful method to train models that require large amounts of data. Indeed, it is often the case that the data is distributed across multiple devices, like mobile phones [9], and it is not only costly to collect all the data in one place, but often also unacceptable due to the requirements of privacy. On the other hand, training the model only on the local data for a particular user is impossible, as the model is large. Hence, a need for a joint training procedure arises. We come to a setting with a server and a number of workers, where each worker has a local dataset. The goal is to train a model on the data from all the workers, without sharing the data between the workers or with the server.

The most obvious way to train a model in a federated setting is to use Stochastic Gradient Descent (SGD) [12] by passing the gradient to the server (add citation). However, when transmitting the gradient itself, the communication cost is unaffordably high and the privacy of the data is compromised. To address this issue, one can use the Sign Stochastic Gradient Descent algorithm [1]. This algorithm transmits only the signs of the coordinates of the gradients, which is much cheaper in terms of communication.

Clipping the norm of gradient estimate before SGD step, which comprises the method called ClipSGD, is an another great idea that demonstrates positive empirical results [11, 4]. Nonetheless, the clipping requires meticulous tuning of the clipping threshold, which was shown to depend on not only iteration number but also the objective function and the noise [13, Theorem 3.1]. In case of federated learning, the search of the threshold is even more complicated, as the data is distributed across the workers and cannot be used for tuning, and the objective function is more complex due to sophistication and size of the model.

The natural simplification of the clipping is normalization of the gradient, which lies at the core of NSGD [5]. NSGD outperformed ClipSGD on the task of sequence labeling via LSTM Language Models [10]. Despite this fact, NSGD, unlike ClipSGD, requires larged batch sizes for convergence, which though can be mended by applying momentum [2]. The major drawback of NSGD is the absence of proofs of convergence with high probability. Moreover, the method seems to be not private.

As to guarantee the privacy is our priority, we opted to base our algorithm on SignSGD.

The convergence with high probability of **SignSGD** for the heavy-tailed noise was recently proved [8], which makes **SignSGD** a perfect candidate for federated learning. Simultaneously, **SignSGD** was already used as a base to create a differentially private algorithm [7]. However, in both cases, no proofs for the convergence of the algorithm in the modern federated learning setting were provided.

In this paper, we propose a modification of **SignSGD** that can be used to train LLMs on real data. We show that our algorithm converges with high probability, even in the presence of heavy-tailed noise, and does not require tuning. We also show that our algorithm is differentially private. Finally, we test the algorithm on the training of LLMs.

2 Problem statement.

To-do: what kind of distribution S is. Change the abstract. It should start with narrow problem. Remove "can" and variants from the paper altogether!

First, we need to state that we work with stochastic optimization.

Stochastic optimization problem. The stochastic optimization problem for a smooth non-convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is:

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim \mathcal{S}}[f(x, \xi)], \quad (1)$$

where random variable ξ is sampled from an unknown distribution \mathcal{S} . The gradient oracle returns unbiased gradient estimate $\nabla f(x, \xi) \in \mathbb{R}^d$. In machine learning, for instance, $f(x, \xi)$ is a loss function on a sample ξ [14].

The most popular algorithm to solve (1) is Stochastic Gradient Descent (SGD) [12]:

$$x^{k+1} = x^k - \gamma_k \cdot g^k, \quad g^k := \nabla f(x^k, \xi^k).$$

For non-convex functions, the algorithm must stop at the point with sufficiently small gradient norm. We will apply the algorithm to the federated optimization problem.

