

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

May 15, 2025

Abstract

Crucial for large-scale models, federated learning faces two major challenges: privacy preservation and high communication costs. While SignSGD addresses the communication issue by transmitting only gradient signs, its only earlier proposed private version lacks proper privacy guarantees and convergence analysis. We construct a new variant of DP-SignSGD that combines Gaussian noise with Bernoulli subsampling to achieve true differential privacy. Our approach satisfies (α, ϵ_R) -Rényi differential privacy, which can be readily converted to standard (ϵ, δ) -privacy guarantees. We demonstrate the algorithm’s performance on logistic regression problem and classification of handwritten digits with MLP and CNN. The main challenge remains the tradeoff between precision of a single iteration and the maximum number of privacy-preserving iterations. Our analysis suggests that the sign mechanism’s binary output and potential gradient privacy may provide additional privacy guarantees beyond our current calculations. The algorithm can be readily adapted to tighter privacy bounds, and we identify the need for theoretical convergence guarantees as the primary direction for future research.

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

Idea: We add Gaussian noise and use privacy amplification by subsampling to make sign-SGD private.

25 **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 is applicable to federated setting.

1 Introduction

30 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 also often 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 and data of one
35 user is insufficient. 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
40 Descent (SGD) [15] by passing the gradient to the server. However, when transmitting the gradient itself, the communication cost is unaffordably high, and again the data privacy might be violated. To address this issue, one can use the Sign Stochastic Gradient Descent (SignSGD) algorithm [1]. This algorithm transmits only the signs of the coordinates of the gradients, which is much cheaper in terms of communication.

45 Communication costs aside, clipping the norm of gradient estimate before SGD step, which comprises the method called ClipSGD, is another great idea that demonstrates positive empirical results [13, 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 [16, Theorem 3.1]. In case of federated learning, the search of the
50 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
 55 LSTM Language Models [10]. Despite this fact, NSGD, unlike ClipSGD, requires large batch sizes for convergence, which, however, 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 far from private.

As to guarantee the privacy is our priority, we opted to base our algorithm on SignSGD.
 60 The convergence with high probability of SignSGD for the heavy-tailed noise has been recently proved [8], which makes SignSGD a perfect candidate for federated learning. Simultaneously, SignSGD has already been used as a base to create an allegedly differentially private algorithm [7] (the authors proved differential privacy of only one step of the algorithm). However, in both cases, no proofs for the convergence of the algorithm in the modern federated learning
 65 setting were provided.

In this paper, based on the concept of Rényi differential privacy, we develop the idea from [7] to derive truly differentially-private modification of SignSGD, which we naturally call DP-SIGNSGD. We show that our algorithm converges with high probability, even in the presence of heavy-tailed noise. Finally, we test the algorithm on the training of LLMs.

70 2 Problem statement.

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
 75 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 [17].

The most popular algorithm to solve (1) is Stochastic Gradient Descent (SGD) [15]:

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

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.

80 **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

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

85 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

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

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,
90 $\mathbb{E}[f_m(w)] = F(w)$ for workers' data distribution.

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

Differential privacy. Additionally, the algorithm must be private. This is guaranteed by (ϵ, δ) -local differential privacy [3]:

95 **Definition 1.** 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 $(D_1, D_2) \in \mathcal{N}_{\mathcal{D}}$,

$$\mathbb{P}(\mathcal{M}(f(D_1)) \in \mathcal{S}) \leq e^\epsilon P(\mathcal{M}(f(D_2)) \in \mathcal{S}) + \delta. \quad (5)$$

100 A key quantity in characterizing local differential privacy for many mechanisms is the sensitivity of the query f in a given norm l_r (we use l_2), which is defined as

$$\Delta_r = \max_{(D_1, D_2) \in \mathcal{N}_{\mathcal{D}}} \|f(D_1) - f(D_2)\|_r. \quad (6)$$

High-probability convergence. For given $\delta_{\text{prob}} \in (0, 1)$, the algorithm must converge with probability at least $1 - \delta_{\text{prob}}$.

3 Theory

105 **Assumption 1** (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
110 variance.

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

Our goal is to find such compressor $q(\cdot)$ that the algorithm is private.

3.2 Incorrectness of earlier proposed DP-SIGN

115 Some time ago, a 1-bit compressor $dp\text{-sign}$ [7] was proposed:

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

$$dp\text{-sign}(\mathbf{g}_m^{(t)}, \epsilon, \delta)_i = \begin{cases} 1, & \text{with probability } \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \\ -1, & \text{with probability } 1 - \Phi\left(\frac{(\mathbf{g}_m^{(t)})_i}{\sigma}\right) \end{cases} \quad (7)$$

120 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 $dp\text{-sign}$ mechanism is (ϵ, δ) -differentially private for any ϵ and $\delta \in (0, 1)$.

