PRIVATE STOCHASTIC CONVEX OPTIMIZATION WITH TYSBAKOV NOISE CONDITION AND LARGE LIPSCHITZ CONSTANT

Anonymous authors

Paper under double-blind review

ABSTRACT

We study Stochastic Convex Optimization in Differential Privacy model (DP-SCO). Unlike previous studies, here we assume the population risk function satisfies the Tysbakov Noise Condition (TNC) with some parameter $\theta > 1$, where the Lipschitz constant of the loss could be extremely large or even unbounded, but the ℓ_2 -norm gradient of the loss has bounded k-th moment with $k \ge 2$. For the Lipschitz case with $\theta \ge 2$, we first propose an (ϵ, δ) -DP algorithms whose utility bound is $\tilde{O}\left(\left(\tilde{r}_{2k}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\epsilon}))^{\frac{k-1}{k}}\right)^{\frac{\theta}{\theta-1}}\right)$ in high probability, where n is the sample size, d is the model dimension, and \tilde{r}_{2k} is a term that only depends on the 2k-th moment of the gradient. It is notable that such an upper bound is independent of the Lipschitz constant. We then extend to the case where $\theta \ge \overline{\theta} > 1$ for some known constant $\overline{\theta}$. Moreover, when the privacy budget ϵ is small enough, we show an upper bound of $\tilde{O}\left(\left(\tilde{r}_k(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\epsilon}))^{\frac{k-1}{k}}\right)^{\frac{\theta}{\theta-1}}\right)$ even if the loss function is not Lipschitz. For the lower bound, we show that for any $\theta \ge 2$, the private minimax rate for ρ -zero Concentrated Differential Privacy is lower bounded by $\Omega\left(\left(\tilde{r}_k(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\sqrt{\rho}}))^{\frac{\theta}{\theta-1}}\right)^{\frac{\theta}{\theta-1}}\right)$.

029 030 031

032

006

008 009 010

011

013

014

015

016

017

018

021

023

024 025

026

027

028

1 INTRODUCTION

Machine learning is increasingly being integrated into daily life, driven by an ever-growing volume of data. This data often includes sensitive information, which raises significant privacy concerns. In response, regulations such as the GDPR mandate that machine learning algorithms not only effectively extract insights from training data but also uphold stringent privacy standards. Differential privacy (DP) (12), a robust framework for ensuring statistical data privacy, has garnered substantial attention recently and has emerged as the leading methodology for conducting privacy-preserving data analysis.

040 Differential Privacy Stochastic Convex Optimization (DP-SCO) and its empirical form, DP Empirical 041 Risk Minimization (DP-ERM), stand as core challenges within the machine learning and differential 042 privacy communities. These methodologies have been the focus of significant research over the past 043 decade, beginning with seminal works like those by Chaudhuri et al. (11) and followed by numerous 044 influential studies (6; 43; 38; 45; 24; 25; 32; 39; 40; 2). For instance, Bassily et al. (5) have provided 045 near-optimal rates for DP-SCO across both convex and strongly convex loss functions. Feldman et al. (14) have developed algorithms that boast linear time complexity, and Su et al. (34) have expanded 046 the discussion to non-Euclidean spaces. 047

However, the majority of existing theoretical frameworks primarily focus on scenarios where the loss function is O(1)-Lipschitz across all data, necessitating assumptions that the underlying data distribution is either bounded or sub-Gaussian. Such assumptions are crucial for the effectiveness of differential privacy methods based on output perturbation (11) and objective or gradient perturbation (6). Yet, these assumptions may not be valid for real-world datasets, particularly those from fields like biomedicine and finance, which are known to exhibit heavy-tailed distributions (44; 7; 19). This discrepancy can compromise the effectiveness of the algorithms in maintaining differential privacy. To bridge this gap, recent research has begun exploring DP-SCO in the context of heavytailed data, where the Lipschitz constant for the loss may be significantly higher or even unbounded (41; 21; 18; 29; 35). These studies typically assume that the gradient of the loss is bounded only in terms of its k-th moment for some k > 0, a much less stringent requirement than O(1)-Lipschitz continuity.

Although DP-SCO with heavy-tailed data has been extensively studied, most research has concen-060 trated on general convex or strongly convex functions. Yet, numerous other problems exist that 061 exceed the complexity of strongly convex functions or do not neatly fit within the convex-to-strongly 062 convex spectrum. In non-private settings, several studies have managed to achieve faster convergence 063 rates by introducing additional constraints on the loss functions. It has been demonstrated that it 064 is possible to exceed the convergence rates of general convex functions (49; 26; 37), and some 065 approaches have even matched the rates typical of strongly convex functions without the function actually being strongly convex (23; 27; 48). Similar advancements have been observed in the context 066 of privacy-preserving algorithms (3; 33). This leads to a compelling question: 067

For the problem of DP-SCO with heavy-tailed data and special classes of population risk functions, is it possible to achieve faster rates of excess population risk than the optimal ones of general convex and (or) strongly convex cases?

In this paper, we affirmatively respond by examining certain classes of population risk functions. Specifically, we focus on the case where the population risk function possesses a large or potentially infinite Lipschitz constant and meets the Tsybakov Noise Condition (TNC)¹, encompassing strongly convex functions, SVM, ℓ_1 -regularized stochastic optimization, and linear regression with heavytailed data as notable examples.

Our contributions are detailed as follows (refer to Table 1 for details).

- 1. We study DP-SCO where the population risk satisfies (θ, λ) -TNC with $\theta > 1$. Here, the loss function is L_f -Lipschitz, and the k-th moment of the loss gradient is small, where $L_f < \infty$ could be extremely large and $k \ge 2$. Based on our newly developed localization method, we propose an (ϵ, δ) -DP algorithm whose utility bound, with high probability, is $\tilde{O}((\tilde{r}_{2k}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\epsilon}))^{\frac{k-1}{k}})^{\frac{\theta}{\theta-1}})$ when $\theta \ge 2$. Here, n is the sample size, d is the model dimension and \tilde{r}_{2k} is a term that only depends on the 2k-th moment of the gradient. It is notable that such an upper bound is independent of the Lipschitz constant.
- 2. We further relax the assumption that $\theta \ge 2$ to $\theta \ge \overline{\theta} > 1$ for some known $\overline{\theta}$ and propose an algorithm that could achieve asymptotically the same bound as the previous one. Moreover, when the privacy budget ϵ is small enough, we show that even if the loss function is not Lipschitz, we can still get an upper bound of $\tilde{O}((\tilde{r}_k(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\epsilon}))^{\frac{k-1}{k}})^{\frac{\theta}{\theta-1}})$.
 - 3. On the lower bound side, for any $\theta \ge 2$, we show that there exists a population risk function satisfying TNC with parameter θ , whose minimax population risk under ρ -zero Concentrated Differential Privacy is always lower bounded by $\Omega((\tilde{r}_k(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{n\sqrt{\rho}}))^{\frac{k-1}{k}})^{\frac{\theta}{\theta-1}}).$

2 RELATED WORK

DP-SCO with Heavy-tailed Data. As we mentioned previously, there is a long list of work for DP-SCO from various perspectives. Here we only focus on the work related to DP-SCO with heavy-tailed 098 data. Generally speaking, there are two ways of modeling heavy-tailedness: The first one considers each coordinate of loss gradient has bounded moments, while the second one assumes the norm 100 of loss gradient has bounded moments, which is stronger than the first one. For the first direction, 101 (41) provides the first study under the assumption of bounded k-th moment $(k \ge 2)$ and proposes 102 three different ways for both convex and strongly convex loss. The bounds were later improved by 103 (20). Specifically, (20) provides improved upper bounds for convex loss and optimal rate for strongly 104 convex loss. Later, there are some works that consider different extensions. For example, (18) extends 105 to the high dimensional and polyhedral cases, (35) extends to the case where the gradient only has 106 (1+v)-th moment with $v \in (0,1]$, (42) considers the ℓ_1 -regression. For the second direction, (29)

078

079

081

082

084

085

090

092 093 094

¹⁰⁷

¹This is also referred to as the Error Bound Condition or the Growth Condition in related literature (27; 48).

109	Table 1: Comparion with previous results on DP-SCO with different assumptions in (ϵ, δ) -DP (we
110	always assume the loss is smooth). All results omit the term of $\log \frac{1}{\delta}$, smoothness and strong
111	convexity. † means the result is for ρ -zCDP. \star indicated the result holds when $\epsilon = \tilde{O}(\frac{1}{n})$.