Federated optimization problem. Let $I = X \times Y$ be a sample space, where X is a space of feature vectors and Y is a label space. For the hypothesis space $\mathcal{W} \subseteq \mathbb{R}^d$, a loss function is defined as $l : \mathcal{W} \times I \rightarrow \mathbb{R}$ which measures the loss of the prediction on the data point

$(x, y) \in I$ based on the hypothesis $w \in \mathcal{W}$. For a dataset $D \subset I$, the global loss function $F : \mathcal{W} \rightarrow \mathbb{R}$ is defined as

$$F(w) = \frac{1}{|D|} \sum_{(x,y) \in D} l(w; (x, y)). \quad (2)$$

In case of distributed optimization, the dataset is split between M workers. Each worker m has a local dataset $D_m \subset I$ and a local function $f_m : \mathbb{R}^d \rightarrow \mathbb{R}$ defined as

$$f_m(w) = \frac{1}{|D_m|} \sum_{(x_n, y_n) \in D_m} l(w; (x_n, y_n)), \quad (3)$$

where $|D_m|$ is the size of worker m 's local dataset D_m .

Thus, our goal is to solve the following federated optimization problem:

$$\min_{w \in \mathbb{R}^d} F(w) \quad \text{where} \quad F(w) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M f_m(w). \quad (4)$$

We assume that the data are distributed over the workers uniformly, consequently, $\mathbb{E}[f_m(w)] = F(w)$ for workers' data distribution.

Now, let us introduce the requirements for the algorithm to solve the federated optimization problem.

Heavy-tailed noise. Noise has bounded κ -th moment for some $\kappa \in (1, 2]$, i.e. $\mathbb{E}_\xi[\|\nabla f(x, \xi) - \nabla f(x)\|_2^\kappa] \leq \sigma^\kappa$. In particular, the noise can have unbounded variance, i.e. $\kappa < 2$.

Differential privacy. Additionally, the algorithm must be private, which means it must satisfy (ϵ, δ) -local differential privacy [3].

High-probability convergence. The algorithm must have convergence guarantees which hold true with probability at least $1 - \delta, \delta \in (0, 1)$.

3 Theory

3.1 The Algorithm and the compressor

The algorithm we are working with can be defined generally as follows:

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

on each worker:

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

As the algorithm must be differentially private, we use as a 1-bit compressor dp-sign compressor [7]:

Definition 1. For any given gradient $\mathbf{g}_m^{(t)}$, the compressor *dp-sign* outputs $\text{dp-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)$. The i -th entry of $\text{dp-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)$ is given by

$$\text{dp-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)_i = \begin{cases} 1, & \text{with probability } \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \\ -1, & \text{with probability } 1 - \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \end{cases} \quad (5)$$

where $\Phi(\cdot)$ is the cumulative distribution function of the normalized Gaussian distribution; $\sigma = \frac{\Delta_2}{\epsilon} \sqrt{2 \ln(\frac{1.25}{\delta})}$, where ϵ and δ are the differential privacy parameters and Δ_2 is the sensitivity measure.

As stated in Theorem 5 from [7], the *dp-sign* mechanism is (ϵ, δ) -differentially private for any ϵ and $\delta \in (0, 1)$.

3.2 Convergence for dp-sign

Theorem 6 from [7] should establish probability of the dissimilarity of majority sign of the gradients and majority sign of the dp-signs:

Theorem 1. 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 $x = \frac{|\sum_{m=1}^M u_m|}{2\sigma M}$.

However, as we recently found out, the theorem is fundamentally flawed, as it makes σ a parameter, while it follow from the definition of the *dp-sign* compressor that σ is a function of ϵ and δ .

Let us suppose that there exist ϵ and δ such that for $\sigma(\epsilon, \delta)$, the inequality on probability from the theorem holds true. Then, let us walk through the proof to find a precise lower bound for σ . If this bound is lower than $\sigma(\epsilon, \delta)$, the theorem is indeed correct.

In the proof of Theorem 6, the authors needed to find the lower bound for the following expression:

$$\frac{1}{\sqrt{2\pi}\sigma} \left[\left| \sum_{m=1}^M u_m \right| e^{-\frac{u_1^2}{2\sigma^2}} + \left| \sum_{m=2}^M u_m \right| \left[e^{-\frac{(\sum_{m=2}^M u_m)^2}{2\sigma^2}} - 1 \right] \right]$$