Theorem 6 from [7] should establish probability of the dissimilarity of majority sign of the gradients and majority sign of the $dp\text{-sign}$ s:

125 **Theorem 1.** Let u_1, u_2, \dots, u_M be M known and fixed real numbers. Further define random variables $\hat{u}_i = dp\text{-sign}(u_i, \epsilon, \delta), \forall 1 \leq i \leq M$. 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 $x = \frac{|\sum_{m=1}^M u_m|}{2\sigma M}$.

130 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 $dp\text{-sign}$ 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 for the case of 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.

The problems with convergence of dp-sign are not limited by the incorrect Theorem 6. The authors have overlooked another issue of utmost importance. Proving the (ϵ, δ) differential privacy of dp-sign, they did so for only one call of dp-sign operator. However, according to [3] (Theorem 3.16), the overall privacy of running dp-sign T times is $(T\epsilon, T\delta)$. Consequently, when we set target $\epsilon \approx 10$, $\delta \approx 10^{-5}$, as is recommended for machine learning purposes [14], for $T = 100$ we get 100 times lower ϵ and δ for dp-sign. Judging by our experiments, this destroys convergence. That is intuitive: high noise for low (ϵ, δ) leads to poorer convergence, and this forces us to use still lower (ϵ, δ) to preserve the required level of privacy during the run of the algorithm.

To the best of our knowledge, there are no differentially private versions of sign-sgd reported except this incorrect dp-sign. As Gaussian noise is still the best mechanism to ensure differential privacy (Laplace noise and exponential mechanism complicate convergence), we preserve the idea of dp-sign. However, we drastically change the privacy analysis to make σ small enough for the algorithm to converge.

3.3 Rényi Differential privacy

To make an economical use of differential privacy, we have to introduce another type of differential privacy, which gives more tight bounds on composition.

Definition 3 (Rényi divergence [11]). *Let P and Q be two distributions on \mathcal{X} defined over the same probability space, and let p and q be their respective densities. The Rényi divergence*

160 of a finite order $\alpha \neq 1$ between P and Q is defined as

$$D_\alpha(P \parallel Q) \triangleq \frac{1}{\alpha - 1} \ln \int_{\mathcal{X}} q(x) \left(\frac{p(x)}{q(x)} \right)^\alpha dx.$$

Rényi divergence at orders $\alpha = 1, \infty$ are defined by continuity.

Definition 4 (Rényi differential privacy (RDP) [11]). *We say that a randomized mechanism $\mathcal{M}: \mathcal{S} \rightarrow \mathcal{R}$ satisfies (α, ε) -Rényi differential privacy (RDP) if for any two adjacent inputs $S, S' \in \mathcal{S}$ it holds that*

$$D_\alpha(\mathcal{M}(S) \parallel \mathcal{M}(S')) \leq \varepsilon.$$

165 The notion of adjacency between input datasets is domain-specific and is usually taken to mean that the inputs differ in contributions of a single individual. In this work, we will use this definition and call two datasets S, S' to be adjacent if $S' = S \cup \{x\}$ for some x (or vice versa).

Definition 5 (Sampled Gaussian Mechanism (SGM) [12]). *Let f be a function mapping*
 170 *subsets of \mathcal{S} to \mathbb{R}^d . We define the Sampled Gaussian mechanism (SGM) parameterized with the sampling rate $0 < q \leq 1$ and the noise $\sigma > 0$ as*

$$\text{SG}_{q,\sigma}(S) \triangleq f(\{x: x \in S \text{ is sampled with probability } q\}) + \mathcal{N}(0, \sigma^2 \mathbb{I}^d),$$

where each element of S is sampled independently at random with probability q without replacement, and $\mathcal{N}(0, \sigma^2 \mathbb{I}^d)$ is spherical d -dimensional Gaussian noise with per-coordinate variance σ^2 .

175 From Section 3.3 (“Numerically stable computations”) of [12], for each integer $\alpha \geq 1$, SGM applied to a function with $\Delta_2 \leq 1$ is (α, ε_R) -RDP for the following ε_R :

$$\varepsilon_R = \frac{1}{\alpha - 1} \log \left(\sum_{k=0}^{\alpha} \binom{\alpha}{k} (1 - q)^{\alpha-k} q^k \exp \left(\frac{k^2 - k}{2\sigma^2} \right) \right) \quad (8)$$

Clearly, if $\Delta_2 > 1$ for a mechanism, $\sigma := \Delta_2 \sigma$ will deliver the same (α, ε_R) -RDP.