	Upper Bound	Lower Bound	Assumption
(5)	$O\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{d}}{n\epsilon}\right)$	$\Omega\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{d}}{n\epsilon}\right)$	O(1)-Lipschitz
(5)	$O\left(\frac{1}{n} + \frac{d}{n^2 \epsilon^2}\right)$	$\Omega\left(\frac{1}{\sqrt{n}} + \frac{d}{n^2\epsilon^2}\right)$	O(1)-Lipschitz
(20)	$\tilde{O}\left(\frac{d}{\sqrt{n}} + \frac{d^2}{n\epsilon}\left(\frac{\epsilon n}{d^{\frac{3}{2}}}\right)^{\frac{1}{k}}\right)$	$\Omega(\sqrt{\frac{d}{n}} + \sqrt{d}(\frac{\sqrt{d}}{n\sqrt{\rho}})^{\frac{k-1}{k}})^{\dagger}$	$O(1)\text{-Lipschitz}$ and bounded $k\text{-th}$ moment ($k\geq 2)$
(20)	$\tilde{O}\left(\frac{d}{n} + d(\frac{\sqrt{d}}{n\epsilon})^{\frac{2(k-1)}{k}}\right)$	$\Omega\left(\frac{d}{n} + d\left(\frac{\sqrt{d}}{n\sqrt{\rho}}\right)^{\frac{2(k-1)}{k}}\right)^{\dagger}$	$O(1)\mbox{-Lipschitz}, strongly convex and bounded k\mbox{-th} moment (k\geq 2)$
(2; 33)	$\tilde{O}\left(\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{d}}{n\epsilon}\right)^{\frac{\sigma}{\sigma-1}}\right)$	$\Omega\left(\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{d}}{n\epsilon}\right)^{\frac{\theta}{\theta-1}}\right)$ when $\theta \ge 2$	$O(1)\text{-Lipschitz}$ under TNC with $\theta>1$
(29)	$O\left(\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)\right)$	$\Omega\left(\widetilde{r}_{k}\left(\frac{1}{\sqrt{n}}+\left(\frac{\sqrt{d}}{\sqrt{\rho n}}\right)^{\frac{k-1}{k}}\right)\right)^{\dagger}$	(large) Lipschitz, bounded $k\text{-th}$ moment ($k\geq 2)$
(29)	$\tilde{O}\left(\tilde{r}_{k}\left[\frac{1}{\sqrt{n}} + \max\left\{\left(\left(\frac{1}{\tilde{r}_{k}}\right)^{1/4}\frac{\sqrt{d}}{en}\right)^{\frac{4(k-1)}{5k-1}}, \left(\frac{\sqrt{d}}{en}\right)^{\frac{k-1}{k}}\right\}\right]\right)$	$\Omega\left(\widetilde{r}_k\left(\tfrac{1}{\sqrt{n}}+\left(\tfrac{\sqrt{d}}{\sqrt{\rho}n}\right)^{\frac{k-1}{k}}\right)\right)^\dagger$	bounded $k\text{-th}$ moment $(k\geq 2)$
(29)	$\tilde{O}\left(\tilde{R}^{2}_{2k,n}\left(\frac{1}{n} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{2(k-1)}{k}}\right)\right)$	$\Omega\left(\tilde{r}_{k}^{2}\left(\frac{1}{n}+\left(\frac{\sqrt{d}}{\sqrt{\rho}n}\right)^{\frac{2(k-1)}{k}}\right)\right)^{\dagger}$	strongly convex, bounded $k\text{-th}$ moment $(k\geq 2)$
This paper	$\tilde{O}\left(\left(\widetilde{R}_{2k,n}(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{k-1}{k}})\right)^{\frac{2}{\theta}}\right)$	$\Omega\left(\left(\tilde{r}_k\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)\right)^{\frac{\theta}{\theta-1}}\right) \text{ when } \theta \ge 2$	(large) Lipschitz function under TNC with $\theta>1$
This paper	$\tilde{O}\left(\tilde{r}_{k}^{\frac{\theta}{\theta-1}}\left(\frac{1}{\sqrt{n}}+\left(\frac{\sqrt{a}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)^{\frac{\theta}{\theta-1}}\right)^{\star}$	$\Omega\left(\left(\tilde{r}_k(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{k-1}{k}})\right)^{\frac{\theta}{\theta-1}}\right) \text{ when } \theta \geq 2$	TNC with $\theta > 1$
	(5) (20) (20) (2; 33) (29) (29) (29) (29) This paper This paper	$ \begin{array}{c c} \hline & \mbox{Upper Bound} \\ \hline (5) & O\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{2}}{n\epsilon}\right) \\ \hline (5) & O\left(\frac{1}{n} + \frac{\sqrt{2}}{n\epsilon}\right) \\ \hline (20) & \tilde{O}\left(\frac{d}{\sqrt{n}} + \frac{\sqrt{2}}{n\epsilon}(\frac{e\pi}{a^2})^{\frac{1}{k}}\right) \\ \hline (20) & \tilde{O}\left(\frac{d}{n} + d(\frac{\sqrt{2}}{a^2})^{\frac{1}{k-1}}\right) \\ \hline (23) & \tilde{O}\left(\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{2}}{n\epsilon}\right)^{\frac{p}{k-1}}\right) \\ \hline (29) & O\left(\tilde{R}_{2k,n}\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{2}}{e\pi}\right)^{\frac{p}{k-1}}\right) \\ \hline (29) & \tilde{O}\left(\tilde{r}_{k}\left[\frac{1}{\sqrt{n}} + \max\left\{\left(\left(\frac{1}{e_{k}}\right)^{\frac{1}{k}}, \frac{\sqrt{2}}{e\pi}\right)^{\frac{k-1}{k}}, \left(\frac{\sqrt{2}}{e\pi}\right)^{\frac{k-1}{k}}\right\}\right)\right) \\ \hline (29) & \tilde{O}\left(\tilde{R}_{2k,n}\left(\frac{1}{n} + \left(\frac{\sqrt{2}}{e\pi}\right)^{\frac{k-1}{k}}\right)\right) \\ \hline \\ \hline (29) & \tilde{O}\left(\tilde{R}_{2k,n}\left(\frac{1}{n} + \left(\frac{\sqrt{2}}{e\pi}\right)^{\frac{k-1}{k}}\right)\right) \\ \hline \\ \hline \\ This paper & \tilde{O}\left(\tilde{r}_{k}^{\frac{p}{k-1}}\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{2}}{e\pi}\right)^{\frac{k-1}{k}}\right)^{\frac{p}{p-1}}\right)^{\frac{k}{k}} \end{array} $	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $

108

127

provides a comprehensive study for both convex and strongly convex loss. In detail, for Lipschitz loss 128 whose gradient has k-th moment, they provide upper bounds that are independent of the Lipschitz 129 constant. Compared to (29), we first extend to the population risks that satisfy TNC (when $\theta = 2$, our 130 results match their results for strongly convex loss). Moreover, the results in (29) are in expectation 131 form while we provide new algorithms, and our results are in the high probability form. 132

DP for Heavy-tailed Data. In addition to DP-SCO, there is also some work on DP for heavy-tailed 133 data. (4) provided the first study on private mean estimation for distributions with the bounded 134 moment, which has been extended by (22; 28; 9) recently. However, these methods cannot be applied 135 to our problem as these results are all in the expectation form. Motivated by (41), we later consider 136 statistical guarantees of DP Expectation Maximization and applies to the Gaussian Mixture Model. 137 (46; 36; 47) considers private reinforcement learning and bandits learning where the reward follows a 138 heavy-tailed distribution. However, since the reward is a scalar, these methods are not applicable to 139 our problem. 140

Loss functions with TNC. While most of this paper focuses on loss functions that are either convex 141 or strongly convex, many loss functions fall between these two categories. That is, they are not 142 strongly convex, but their statistical rate is better than purely convex losses. For TNC and Lipschitz 143 loss functions, the best-known current rate is $O((\frac{1}{\sqrt{n}})^{\frac{\theta}{\theta-1}})$ (27), which corresponds to the first term 144 in our upper bounds. In Theorem 5, we demonstrate that this upper bound is tight for $\theta > 2$. A 145 comparison to the non-private setting will be included in the final version of the paper. 146

147 Our methods introduce novel technical challenges compared to non-private approaches. The key 148 innovation lies in our analysis, which is based on algorithmic stability and a newly developed 149 localized and clipped algorithm (Algorithm 3), which has not been previously studied. Specifically, Algorithm 4 is inspired by Algorithm 2 in (27). However, while the base algorithm in (27) is a 150 simple averaged version of projected SGD, our Algorithm 3 is significantly more complex. One 151 major technical challenge is that Algorithm 2 in (27) assumes a Lipschitz loss function with a fixed 152 Lipschitz constant. Consequently, their bounds rely on this constant. In contrast, we address scenarios 153 where the Lipschitz parameter can be extremely large. Therefore, we developed a new base algorithm 154 that removes dependence on this parameter and instead utilizes moments. 155

156 157

3 **PRELIMINARIES**

158 159

Definition 1 (Differential Privacy (12)). Given a data universe \mathcal{X} , we say that two datasets $S, S' \subseteq \mathcal{X}$ 160 are neighbors if they differ by only one entry, which is denoted as $S \sim S'$. A randomized algorithm 161 A is (ϵ, δ) -differentially private (DP) if for all neighboring datasets S, S' and for all events E in the 162 output space of \mathcal{A} , the following holds

 $\mathbb{P}(\mathcal{A}(S) \in E) \leqslant e^{\epsilon} \mathbb{P}(\mathcal{A}(S') \in E) + \delta.$

165 If $\delta = 0$, we call algorithm \mathcal{A} is ϵ -DP.

164

Definition 2 (zCDP (10)). A randomized algorithm \mathcal{A} is ρ -zero-concentrate-differentially private (zCDP) if for all neighboring datasets S, S' and $\alpha \in (1, \infty)$, we have $D_{\alpha}(\mathcal{A}(S) || \mathcal{A}(S')) \leq \rho \alpha$, where D_{α} is the α -Rényi divergence between $\mathcal{A}(S)$ and $\mathcal{A}(S')$.

Remark 1. In this paper, we focus on (ϵ, δ) -DP for upper bounds and ρ -zCDP for lower bounds, and we mainly use the Gaussian mechanism to guarantee the DP property. For Algorithms 1-5, which are based on stability analysis and the Gaussian mechanism, they operate as one-pass algorithms without sub-sampling. As a result, they can easily meet the requirements for CDP. However, a challenge arises with Algorithm 6. In this case, we employ privacy amplification via shuffling to reduce the noise. Currently, privacy amplification via shuffling is only applicable to ϵ and (ϵ, δ) -DP, and no version exists for zCDP. To maintain consistency throughout the paper, we use (ϵ, δ) -DP for all our upper bounds.