Instead of taking limit $\sigma \rightarrow \infty$ as they did, we will use the well known relation $e^{-x^2} \geq 1 - x^2$. After applying it to the earlier mentioned expression and making some trivial transformations, we get the following:

$$\sigma^2 \geq \frac{7}{5} \left(u_1^2 + \left| \sum_{m=1}^M u_m \right|^2 + \frac{|u_1|^3}{\left| \sum_{m=1}^M u_m \right|} \right).$$

It is a sufficient condition, not a necessary one. However, it reflects the key features of the condition on σ . $\frac{|u_1|^3}{\left| \sum_{m=1}^M u_m \right|}$ is an extremely unreliable term, as $|u_1|$ may be high, especially for heavy-tailed noise, while $\left| \sum_{m=1}^M u_m \right|$ may be small (it is easy to construct an appropriate example, with 3 workers). Hence, the bound on σ is not only high, but also unstable. Consequently, we have no proofs whatsoever of the convergence of the algorithm.

Moreover, our misgivings are supported by the fact that in an updated version [6] of the article [7], the authors have removed all mentions of dp-sign. Right now, we are facing the problem that the algorithm for which we sought proofs for more general type of noise, might not make sense at all. This Friday, we are going to discuss the results and update this section.

4 Experiments

In this section, we present the experimental results for the methods we discussed in Sections 2 and 3. First, we applied the algorithms to a classic machine learning problem. Then, we tested the algorithms on the training of Large Language Models.

4.1 Synthetic noise.

We test our algorithm on the method of logistic regression for UCI Mushroom Dataset. The dataset consists of 6,449 training samples and 1,625 testing samples. Each sample has 112 features, and represents a mushroom either poisonous or edible. We apply different algorithms to train the logistic regression model. The algorithms are run on 20 workers, with training

data distributed between the workers equally. We compare the cases when there is no noise and when the noise is normal with $\sigma^2 = 1/2$ (to-do: when we add Student's distribution, we'll have a heavy-tailed noise). We set the learning rate $\gamma = 0.05$ for SGD and 0.02 for SignSGD, and there are 2000 iterations for each algorithm. Additionally, we model the time spent on sending the data to the server by ascribing 0.5 ms spent time to each sent bit. The results are presented in fig. 1.

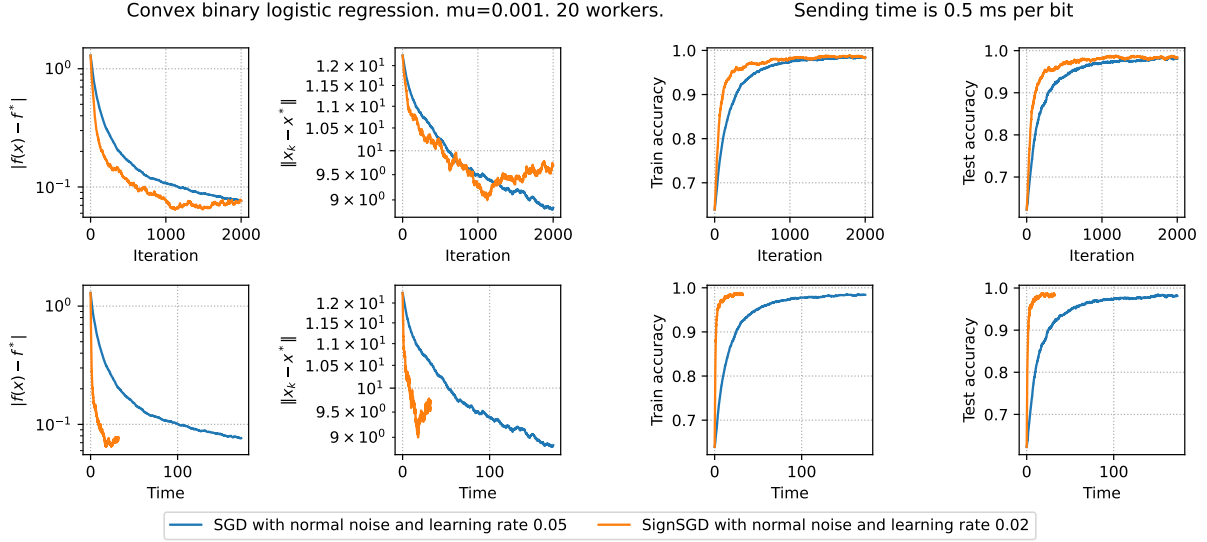


Figure 1: Logistic regression on UCI Mushroom Dataset. Performance of SGD and SignSGD with majority voting.