If an (α, ε_R) -RDP mechanism is applied T times, by adaptive sequential composition [11], we get $(\alpha, T\varepsilon_R)$ -RDP of the procedure. From *Proposition 3* from [11], it follows that the
 180 procedure is also $(\varepsilon_R + \frac{\log 1/\delta}{\alpha-1}, \delta)$ -differentially private for any $0 \leq \delta \leq 1$. Thus, the procedure is (ε, δ) -differential privacy, if (α, ε_R) satisfy:

$$\varepsilon_R \leq \varepsilon/T - \frac{\log 1/\delta}{T(\alpha - 1)} \quad (9)$$

As convergence of the algorithm depends only on σ , while α is an internal parameter, we are interested in finding the minimal σ . This we achieve via grid search on integer $1 \leq \alpha \leq 20$ and $0.5 \leq \sigma \leq 3$. An example can be seen in Figure 1.

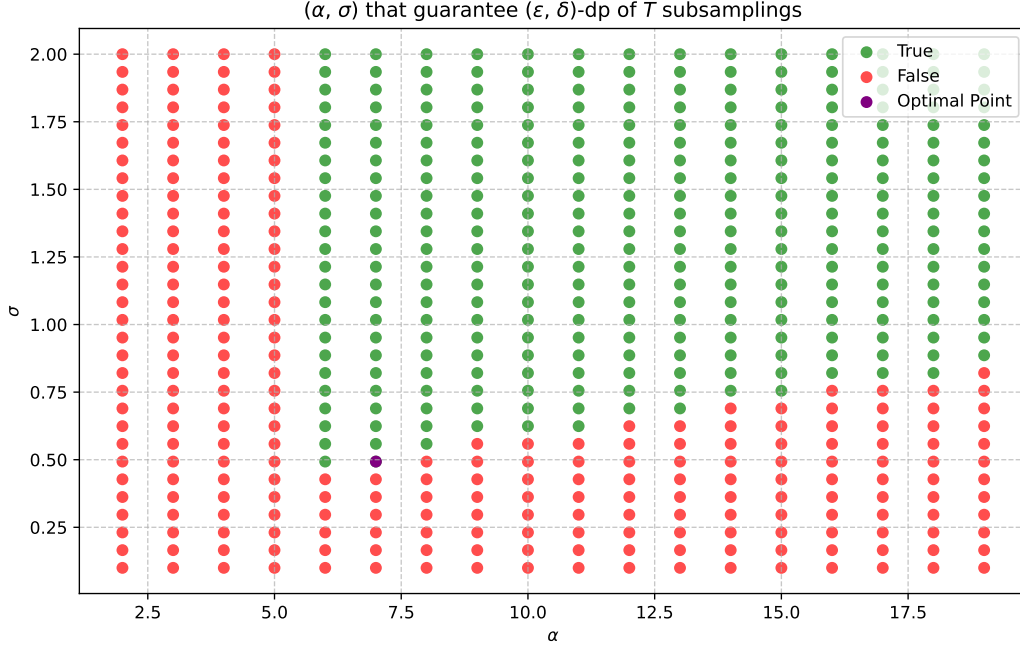


Figure 1: Finding the minimal σ for SGM with $\Delta_2 \leq 1$, sampling rate $q = 1/300$, and $T = 1000$

Judging by the chart, there is no need in searching among $\alpha \in \mathbb{R}$: *Case II. Fractional α* from [12] gives a series expression to obtain a bound for ε_R , and the series is continuous on α . Continuous borderline on 1 implies that σ we get by the grid search is close to the optimum.

3.4 Constructing DP-SIGN

Let us define DP-SIGN 2. From the previous section, (ε, δ) -privacy of DP-SIGN immediately follows.

Theorem 2. *Let $l(\cdot; \cdot)$ be a loss function and let $\varepsilon > 0$, $\delta \in (0, 1)$. Then applying T times $dp\text{-sign}(w_i, l, D, (\varepsilon, \delta), T, q, C)$ in different points w_i on a user database D with a fixed sampling rate q and fixed clipping level $C > 0$ is (ε, δ) -differentially private.*

Algorithm 2 DP-SIGN compressor

Input: coordinate w , loss function l , user database D , (ε, δ) -privacy requirement, number of iterations T , sampling rate q , clipping level C .

Prepare subsample S : add each element $(x, y) \in D$ with probability q .

Compute the gradient \mathbf{g} of the subsample: $\frac{1}{|S|} \sum_{(x,y) \in S} l(w; (x, y))$. If S is empty, let $\mathbf{g} = 0$.