177 178 178 179 180 181 Definition 3 (Gaussian Mechanism). Given any function $q: \mathcal{X}^n \to \mathbb{R}^d$, the Gaussian mechanism is defined as $q(S) + \xi$ where $\xi \sim \mathcal{N}(0, \frac{16\Delta_2^2(q)\log(1/\delta)}{\epsilon^2}\mathbb{I}_d)$, where $\Delta_2(q)$ is the ℓ_2 -sensitivity of the function q, *i.e.*, $\Delta_2(q) = \sup_{S \sim S'} ||q(S) - q(S')||_2$. Gaussian mechanism preserves (ϵ, δ) -DP for $0 < \epsilon, \delta \leq 1$.

Definition 4 (DP-SCO (6)). Given a dataset $S = \{x_1, \dots, x_n\}$ from a data universe \mathcal{X} where x_i are i.i.d. samples from some unknown distribution \mathcal{D} , a convex loss function $f(\cdot, \cdot)$, and a convex constraint set $\mathcal{W} \subseteq \mathbb{R}^d$, Differentially Private Stochastic Convex Optimization (DP-SCO) is to find w^{priv} so as to minimize the population risk, *i.e.*, $F(w) = \mathbb{E}_{x \sim \mathcal{D}}[f(w, x)]$ with the guarantee of being differentially private. The utility of the algorithm is measured by the (*expected*) excess population *risk*, that is $\mathbb{E}\left[E(w^{\text{priv}})\right] = \min_{x \in \mathcal{D}} E(w)$

$$\mathbb{E}_{\mathcal{A}}[F(w^{\text{priv}})] - \min_{w \in \mathcal{W}} F(w),$$

where the expectation of \mathcal{A} is taken over all the randomness of the algorithm. Besides the population risk, we may also measure the *empirical risk* of dataset $S: \overline{F}(w, S) = \frac{1}{n} \sum_{i=1}^{n} f(w, x_i)$.

Definition 5 (Lipschitz). A function $f : W \mapsto \mathbb{R}$ is L-Lipschitz over the domain W if for all $w, w' \in \mathcal{W}, |f(w) - f(w')| \leq L ||w - w'||_2$.

193 **Definition 6** (Smoothness). A function $f : \mathcal{W} \mapsto \mathbb{R}$ is β -smooth over the domain \mathcal{W} if for all $w, w' \in \mathcal{W}, f(w) \leq f(w') + \langle \nabla f(w'), w - w' \rangle + \frac{\beta}{2} ||w - w'||_2^2$.

Definition 7 (Strongly Convex). A function $F : \mathcal{W} \mapsto \mathbb{R}$ is λ -strongly convex over the domain \mathcal{W} if, for all $w, w' \in \mathcal{W}, F(w) + \langle \nabla F(w), w' - w \rangle + \frac{\lambda}{2} ||w' - w||_2^2 \leq F(w')$.

Previous work on DP-SCO only focused on cases where the loss function is either convex or strongly convex (5; 14). In this paper, we mainly examine the case where the population risk satisfies the Tysbakov Noise Condition (TNC) (30; 27; 31), which has been extensively studied and has been shown that it could achieve faster rates than the optimal one of general convex loss functions in the non-private case. Below, we introduce the definition of TNC.

Definition 8 (Tysbakov Noise Condition). For a convex function $F(\cdot)$, let $\mathcal{W}_* = \arg\min_{w \in \mathcal{W}} F(w)$ denote the optimal set and for any $w \in \mathcal{W}$, let $w^* = \arg\min_{u \in \mathcal{W}_*} ||u - w||_2$ denote the projection of w onto the optimal set \mathcal{W}_* . The function F satisfies (θ, λ) -TNC for some $\theta > 1$ and $\lambda > 0$ if, for any $w \in \mathcal{W}$, the following inequality holds:

207 208

188

$$F(w) - F(w^*) \ge \lambda \|w - w^*\|_2^{\theta}.$$
 (1)

From the definition of TNC and Definition 7, we can see that a λ -strong convex function is $(2, \frac{\lambda}{2})$ -TNC. Moreover, if a function is (θ, λ) -TNC, then it is also (θ', λ) -TNC for any $\theta < \theta'$. Throughout the paper, we assume that θ is a constant and thus we omit the term of c^{θ} in the Big-O notation if c is a constant.

Lemma 1 (Lemma 2 in (30)). If the function $F(\cdot)$ is (θ, λ) -TNC and L_f -Lipschitz, then we have $\|w - w^*\|_2 \leq (L_f \lambda^{-1})^{\frac{1}{\theta-1}}$ and $F(w) - F(w^*) \leq (L_f^{\theta} \lambda^{-1})^{\frac{1}{\theta-1}}$ for all $w \in W$, where w^* is defined as in Definition 8. As mentioned earlier, our primary focus here is on cases where the loss function's Lipschitz constant is sufficiently large or even infinite. In such cases, we may seek alternative terms to replace the Lipschitz constant. Motivated by previous work on DP-SCO with heavy-tailed gradients, we consider the moments of the gradient. Specifically, we assume that the stochastic gradient distributions have bounded k-th moment for some $k \ge 2$:

Assumption 1. There exists $k \ge 2$ and $\tilde{r}^{(k)} > 0$ such that $\mathbb{E}\left[\sup_{w \in \mathcal{W}} \|\nabla f(w, x)\|_2^k\right] \le \tilde{r}^{(k)}$, where $\tilde{r}_k := (\tilde{r}^{(k)})^{1/k}$. Moreover, we assume the constrained set \mathcal{W} is bounded with diameter D.

If the loss function is L_f -Lipschitz, we can always observe that $\tilde{r}_k \leq L_f = \sup_{w,x} \|\nabla f(w,x)\|_2$. Moreover, \tilde{r}_k could be far less than the Lipschitz constant.

To state our subsequent theoretical results more clearly, we introduce some additional notations. For a batch of data $X \in \mathcal{X}^m$, we define the k-th empirical moment of $f(w, \cdot)$, by

$$\widehat{r}_m(X)^{(k)} = \sup_{w \in \mathcal{W}} \frac{1}{m} \sum_{i=1}^m \|\nabla f(w, x_i)\|_2^k$$

For $X \sim \mathcal{D}^m$, we denote the k-th expected empirical moment by

$$\widetilde{e}_m^{(k)} := \mathbb{E}[\widehat{r}_m(X)^{(k)}],$$

 $\widetilde{r}_{k m} := (\widetilde{e}_m^{(k)})^{1/k}.$

and let

236 237 238

221

222

223

233 234 235

Note that $\tilde{r}_{k,1} = \tilde{r}_k$. We define $\tilde{R}_{k,n} := \sqrt{\sum_{i=1}^l 2^{-i} \tilde{r}_{k,n_i}^2}$, where $n_i = 2^{-i}n$ and $l = \log_2 n$. Actually, $\tilde{R}_{k,n}$, a weighted average of the expected empirical moments for distinct batch sizes, is a constant used in the excess risk upper bounds, where we give more weight to \tilde{r}_m for large m. The following lemma indicates that it is smaller than \tilde{r}_k .

Lemma 2 ((29)). Under Assumption 1, we have: $\tilde{r}^{(k)} = \tilde{e}_1^k \ge \tilde{e}_2^k \ge \tilde{e}_4^k \ge \cdots \ge r^{(k)}$. Thus, we have $\tilde{R}_{k,n} \le \tilde{r}_k$.

245 246 247

248

243

244

4 LARGE LIPSCHITZ CONSTANT CASE

In this section, we will focus on the population risk function satisfying (θ, λ) -TNC, and the Lipschitz constant of the loss is extremely large (but finite). Before that, we first propose a novel localized noisy stochastic gradient method whose excess population risk is independent of the Lipschitz constant for general convex loss. See Algorithm 3 for details.