As we see, Sign-SGD with majority voting has run correctly. It not only has run faster than SGD, but also has delivered the same test accuracy as SGD did. The algorithm has also been more robust to the noise, as it was able to converge even with the heavy-tailed noise. Crucially, the algorithm has been differentially private, as it has run with dp-sign compressor with parameters $\epsilon = 1$ and $\delta = 10^{-5}$ (this is a TO-DO).

Sidenote: I haven't created a draft of the plot or error analysis. The case of training LLMs is for now beyond my competence, and is not essential to the theoretical problem I'm solving. Hence, I start only with training on data with synthetic heavy-tailed noise. There, I already got the computer plot I planned to obtain, so there is no need in draft again.

4.2 Large models Pre-Training.

To-do: training LLMs, using torch etc.

5 Conclusion

Summarize your findings and discuss future work.

6 Acknowledgments

Optional acknowledgments section.

A Additional Proofs and Results

Include detailed proofs and supplementary materials here.

References

- [1] Jeremy Bernstein et al. “signSGD: Compressed Optimisation for Non-Convex Problems”. In: *International Conference on Machine Learning*. 2018, pp. 560–569.
- [2] Ashok Cutkosky and Harsh Mehta. “Momentum improves normalized sgd”. In: *International conference on machine learning*. PMLR. 2020, pp. 2260–2268.
- [3] C. Dwork, A. Roth, et al. “The algorithmic foundations of differential privacy”. In: *Foundations and Trends® in Theoretical Computer Science* 9.3–4 (2014), pp. 211–407.
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [5] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. “Beyond convexity: Stochastic quasi-convex optimization”. In: *Advances in Neural Information Processing Systems*. 2015, pp. 1594–1602.
- [6] Richeng Jin et al. “Sign-Based Gradient Descent With Heterogeneous Data: Convergence and Byzantine Resilience”. In: *IEEE Transactions on Neural Networks and Learning Systems* 36.2 (2025), pp. 3834–3846. DOI: 10.1109/TNNLS.2023.3345367.
- [7] Richeng Jin et al. “Stochastic-Sign SGD for Federated Learning with Theoretical Guarantees”. In: *Part of this work is published in IEEE Transactions on Neural Networks and Learning Systems, 2024* 36.2 (Feb. 25, 2020), pp. 3834–3846. ISSN: 2162-2388. DOI: 10.1109/tnnls.2023.3345367. arXiv: 2002.10940 [cs.LG].

- [8] Nikita Kornilov et al. *Sign Operator for Coping with Heavy-Tailed Noise: High Probability Convergence Bounds with Extensions to Distributed Optimization and Comparison Oracle*. 2025. DOI: 10.48550/ARXIV.2502.07923.
- [9] Brendan McMahan et al. “Communication-Efficient Learning of Deep Networks from Decentralized Data”. In: *Artificial Intelligence and Statistics*. 2017, pp. 1273–1282.
- [10] Stephen Merity, Nitish Shirish Keskar, and Richard Socher. “Regularizing and optimizing LSTM language models”. In: *arXiv preprint arXiv:1708.02182* (2017).
- [11] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. “On the difficulty of training recurrent neural networks”. In: *International conference on machine learning*. 2013, pp. 1310–1318.
- [12] Herbert Robbins and Sutton Monro. “A stochastic approximation method”. In: *The annals of mathematical statistics* (1951), pp. 400–407.
- [13] Abdurakhmon Sadiev et al. “High-Probability Bounds for Stochastic Optimization and Variational Inequalities: the Case of Unbounded Variance”. In: *arXiv preprint arXiv:2302.00999* (2023).
- [14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.