If $\|\mathbf{g}\|_2 > C$: $\mathbf{g} = C \frac{\mathbf{g}}{\|\mathbf{g}\|_2}$.

Grid search $\sigma(q, T, \varepsilon, \delta)$ to satisfy (8) and (9).

$\mathbf{g}_{priv} = \mathbf{g} + \mathcal{N}(0, (C\sigma)^2 \mathbb{I}^d)$

Output: $\text{sign}(\mathbf{g}_{priv})$

Proof. \mathbf{g}_{priv} is private due to the satisfaction of (8) and (9). As RDP preserves under
195 post-processing ([11]), $\text{sign}(\mathbf{g}_{priv})$ is also (α, ε_R) -RPD, which entails (ε, δ) -dp for all T
iterations. \square

4 Experiments

In this section, we present the experimental results for the methods we discussed in Sections
2 and 3. First, we apply the algorithms to a classic machine learning problem. Then, we test
200 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 $d = 112$
features, and represents a mushroom either poisonous or edible. We apply different algorithms
205 to train the logistic regression model. The algorithms are run for $T = 100000$ iterations on
10 workers, with training data distributed between the workers equally. We compare the
cases when there is no noise, when the noise is normal with $\sigma = 1/4$, and when the noise
is coordinatewise Lévy stable with $\sigma_l = 1/4$, $\alpha_l = 1.6$, which corresponds to $\kappa = 1.5$, and
 $\beta_l = 0$ (this distribution is defined by its characteristic function $\varphi(t) = \exp(-0.25^{1.6}|t|^{1.6})$).
210 For private algorithms, we set $\epsilon \approx 10$, $\delta \approx 1/n^{1.1}$, where n is the size of local user database,
as is recommended for machine learning purposes [14]. We set the learning rate $\gamma = 1/\sqrt{dT}$
for SGD, sign-SGD, and DP-SIGNSGD. The results are presented in 2.

As we see, DP-SIGNSGD does converge, even under heavy-tailed noise.

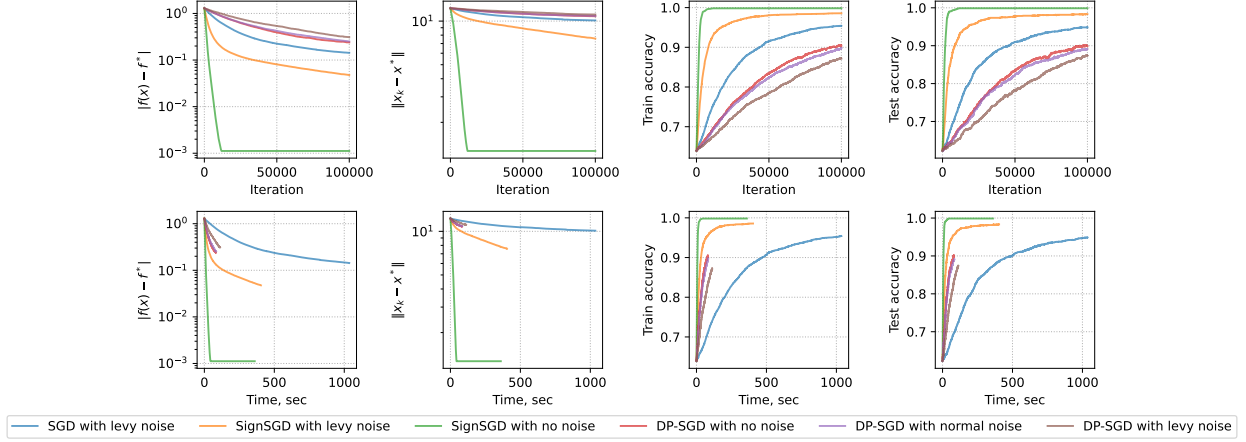


Figure 2: Logistic regression on UCI Mushroom Dataset.
Performance of SGD, SignSGD, and DP-SignSGD with majority voting.

Table 1: DP-Sign SGD methods

Method	Citation	Is (ϵ, δ) private	Converges
ChineseDP-SignSGD-cheat	[7]	No	Yes
ChineseDP-SignSGD rectified	[7]	Yes	No
Rényi DP-SignSGD	This paper	Yes	Yes

4.2 DP-Sign (outdated)

Having verified Sign-SGD, we test dp-sign on the same dataset. We present three variants of the dp-sign:

1. DP-SignSGD from [7] with unrectified (ϵ, δ) . The algorithm converges, but is not sufficiently private, as we highlighted earlier.
2. DP-SignSGD from [7] with rectified $(\epsilon/T, \delta/T)$, where T is 1000. The algorithm is now private, but does not converge.
3. Our modification of DP-SignSGD. We used Rényi differential privacy [3]. This type of privacy can be easily converted to classical differential privacy, but has a more temperate deterioration of privacy during the sequential composition: the noise $\sigma \sim \sqrt{T}$ instead of $\sigma \sim T$. We will expand Theory and introduce there this algorithm, when it will be ready. For now, it as dp-sign does not converge.