253 In Algorithm 3, we first partition the dataset into $O(\log_2 n)$ subsets where the *i*-th set has $O(2^{-i}n)$ 254 samples. In the *i*-th iteration, we use the *i*-th set and construct an ℓ_2 -regularized empirical risk 255 function F_i with hyperparameter λ_i in step 5. Moreover, based on the current model w_{i-1} , we 256 construct the constrained set \mathcal{W}_i with diameter exponential decay D_i . To handle large Lipschitz 257 constant and to solve the ℓ_2 -regularized empirical risk, we adopt a clipped gradient descent method 258 (Algorithm 2) with clip threshold C_i , where we use clipped gradients (Algorithm 1) to update our 259 model instead of the original gradient. After T_i iterations, we add Gaussian noise based on the stability of our clipped gradient descent to ensure (ϵ, δ) -DP. In the following we show Algorithm 3 260 could achieve a rate $\tilde{O}(\max\{\frac{1}{\sqrt{n}}, (\frac{d\log \frac{1}{\delta}}{\epsilon n})^{\frac{k-1}{k}}\})$ with specific parameters λ_i, T_i and C_i . 261 262

263
264Algorithm 1 ClippedMean($\{z_i\}_{i=1}^n, n, C$)265Input: $Z = \{z_i\}_{i=1}^n, C > 0$,266
2671: Compute $\widetilde{v} := \frac{1}{n} \sum_{i=1}^n \prod_C (z_i)$, where $\prod_C (z) := \operatorname{argmin}_{y \in \mathbb{B}_C} \|y - z\|_2^2$ denotes the projection268
269onto the ℓ_2 ball \mathbb{B}_C .269Return \widetilde{v}

270

Algorithm 2 Clipped Regularized Gradient Method

Input: Dataset $S \in \mathcal{X}^n$, iteration number T, stepsize η , clipping threshold C, regularization $\lambda \ge 0$, constraint set \mathcal{W} and initialization $w_0 \in \mathcal{W}$. 1: for all $t \in [T - 1]$ do $\nabla F_t(w_t) :=$ Clipped Mean $(\{\nabla f(w_t, x_i)\}_{i=1}^n; C)$ for gradients $\nabla f(w_t, x_i)$. 2: $w_{t+1} = \prod_{\mathcal{W}} [w_t - \eta (\nabla F_t(w_t) + \lambda (w_t - w_0))]$ 3: 276 4: **end for** 277 Return w_T

278 279

281 282

283 284

285 286

287 288

291

292

295

296

297

298 299

300

301 302

303

305 306

307

Theorem 1. Under Assumption 1, suppose that $f(\cdot, x)$ is α -smooth and L_f -Lipschitz with $L_f < \infty$ for every x. Then, for any $0 < \epsilon \leq \sqrt{\log(1/\delta)}, 0 < \delta < 1$ and $\eta_i \leq \frac{1}{\alpha}$ for all i, Algorithm 3 is (ϵ, δ) -DP. Let $p \ge 1$ such that $L_f \leqslant n^{p/2} \widetilde{R}_{2k,n}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d \log n}}{\epsilon n})^{\frac{k-1}{k}})$. For any $0 < \beta \leqslant \frac{1}{n}$, with probability at least $1 - \beta$, it holds that

$$F(w_l) - F(w^*) \leqslant \tilde{O}\left(\widetilde{R}_{2k,n}D(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{\epsilon n})^{\frac{k-1}{k}})\right),$$

where the Big- \tilde{O} notation omits all logarithmic terms (it is the same for other upper bounds).

289 **Remark 2.** Previous work on DP-SCO such as (43; 6), Lipschitz is still required for the loss function, though, it disappears in the final excess risk upper bound. And due to the property of worst-case stability and our assumption that L_f can be controlled by $n^{p/2} \widetilde{R}_{2k,n}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d \log n}}{\epsilon n})^{\frac{k-1}{k}})$ for sufficiently large p, we reach the upper bound with high probability without L_f in the final result. 293 Compared to (29), the main difference is that our result is in the high probability form while (29) is only in the expectation form. Specifically, to achieve a high probability result, instead of adding Gaussian noise to the gradient, we use the stability of the gradient descent. However, we cannot directly use the stability result in (17) here, which depends on the Lipschitz constant, making a large noise, we show that by using clipping, the stability now only depends on the clipping threshold.

Algorithm 3 Localized Noisy Clipped Gradient Method for DP-SCO(LNC-GM)(w_0, η, n, W)

- **Input:** Dataset $S \in \mathcal{X}^n$, stepsize η , clipping threshold $\{C_i\}_{i=1}^{\log_2 n}$, privacy parameter ϵ, δ , hyperparameter p, initialization $w_0 \in \mathcal{W}$.
- 1: Let $l = \log_2 n$.

304 2: for all $i \in [l]$ do

Set $n_i = 2^{-i}n, \eta_i = 4^{-i}\eta, \lambda_i = \frac{1}{\eta_i n_i^p}$ when $i \ge 2$, and $\lambda_1 = \frac{1}{\eta_i n_i^{2p}}, T_i = \frac{1}{\eta_i n_i^{2p}}$ 3: $\widetilde{\Theta}\left(\frac{1}{\lambda_i\eta_i}\right)$, and $D_i = \frac{2L_f}{\lambda_i}$. Draw a new batch \mathcal{B}_i of $n_i = |\mathcal{B}_i|$ samples from S without replacement.

4: 308

5: Denote
$$F_i(w) := \frac{1}{n_i} \sum_{j \in \mathcal{B}_i} f(w, x_j) + \frac{\lambda_i}{2} \|w - w_{i-1}\|^2$$
.

Use Algorithm 2 with initialization w_{i-1} to minimize \widehat{F}_i over $\mathcal{W}_i := \{w \in \mathcal{W} | \|w - w_{i-1}\| \leq D_i\}$ for T_i iterations with clipping threshold $C_i = \widetilde{r}_{2k,n_i} (\frac{\epsilon n_i}{\sqrt{d \log(n)}})^{1/k}$ and stepsize η_i . Let 310 6: 311 312 \hat{w}_i be the output of Algorithm 2

313
314 7: Set
$$\xi_i \sim \mathcal{N}(0, \sigma_i^2 \mathbb{I}_d)$$
 where $\sigma_i = \frac{8C_i \sqrt{\log}}{n_i \lambda_i \epsilon}$

- Set $w_i = \widehat{w}_i + \xi_i$. 8: 315
- 9: end for 316
- 10: **Return** the final iterate w_l 317

318

Theorem 2. Under Assumption 1 and suppose that the population risk function $F(\cdot)$ is (θ, λ) -TNC 323 with $\theta \geq 2$, and $f(\cdot, x)$ is α -smooth and L_f -Lipschitz for each x. Additionally, take $p \geq 1$ such that

Based on our novel locality algorithm, we then apply it to TNC functions. See Algorithm 4 for details. 319 Specifically, we partition the dataset into several subsets of equal size. As the iteration number 320 increases, we consider a constrained set centered at the current parameter with a smaller diameter 321 and learning rate in Algorithm 3. 322

324 Algorithm 4 Private Stochastic Approximation (w_1, n, R_0) 325 **Input:** Dataset $S \in \mathcal{X}$, initial point $w_1 \in \mathcal{W}$, privacy parameter ϵ and δ , hyperparameter p, initial 326 diameter R_0 . 1: Set $\hat{w}_0 = w_1$, $m = \lfloor \frac{1}{2} \log_2 \frac{2n}{\log_2 n} \rfloor - 1$, $n_0 = \lfloor \frac{n}{m} \rfloor$. Then partition the dateset S into m disjoint 327 328 subsets, namely, $\{S_1, \dots, S_m\}$ with each $|S_i| = n_0$. 2: for all $l \in [m]$ do Set $\gamma_l = \frac{R_{l-1}}{n_0^2} \min\{\frac{1}{L_f}, \frac{1}{\tilde{R}_{2k,n}n_0^{\frac{p+1}{2}}} (\frac{\epsilon n_0}{\sqrt{d\log n}})^{\frac{k-1}{k}}, \frac{1}{n_0^{\frac{p-1}{2}}L_f^2 \sqrt{\log n_0 \log(1/\beta)}}\}$ and $R_l = \frac{R_{l-1}}{2}$. Denote $\hat{w}_l = \text{LNC-GM}(\hat{w}_{l-1}, \gamma_l, n_l, \mathcal{W})$, and constrained set $\mathcal{W} \cap \mathbb{B}(\hat{w}_{l-1}, R_{l-1})$. 330 3: 331 332 4: 333 5: end for 334 **Return** \hat{w}_m 335

336 337 338

358

359

 $L_f \leq n^{p/2} \widetilde{R}_{2k,n} \left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d \log n}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)$, then algorithm 4 is (ϵ, δ) -DP. Moreover, for sufficiently large n such that $\gamma_l \leq \frac{1}{\alpha}$, with probability at least $1 - \beta$, we have

$$F(\hat{w}_m) - F(w^*) \leqslant \tilde{O}\left(\frac{1}{\lambda^{\frac{1}{\theta-1}}} (\tilde{R}_{2k,n}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{\epsilon n})^{\frac{k-1}{k}}))^{\frac{\theta}{\theta-1}}\right).$$

We note that there is no dependence on p in the final bound in Theorem 1 and 2. p is used to control the Lipschitz constant thus we can remove the Lipschitz constant from the final bound. We can see that in the proof of Theorem 1, there exists a term with n^p both in the numerators and denominators. By assuming that L_f is controlled by the $O(n^p/2)$ and choosing specific η , we can eliminate the p in the final bound. A similar result holds for Theorem 2.

Remark 3. In the case of O(1)-Lipschitz loss under TNC, compared with the optimal rate $\Theta(((\frac{1}{\sqrt{n}} +$

 $\frac{(\sqrt{d})}{(\epsilon_n)} \frac{k-1}{k} = 0$ in (3), our improved result gets rid of the dependence of Lipschitz constant, which could be extremely large. Moreover, when $\theta = 2$, i.e., the population risk is strongly convex, our result covers the result in (29). Thus, our result is a generalized upper bound. It is also notable that our upper bound is independent of the diameter of the constrained set and the Lipschitz-smoothness parameter. In Algorithm 4, one need the projection onto the ball $W \cap \mathbb{B}(\hat{w}_{l-1}, R_{l-1})$ in each iteration of the Phased-SGD in each phase. This could be solved using Dykstra's algorithm (13; 8).

Example. We consider the ℓ_1 constrained ℓ_4 -norm linear regression, which has been studied in (48) and satisfies TNC with $\theta = 4$ (27). Specifically, it can be written as the following.

$$\min_{\|w\|_1 \leqslant 1} F(w) \stackrel{\Delta}{=} \mathbb{E}[(\langle w, x \rangle - y)^4].$$
⁽²⁾

When y is bounded by O(1) and x follows a truncated normal Gaussian distribution at $[-n, n]^d$. Then we can see that the loss function is Poly(n)-Lispchitz, but its 2k-th moment is O(1). In this case, our bound in equation 2 is much smaller than the previous results in (3; 33).

So far, we have proposed an algorithm for TNC. Nevertheless, we also find that Theorem 2 needs a strong assumption on θ , i.e. $\theta \ge 2$. Thus, a direct question that occurs to us is whether we can further improve the upper bound. To conquer the disadvantage of the above algorithm, we propose the following. We assume θ is unknown but bigger than some definite $\bar{\theta} > 1$. Then we divide the whole dataset into subsets with distinct elements, detailly $l = \lfloor (\log_{\bar{\theta}} 2) \cdot \log_{\bar{\theta}} n \rfloor$ with $n_i = \lfloor 2^{i-1}n/(\log n)^{\log_{\bar{\theta}}^2 2} \rfloor$ samples for each subset. Then we run the Algorithm 1 for *l* times while each phase implements on the *i*-th subset and is initialized at the output of the previous one.

Theorem 3. Under Assumption 1 and assume that the loss function $F(\cdot)$ satisfies (θ, λ) -TNC with parameter $\theta \ge \overline{\theta} > 1$ for some definite constant $\overline{\theta}$, and $f(\cdot, x)$ is convex, α smooth and L_f -Lipschitz for each x. Algorithm 5 is (ϵ, δ) -DP for any $\epsilon \le 2\log(1/\delta)$, and take $p \ge 1$ such that $L_f \le n^{p/2} \widetilde{R}_{2k,n}(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d\log n}}{\epsilon n})^{\frac{k-1}{k}})$. Moreover, if the sample size n is sufficiently large such that $\overline{\theta} \ge 2^{\frac{\log\log n}{\log n-1}}$ and $\eta_t \le \frac{1}{\alpha}$, we have with probability at least $1 - \beta$

377
$$F(w_l) - F(w^*) \leq \tilde{O}\left(\left(\frac{1}{\lambda}\right)^{\frac{1}{\theta-1}} \left(\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)\right)^{\frac{\theta}{\theta-1}}\right),$$

378	Algorithm 5 Iterated Localized Noisy Clipped Gradient Method
379	Input: Dataset $S \in \mathcal{X}^n$ initial point $w_0 \in \mathcal{W}$ privacy parameter ϵ and δ parameter n initial
380	diameter R_0 .
381	1: Partite the data S into l disjoint subsets $\{S_1, \dots, S_l\}$, where $l = \lfloor (\log_{\bar{\theta}} 2) \cdot \log \log n \rfloor$ and for
382 383	each $i \in [l], S_i = n_i = \left 2^{i-1} n / (\log n)^{\log_{\overline{\theta}}^2 2} \right .$
384	2: for all $t = 1, \cdots, l$ do
385	3: Let w_t = Algorithm 3 with input $(S_i, w_{t-1}, \eta_t, \mathcal{W})$, where η_t =
386	$\frac{R_{t-1}}{\frac{p}{2}}\min\{\frac{1}{L_f}, \frac{1}{\tilde{\mu}_{t-1}}, \frac{en_i}{\frac{p+1}{2}}(\frac{\epsilon n_i}{\sqrt{d\log n}})^{\frac{k-1}{k}}, \frac{1}{\frac{p-1}{2}L^2\sqrt{\epsilon n}} + \frac{1}{\ln \frac{p}{2}}\} \text{ and } R_l = \frac{R_{l-1}}{2}.$
387	$n_0 \qquad R_{2k,n}n_i \stackrel{*}{\longrightarrow} V \stackrel{*}{\longrightarrow} n_i \stackrel{*}{\longrightarrow} L_{\widehat{fV}} \log n_i \log(1/\beta)$ 4. end for
388	Return w_k
000	

Remark 4. We pause to have another glimpse of Algorithm 4 and Algorithm 5. Note that they have a similar procedure to take the dataset apart, while the number of each subset is the same in Algorithm 5 and increases in Algorithm 5 as the iteration grows. And the set we project on also varies between Algorithm 4 and 5.

LOWER BOUNDS

In this section, we will show that the above upper bounds is nearly optimal (if \tilde{r}_{2k} and \tilde{r}_k are asymptotically the same) by providing lower bounds of the private minimax rate for ρ -zCDP. Specifically, for a sample space $\mathcal{X} \subseteq \mathbb{R}^d$ and collection of distributions \mathcal{P} over \mathcal{X} , we define the function class $\mathcal{F}_k^{\theta}(\mathcal{P}, \tilde{r}^{(k)})$ as the set of population risk functions from $\mathbb{R}^d \mapsto \mathbb{R}$ that satisfy $(\theta, 1)$ -TNC and their loss satsifies Assumption 1. We define the constrained minimax risk

$$\mathcal{M}(\mathcal{W}, \mathcal{P}, \mathcal{F}_k^{\theta}(\mathcal{P}, \tilde{r}^{(k)}), \rho) = \inf_{\mathcal{A} \in \mathcal{Q}(\rho)} \max_{F \times P \in \mathcal{F}_k^{\theta}(\mathcal{P}, \tilde{r}_k) \times \mathcal{P}} \mathbb{E}_{\mathcal{A}, D \in P^n}[F(\mathcal{A}(D)) - \min_{w \in \mathcal{W}} F(w)],$$

where $Q(\rho)$ is the set of all ρ -zCDP algorithms. We will show the following two results for different sample spaces and constraint sets.

Theorem 4. For any $\theta, k \geq 2$, $\tilde{r}^{(k)} > 0$, denote $\mathcal{X} = \{\pm p^{-\frac{1}{k}} \frac{\tilde{r}_k}{2\sqrt{d}}\}^d \cup \{0\}$ with $\tilde{r}_k = (\tilde{r}^{(k)})^{\frac{1}{k}}$, and $\mathcal{W} = \mathbb{B}_r$ with $r = \left(\frac{p^{-\frac{1}{k}}\tilde{r}_k}{2}\right)^{\frac{1}{\theta-1}}$ and $p = \frac{d}{n\sqrt{\rho}}$, then, if n is large enough such as $n \ge \Omega(\frac{\sqrt{d}}{\sqrt{\rho}})$, we have the following lower bound

 $\mathcal{M}(\mathcal{W}, \mathcal{P}, \mathcal{F}_k^{\theta}(\mathcal{P}, \tilde{r}^{(k)}), \rho) \ge \Omega\left((\tilde{r}_k((\frac{\sqrt{d}}{\sqrt{\rho}n})^{\frac{k-1}{k}}))^{\frac{\theta}{\theta-1}} \right).$

Theorem 5. For any $\theta, k \geq 2, \tilde{r}_k > 0$, denote $\mathcal{X} = \{\pm \frac{\tilde{r}_k}{2\sqrt{d}}\}^d$, and $\mathcal{W} = \mathbb{B}_r$ with $r = (\frac{\tilde{r}_k}{2})^{\frac{1}{\theta-1}}$, then, if $n \geq \Omega(\sqrt{d})$, we have the following lower bound

$$\mathcal{M}(\mathcal{W}, \mathcal{P}, \mathcal{F}_k^{\theta}(\mathcal{P}, \tilde{r}^{(k)}), \rho) \ge \Omega\left((\frac{\tilde{r}_k}{\sqrt{n}})^{\frac{\theta}{\theta-1}} \right)$$

Remark 5. First, it is notable that although the upper bounds in Section 4 are for (ϵ, δ) -DP, we can easily extend to the ρ -zCDP case as we used the Gaussian mechanism and parallel theorem to guarantee DP, which also hold for zCDP (10). The only difference is replacing the term $O(\frac{\sqrt{\log \frac{1}{\delta}}}{\epsilon})$ by $O(\frac{1}{\sqrt{a}})$. Thus, from this side, combining with Theorem 4 and 5, we can see the upper bound is nearly optimal for ρ -zCDP in the general case if \tilde{r}_{2k} (since $\tilde{R}_{2k,n} \leq \tilde{r}_{2k}$) and \tilde{r}_k are asymptotically the same. Secondly, in the Lipschitz case for (ϵ, δ) -DP, (2) proved the lower bound result via a reduction to the ERM problem for general convex loss. However, their reduction cannot be applied to our problem as their proof heavily relies on the O(1)-Lipschizt condition, which is not satisfied for our loss. For ϵ -DP, (2) considered the empirical risk and used the packing argument for the lower bound, which cannot be applied to our problem as our loss is not constant Lipschitz. In our proof, we directly considered the population risk $F_P(w) = -\langle w, \mathbb{E}_P[x] \rangle + \frac{1}{d} \|w\|_{\ell}^{\theta}$ for some data distribution P and used private Fano's lemma to prove the lower bound.

6 RELAX THE LIPSCHITZ ASSUMPTION

> private paratemeter ϵ, δ , initialization w_0 . 1: Randomly permute the data and denote the permuted data as $\{x_1, \dots, x_n\}$. 2: Initialize $w_0^{ag} = w_0$. 3: for all $t \in [T]$ do 4: $w_t^{md} := (1 - \alpha_t) w_{t-1}^{ag} + \alpha_t w_{t-1}$. 5: Draw new batch \mathcal{B}_t (without replacement) of n/T samples from S. 6: $\widetilde{\nabla}F_t(w_t^{md}) := \text{ClippedMean}\left(\{\nabla f(w_t^{md}, x)\}_{x \in \mathcal{B}_t}; \frac{n}{T}; C\right) + \zeta_i$, where $\zeta_i \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_d)$, $\sigma^2 = O(\frac{C^2 T \log \frac{1}{\delta}}{n^2 \epsilon^2})$ and $C = \widetilde{r}_k \left(\frac{\epsilon n}{\sqrt{d \log(1/\delta)}}\right)^{1/k}$. 7: $w_t := \arg_{w \in \mathcal{W}} \left\{\alpha_t \left\langle \widetilde{\nabla}F_t(w_t^{md}), w \right\rangle + \frac{\eta_t}{2} \|w_{t-1} - w\|^2 \right\}$. 8: $w_t^{ag} := \alpha_t w_t + (1 - \alpha_t) w_{t-1}^{ag}$.