Thus, we have the following table 1.

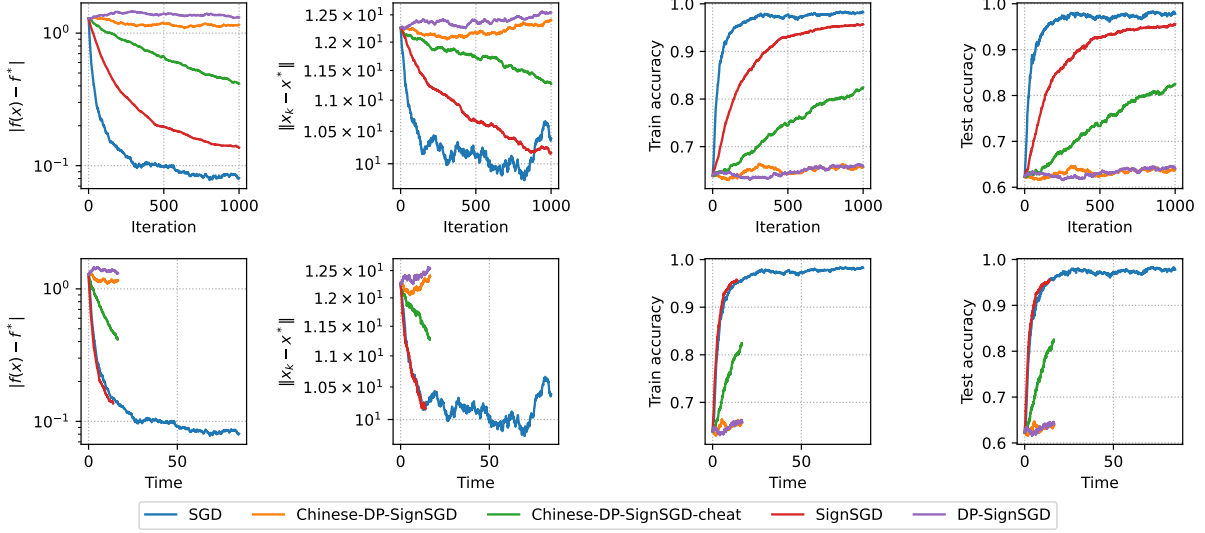


Figure 3: Logistic regression on UCI Mushroom Dataset. Performance of SGD, SignSGD and different variants of DP-SignSGD with majority voting.

5 Conclusion

We have highlighted the deficiencies of the previous approach of working with dp-sign-SGD. With a use of Rényi differential privacy and advanced numerically-produced bounds for
230 Sampled Gaussian Mechanism, we have consturcted a privacy accountant that makes dp-sign-SGD converge and private, which, to the best of our knowledge, has never been done before.

6 Acknowledgments

We will fill this section when we are done with the article.

235 A Additional Proofs

We will present them after we complete the tests and finalize the algorithm.

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] Ilya Mironov. “Rényi Differential Privacy”. In: *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, Aug. 2017, pp. 263–275. DOI: 10.1109/csf.2017.11. URL: <http://dx.doi.org/10.1109/CSF.2017.11>.
- [12] Ilya Mironov, Kunal Talwar, and Li Zhang. *Rényi Differential Privacy of the Sampled Gaussian Mechanism*. 2019. arXiv: 1908.10530 [cs.LG]. URL: <https://arxiv.org/abs/1908.10530>.
- [13] 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.

- [14] Natalia Ponomareva et al. “How to DP-fy ML: A Practical Guide to Machine Learning with Differential Privacy”. In: *Journal of Artificial Intelligence Research* 77 (July 2023). arXiv:2303.00654 [cs], pp. 1113–1201. ISSN: 1076-9757. DOI: 10.1613/jair.1.14649. URL: <http://arxiv.org/abs/2303.00654> (visited on 04/09/2025).
- 275 [15] Herbert Robbins and Sutton Monro. “A stochastic approximation method”. In: *The annals of mathematical statistics* (1951), pp. 400–407.
- [16] 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).
- 280 [17] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.