Algorithm 6 Permuted Noisy Clipped Accelerated SGD for Heavy-Tailed DP SCO (PNCA-SGD)

Input: Data $S \in \mathcal{X}^n$, iteration number T, stepsize parameters $\{\eta_t\}_{t \in [T]}, \{\alpha_t\}_{t \in [T]}$ with $\alpha_1 = 1$,

Return w_T^{ag}

 In the previous section, we have considered the Lipschitz case and show that under the TNC, compared to that for the general convex loss, it is possible to get improved excess population risk that is independent of the Lipschitz constant. There are still two questions left: (1) Compared to the previous studies on DP-SCO with heavy-tailed gradient such as (41; 20), our above upper bounds still need the finite Lipschitz condition; (2) We can see our upper bounds depend on $\tilde{R}_{2k,n} \leq \tilde{r}_{2k}$ while the lower bounds only depend on \tilde{r}_k . Thus, there is a gap for the moment term. In this section, we aim to address these two issues. Specifically, we will show that even if the loss function is not Lipschitz, it is still possible to get the same upper bound as in the above section when ϵ is small enough. Moreover, we can improve the dependency from $\tilde{R}_{2k,n}$ to \tilde{r}_k .

Specifically, our main method, Algorithm 7, shares a similar idea as in Algorithm 5 with different parameters and base algorithm. Specifically, rather than using Algorithm 3, here we propose Algorithm 6 as our base algorithm, which is a shuffled, clipped, and private version of the accelerated SGD. Specifically, in step 1 we randomly shuffle the data for privacy amplification (15). Then, in each iteration, we clipped the gradients and added Gaussian noise to ensure DP. We can show that with some parameters, the output could achieve an upper bound similar to Theorem 1 even if the loss is not Lipschitz.

Algorithm 7 Iterated PNCA-SGD $(w_0, n, W, \bar{\theta})$ **Input:** Dataset $S \in \mathcal{X}^n$, initial point $w_0 \in \mathcal{W}$, privacy parameter ϵ and δ . 1: Partite the data S into k disjoint subsets $\{S_1, \dots, S_k\}$, where $k = \lfloor (\log_{\bar{\theta}} 2) \cdot \log \log n \rfloor$, and for each $i \in [k], |S_i| = n_i = \left| 2^{i-1} n / (\log n)^{\log_{\tilde{\theta}}^2 2} \right|.$ 2: for all $t = 1, \cdots, k$ do Let $w_t = \text{PNCA-SGD}(w_{t-1}, \eta_t, n_t, \mathcal{W})$, where the AC-SA runs on the t-th subset S_i .For (ϵ, δ) -DP, $\eta_t = \frac{4\eta}{t(t+1)}$, $\alpha_t = \frac{2}{t+2}$ and $R_l = \frac{R_{l-1}}{2}$. 4: end for **Return** w_k

Theorem 6. For any $\epsilon = O(\sqrt{\frac{\log n/\delta}{n}})$, and $0 < \delta < 1$, Algorithm 6 is (ϵ, δ) -DP. Moreover, under Assumption 1 and assume function F is β -smooth with the diameter D over $w \in W$, then the output of Algorithm 6, by selecting the following T,

$$T = \min\{\sqrt{\frac{\beta D}{\widetilde{r}_k}} \cdot (\frac{\epsilon n}{\sqrt{d\log(1/\delta)}})^{\frac{k-1}{2k}}, \sqrt{\frac{\beta D}{\widetilde{r}_k}} \cdot n^{1/4}\},$$

we have

$$\mathbb{E}F\left(w_T^{ag}\right) - F^* \leqslant O\left(\widetilde{r}_k D\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d\log(1/\delta)}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)\right).$$

Note that (29) also proposes a private accelerated SGD. However, their bound is sub-optimal (see the second row is Table 1). Here, we leverage privacy amplification via shuffling to reduce the noise added to each iteration. Thus, we can get the optimal rate here. We note that this is also the first result that can achieve the optimal rate for the general convex function without the Lipschitz assumption. Based on this result, we have the following theorem for Algorithm 7.

Theorem 7. For any $\epsilon = O(\sqrt{\frac{\log n/\delta}{n}})$, and $0 < \delta < 1$, Algorithm 7 is (ϵ, δ) -DP. Moreover, under Assumption 1 and assume function F is β -smooth, then we have

$$\mathbb{E}F(\hat{w}_m) - F(w^*) \leqslant \tilde{O}\left(\frac{1}{\lambda^{\frac{1}{\theta-1}}} (\tilde{r}_k(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d}}{\epsilon n})^{\frac{k-1}{k}}))^{\frac{\theta}{\theta-1}}\right).$$
(3)

Compared with the results in the above section, we can see the result in Theorem 7 is in the expectation form, which is due to the noisy clipped gradient in Algorithm 6. Moreover, the constraint of $\epsilon = O(\sqrt{\frac{\log n/\delta}{n}})$ comes from the results of privacy amplification via shuffling (15). We leave these two assumptions to be relaxed for future research. Moreover, the improvement from $\tilde{R}_{2k,n}$ to \tilde{r}_k is due to the different results between Theorem 6 and 1.

540 REFERENCES

- [1] Jayadev Acharya, Ziteng Sun, and Huanyu Zhang. Differentially private assouad, fano, and le cam. In *Algorithmic Learning Theory*, pages 48–78. PMLR, 2021.
 - [2] Hilal Asi, John Duchi, Alireza Fallah, Omid Javidbakht, and Kunal Talwar. Private adaptive gradient methods for convex optimization. In *International Conference on Machine Learning*, pages 383–392. PMLR, 2021.
- [3] Hilal Asi, Daniel Lévy, and John C Duchi. Adapting to function difficulty and growth conditions in private optimization. *Advances in Neural Information Processing Systems*, 34:19069–19081, 2021.
 - [4] Rina Foygel Barber and John C. Duchi. Privacy and statistical risk: Formalisms and minimax bounds. *arXiv preprint arXiv:1412.4451*, 2014.
 - [5] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Thakurta. Private stochastic convex optimization with optimal rates. *arXiv preprint arXiv:1908.09970*, 2019.
 - [6] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In 2014 IEEE 55th Annual Symposium on Foundations of Computer Science, pages 464–473. IEEE, 2014.
- [7] Atanu Biswas, Sujay Datta, Jason P Fine, and Mark R Segal. *Statistical advances in the biomedical science*. Wiley Online Library, 2007.
- [8] James P Boyle and Richard L Dykstra. A method for finding projections onto the intersection of convex sets in hilbert spaces. In *Advances in order restricted statistical inference*, pages 28–47. Springer, 1986.
- [9] Victor-Emmanuel Brunel and Marco Avella-Medina. Propose, test, release: Differentially private estimation with high probability. *arXiv preprint arXiv:2002.08774*, 2020.
- [10] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer, 2016.
- [11] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- [12] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [13] Richard L Dykstra. An algorithm for restricted least squares regression. Journal of the American Statistical Association, 78(384):837–842, 1983.
- [14] Vitaly Feldman, Tomer Koren, and Kunal Talwar. Private stochastic convex optimization: optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 439–449, 2020.
- [15] Vitaly Feldman, Audra McMillan, and Kunal Talwar. Hiding among the clones: A simple and nearly optimal analysis of privacy amplification by shuffling. In 2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS), pages 954–964. IEEE, 2022.
- [16] Vitaly Feldman and Jan Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279. PMLR, 2019.
- [17] Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. arXiv e-prints, pages arXiv–1509, 2015.
- [18] Lijie Hu, Shuo Ni, Hanshen Xiao, and Di Wang. High dimensional differentially private
 stochastic optimization with heavy-tailed data. In *Proceedings of the 41st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages 227–236, 2022.

598

600

601

602

603

604

605

606

607

608 609

610

611

615

616

617

618

619

620

621

622 623

624

625

626

627

628

629 630

631

632 633

634

635

636

637

638 639

640

641

642

643

644

- [19] Marat Ibragimov, Rustam Ibragimov, and Johan Walden. *Heavy-tailed distributions and robustness in economics and finance*, volume 214. Springer, 2015.
 - [20] Gautam Kamath, Xingtu Liu, and Huanyu Zhang. Improved rates for differentially private stochastic convex optimization with heavy-tailed data. arXiv preprint arXiv:2106.01336, 2021.
 - [21] Gautam Kamath, Xingtu Liu, and Huanyu Zhang. Improved rates for differentially private stochastic convex optimization with heavy-tailed data. In *International Conference on Machine Learning*, pages 10633–10660. PMLR, 2022.
 - [22] Gautam Kamath, Vikrant Singhal, and Jonathan Ullman. Private mean estimation of heavytailed distributions. In *Proceedings of 33rd Conference on Learning Theory (COLT)*, pages 2204–2235, 2020.
 - [23] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximalgradient methods under the polyak-lojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.
 - [24] Shiva Prasad Kasiviswanathan and Hongxia Jin. Efficient private empirical risk minimization for high-dimensional learning. In *International Conference on Machine Learning*, pages 488–497, 2016.
- [25] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pages 25–1, 2012.
 - [26] Tomer Koren and Kfir Y Levy. Fast rates for exp-concave empirical risk minimization. In NIPS, pages 1477–1485, 2015.
 - [27] Mingrui Liu, Xiaoxuan Zhang, Lijun Zhang, Rong Jin, and Tianbao Yang. Fast rates of erm and stochastic approximation: Adaptive to error bound conditions. *arXiv preprint arXiv:1805.04577*, 2018.
 - [28] Xiyang Liu, Weihao Kong, Sham Kakade, and Sewoong Oh. Robust and differentially private mean estimation. *arXiv preprint arXiv:2102.09159*, 2021.
 - [29] Andrew Lowy and Meisam Razaviyayn. Private stochastic optimization with large worst-case lipschitz parameter: Optimal rates for (non-smooth) convex losses and extension to non-convex losses. In *International Conference on Algorithmic Learning Theory*, pages 986–1054. PMLR, 2023.
 - [30] Aaditya Ramdas and Aarti Singh. Optimal rates for first-order stochastic convex optimization under tsybakov noise condition. *arXiv preprint arXiv:1207.3012*, 2012.
 - [31] Aaditya Ramdas and Aarti Singh. Algorithmic connections between active learning and stochastic convex optimization. In *International Conference on Algorithmic Learning Theory*, pages 339–353. Springer, 2013.
 - [32] Adam Smith, Abhradeep Thakurta, and Jalaj Upadhyay. Is interaction necessary for distributed private learning? In 2017 IEEE Symposium on Security and Privacy (SP), pages 58–77. IEEE, 2017.
 - [33] Jinyan Su and D Wang. Faster rates of differentially private stochastic convex optimization. *arXiv preprint arXiv*, 2108, 2021.
 - [34] Jinyan Su, Changhong Zhao, and Di Wang. Differentially private stochastic convex optimization in (non)-euclidean space revisited. In *Uncertainty in Artificial Intelligence*, pages 2026–2035. PMLR, 2023.
 - [35] Youming Tao, Yulian Wu, Xiuzhen Cheng, and Di Wang 0015. Private stochastic convex optimization and sparse learning with heavy-tailed data revisited. In *IJCAI*, pages 3947–3953, 2022.
- [36] Youming Tao, Yulian Wu, Peng Zhao, and Di Wang. Optimal rates of (locally) differentially private heavy-tailed multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 1546–1574. PMLR, 2022.

- [37] Tim van Erven, Peter D Grünwald, Nishant A Mehta, Mark D Reid, and Robert C Williamson.
 Fast rates in statistical and online learning. *Journal of Machine Learning Research*, 16:1793–1861, 2015.
 - [38] Di Wang, Changyou Chen, and Jinhui Xu. Differentially private empirical risk minimization with non-convex loss functions. In *International Conference on Machine Learning*, pages 6526–6535, 2019.
 - [39] Di Wang, Marco Gaboardi, and Jinhui Xu. Empirical risk minimization in non-interactive local differential privacy revisited. In Advances in Neural Information Processing Systems, pages 965–974, 2018.
 - [40] Di Wang, Adam Smith, and Jinhui Xu. Noninteractive locally private learning of linear models via polynomial approximations. In *Algorithmic Learning Theory*, pages 897–902, 2019.
 - [41] Di Wang, Hanshen Xiao, Srinivas Devadas, and Jinhui Xu. On differentially private stochastic convex optimization with heavy-tailed data. In *International Conference on Machine Learning*, pages 10081–10091. PMLR, 2020.
 - [42] Di Wang and Jinhui Xu. Differentially private 1-norm linear regression with heavy-tailed data. In 2022 IEEE International Symposium on Information Theory (ISIT), pages 1856–1861. IEEE, 2022.
 - [43] Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. In Advances in Neural Information Processing Systems, pages 2722– 2731, 2017.
 - [44] Robert F Woolson and William R Clarke. *Statistical methods for the analysis of biomedical data*, volume 371. John Wiley & Sons, 2011.
 - [45] Xi Wu, Fengan Li, Arun Kumar, Kamalika Chaudhuri, Somesh Jha, and Jeffrey Naughton. Bolton differential privacy for scalable stochastic gradient descent-based analytics. In *Proceedings* of the 2017 ACM International Conference on Management of Data, pages 1307–1322. ACM, 2017.
 - [46] Yulian Wu, Xingyu Zhou, Sayak Ray Chowdhury, and Di Wang. Differentially private episodic reinforcement learning with heavy-tailed rewards. In *International Conference on Machine Learning*, pages 37880–37918. PMLR, 2023.
 - [47] Yulian Wu, Xingyu Zhou, Youming Tao, and Di Wang. On private and robust bandits. *Advances in Neural Information Processing Systems*, 36, 2024.
 - [48] Yi Xu, Qihang Lin, and Tianbao Yang. Stochastic convex optimization: Faster local growth implies faster global convergence. In *International Conference on Machine Learning*, pages 3821–3830. PMLR, 2017.
 - [49] Tianbao Yang, Zhe Li, and Lijun Zhang. A simple analysis for exp-concave empirical minimization with arbitrary convex regularizer. In *International Conference on Artificial Intelligence and Statistics*, pages 445–453. PMLR, 2018.

702 A OMITTED PROOFS

Proof of Lemma 2. Let $l \in \mathbb{N}$, and $n = 2^{l}$ and consider

$$\widehat{r}_{n}(X)^{(k)} = \frac{1}{n} \sup_{w} \left(\sum_{i=1}^{n/2} \|\nabla f(w, x_{i})\|^{k} + \sum_{i=n/2+1}^{n} \|\nabla f(w, x_{i})\|^{k} \right)$$
$$\leqslant \frac{1}{n} \left(\sup_{w} \sum_{i=1}^{n/2} \|\nabla f(w, x_{i})\|^{k} + \sup_{w} \sum_{i=n/2+1}^{n} \|\nabla f(w, x_{i})\|^{k} \right)$$

Taking expectations over the random draw of $X \sim \mathcal{D}^n$ and we have $\tilde{e}_n^{(k)} \leq \tilde{e}_{n/2}^{(k)}$. Thus, $\tilde{R}_{k,n} \leq \tilde{r}_k$.

Proof of Theorem 1. Privacy. Since in each epoch of Algorithm 3 we use a disjoint dataset, it is sufficient for us to show each w_i is (ϵ, δ) -DP. Recall the following lemma:

Lemma 3. [Theorem 3.9 in ((17))] Assume the loss function $f(\cdot, x)$ is L-Lipschitz, λ -strongly convex and β -smooth with respect to $w \in W$ for all x. Let S_i and S'_i be two samples of size n_i differing in only a single element. Denote w_i^t and w'_i as the outputs of the projected stochastic gradient method on datasets S_i and S'_i respectively, then if $\eta \leq \frac{1}{\beta}$ we have

$$\left\| w_i^t - w_i'^t \right\| \leqslant \frac{2L}{\lambda n_i}$$

We can see F_i is λ_i -strongly convex and the gradient after clipping $\Pi_C(\nabla f(\cdot, x))$ is also α -Lipschitz if $f(\cdot, x)$ is α -smooth, which is due to that

$$||\Pi_C(\nabla f(w,x)) - \Pi_C(\nabla f(w',x))| \le ||\nabla f(w,x) - \nabla f(w',x)||_2 \le \alpha ||w - w'||_2.$$
(4)

Thus, by a similar proof for Lemma 3, we can see that $\|\hat{w}_i - \hat{w}'_i\|_2 \leq \frac{2C_i}{n_i\lambda_i}$, where \hat{w}'_i is the output of Algorithm 2 when the input is S'_i . Thus, by the Gaussian mechanism we can see it is (ϵ, δ) -DP.

Excess risk: We finish our proof through several parts. We first recall the following lemma.

1734 Lemma 4. [(16)] Assume diam₂(\mathcal{X}) $\leq D$. Let $\mathcal{S} = (S_1, \ldots, S_n)$ where $S_1^n \stackrel{iid}{\sim} P$ and f(w, x) **1735** is L-Lipschitz and λ -strongly convex for all $x \in \mathcal{X}$. Let $\hat{x} = \operatorname{argmin}_{x \in \mathcal{X}} \bar{F}(w)$ be the empirical **1736** minimizer. For $0 < \beta \leq 1/n$, with probability at least $1 - \beta$

$$F(\hat{x}) - F(x^{\star}) \leqslant \frac{cL^2 \log(n) \log(1/\beta)}{\lambda n} + \frac{cLD\sqrt{\log(1/\beta)}}{\sqrt{n}}$$

Theorem 8. We have the following bound for $||w_T - \hat{w}||^2$ for T iterations:

$$\|w_T - \hat{w}\|^2 \leq \exp\{-\frac{\lambda\eta T}{2}\} \|w_0 - \hat{w}\|^2 + \frac{8\eta \hat{r}_n^2(x)}{\lambda} + 8\eta\lambda D^2 + \frac{20\hat{B}^2}{\lambda^2}.$$

Proof. Detailly,

$$\begin{split} \left\| \tilde{\nabla} F_{\lambda}\left(w_{t}\right) \right\|^{2} &\leq 2 \left(\left\| \nabla \widehat{F}_{\lambda}\left(w_{t}\right) \right\|^{2} + \left\| b_{t} \right\|^{2} \right) \\ &\leq 2 \left(2 \widehat{r}_{n}(X)^{2} + 2\lambda^{2}D^{2} + \hat{B}^{2} \right), \end{split}$$

And also, by Young's inequality,

$$\langle b_t, w_t - \hat{w} \rangle | \leq \frac{\widehat{B}^2}{\lambda} + \frac{\lambda}{4} \|w_t - \hat{w}\|^2.$$

755 Set $\widetilde{\nabla}F_{\lambda}(w_t) = \nabla \hat{F}_{\lambda}(w_t) + b_t = \frac{1}{n} \sum_{i=1}^{n} \prod_C (\nabla f(w, x_i)) + \lambda (w - w_0)$ as the biased, noisy subgradients of the regularized empirical loss in Algorithm 3, with $N_t \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ and $b_t =$ $\frac{1}{n}\sum_{i=1}^{n}\Pi_{C}\left(\nabla f\left(w_{t},x_{i}\right)\right)-\frac{1}{n}\sum_{i=1}^{n}\nabla f\left(w_{t},x_{i}\right). \text{ Denote } y_{t+1}=w_{t}-\eta\tilde{\nabla}F_{\lambda}\left(w_{t}\right), \text{ so that } w_{t+1}=\Pi_{\mathcal{W}}\left(y_{t+1}\right). \text{ For now, by strong convexity, we have}$ $\widehat{F}_{\lambda}(w_{t}) - \widehat{F}_{\lambda}(\hat{w}) \leq \left\langle \nabla \widehat{F}_{\lambda}(w_{t}), w_{t} - \hat{w} \right\rangle - \frac{\lambda}{2} \left\| w_{t} - \hat{w} \right\|^{2}$ $=\left\langle \widetilde{\nabla}F_{\lambda}\left(w_{t}\right),w_{t}-\hat{w}\right\rangle -\frac{\lambda}{2}\left\|w_{t}-\hat{w}\right\|^{2}+\left\langle \nabla\hat{F}_{\lambda}\left(w_{t}\right)-\widetilde{\nabla}F_{\lambda}\left(w_{t}\right),w_{t}-\hat{w}\right\rangle$ $=\frac{1}{2n}\left(\|w_{t}-\hat{w}\|^{2}+\|w_{t}-y_{t+1}\|^{2}-\|y_{t+1}-\hat{w}\|^{2}\right)-\frac{\lambda}{2}\|w_{t}-\hat{w}\|^{2}$ $+\left\langle \nabla \hat{F}_{\lambda}\left(w_{t}\right)-\tilde{\nabla}F_{\lambda}\left(w_{t}\right),w_{t}-\hat{w}\right\rangle$ $=\frac{1}{2n}\left(\left\|w_{t}-\hat{w}\right\|^{2}(1-\lambda\eta)-\left\|y_{t+1}-\hat{w}\right\|^{2}\right)+\frac{\eta}{2}\left\|\widetilde{\nabla}F_{\lambda}\left(w_{t}\right)\right\|^{2}$ $+\left\langle \nabla \widehat{F}_{\lambda}\left(w_{t}\right)-\tilde{\nabla}F_{\lambda}\left(w_{t}\right),w_{t}-\hat{w}\right\rangle$ $\leq \frac{1}{2n} \left(\left\| w_t - \hat{w} \right\|^2 (1 - \lambda \eta) - \left\| w_{t+1} - \hat{w} \right\|^2 \right) + \frac{\eta}{2} \left\| \widetilde{\nabla} F_\lambda \left(w_t \right) \right\|^2 - \left\langle b_t, w_t - \hat{w} \right\rangle,$

where we used non-expansiveness of projection and the definition of $\tilde{\nabla}F_{\lambda}(w_t)$ in the last line. Now, re-arranging this inequality,

$$\leqslant \|w_t - \hat{w}\|^2 \left(1 - \frac{\lambda\eta}{2}\right) + \eta^2 \cdot 2(2\hat{r}_n^2(x) + 2\lambda^2 D^2 + \hat{B}^2) + \frac{2\eta B}{\lambda}$$

$$\leqslant \|w_t - \hat{w}\|^2 \left(1 - \frac{\lambda\eta}{2}\right) + 4\eta^2 (\hat{r}_n^2(x) + \lambda^2 D^2 + \hat{B}^2) + \frac{2\eta \hat{B}^2}{\lambda},$$

where \widehat{B} is defined as below,

$$\widehat{B} = \sup_{t \in T} \|b_t\| \leqslant \frac{\widehat{r}_n(X)^{(k)}}{(k-1)C^{k-1}}$$

Thus, iterating the above equation, we get

$$\begin{split} \|w_{T} - \hat{w}\|^{2} &\leq (1 - \frac{\lambda\eta}{2})^{T} \|w_{0} - \hat{w}\|^{2} + (4\eta^{2}(\hat{r}_{n}^{2}(x) + \lambda^{2}D^{2} + \hat{B}^{2}) + \frac{2\eta\hat{B}^{2}}{\lambda})\sum_{t=1}^{T-1}(1 - \frac{\lambda\eta}{2})^{t} \\ &\leq (1 - \frac{\lambda\eta}{2})^{T} \|w_{0} - \hat{w}\|^{2} + (4\eta^{2}(\hat{r}_{n}^{2}(x) + \lambda^{2}D^{2} + \hat{B}^{2}) + \frac{2\eta\hat{B}^{2}}{\lambda})\frac{2}{\lambda\eta} \\ &= (1 - \frac{\lambda\eta}{2})^{T} \|w_{0} - \hat{w}\|^{2} + \frac{8\eta}{\lambda}(\hat{r}_{n}^{2}(x) + \lambda^{2}D^{2} + \hat{B}^{2}) + \frac{4\hat{B}^{2}}{\lambda^{2}} \\ &\leq \exp\{-\frac{\lambda\eta T}{2}\} \|w_{0} - \hat{w}\|^{2} + \frac{8\eta\hat{r}_{n}^{2}(x)}{\lambda} + 8\eta\lambda D^{2} + \frac{8\eta\hat{B}^{2}}{\lambda} + \frac{4\hat{B}^{2}}{\lambda^{2}} \\ &\leq \exp\{-\frac{\lambda\eta T}{2}\} \|w_{0} - \hat{w}\|^{2} + \frac{8\eta\hat{r}_{n}^{2}(x)}{\lambda} + 8\eta\lambda D^{2} + \frac{20\hat{B}^{2}}{\lambda^{2}}. \end{split}$$

The last inequality holds due to the assumption that $\eta \leq \frac{2}{\lambda}$.

Theorem 9. We have the following bound for $f(w_l) - f(\hat{w}_l)$:

$$F(w_l) - F(\hat{w}_l) \leqslant \widetilde{O}\left(\frac{\widetilde{r}_{2k,n_l}}{\widetilde{R}_{2k,n}} \cdot \frac{DL_f}{\sqrt{n}}\right).$$

Proof. Firstly, the choise of D_i ensures that $\hat{w}_i \in \mathcal{W}_i$.

Then by the above lemma, and choosing specific T_i ,

$$\|w_i - \hat{w}_i\|^2 \leqslant \exp\{-\frac{\lambda_i \eta_i T_i}{2}\} \|w_{i-1} - \hat{w}_i\|^2 + \frac{8\eta_i \hat{r}_{n_i}^2 (B_i)^{(2)}}{\lambda_i} + 8\eta_i \lambda_i D_i^2 + \frac{20\hat{r}_{n_i}^2 (B_i)^{(2)}}{\lambda_i^2 (k-1)C_i^{k-1}}.$$

$$\|w_i - \hat{w}_i\|^2 \lesssim \frac{\eta_i}{2}L_f^2 + \frac{1}{2}$$

$$\|w_i - \hat{w}_i\|^2 \lesssim \frac{\eta_i}{\lambda_i} L_f^2 + \frac{\widetilde{r}_{n_i}^{(2k)}}{\lambda_i^2 C_i^{2k-2} 4^i} \lesssim \frac{\eta^2 n}{16^i 4^i} (L_f^2 + \frac{n \widetilde{r}_{n_i}^{(2k)}}{C_i^{2k-2} 4^i}).$$
(5)

Then by setting $L = \sup_{w \in \mathcal{W}} \|\nabla F(w)\| \leq r$. Therefore,

$$F(w_{l}) - F(\hat{w}_{l}) \leqslant \sqrt{\|w_{l} - \hat{w}_{l}\|^{2}} \\ \leqslant L \sqrt{\eta_{l}^{2} (L_{f}^{2} + \frac{\tilde{e}_{n_{i}}^{(2k)}}{C_{l}^{2k-2} 4^{i}})} \\ \lesssim L \frac{\eta}{n^{2}} (L_{f} + \frac{\tilde{r}_{2k}^{k}}{C_{l}^{k-1}}) \\ \lesssim L \frac{\eta}{n^{2}} (L_{f} + \frac{\tilde{r}_{2k}^{k}}{C_{l}^{k-1}}) \\ \leqslant L \frac{\eta}{n^{2}} \left(L_{f} + \tilde{r}_{2k} (\frac{\sqrt{d}}{\epsilon})^{\frac{k-1}{k}} \right).$$

We know that $\xi_i \sim \mathcal{N}(0, \sigma_i^2)$ and ξ is sub-Gaussian, thus, we can derive that

$$\mathbb{P}\{\|\xi_i\| \ge t\sqrt{d}\} \le 2\exp\{-\frac{t^2}{16\sigma_i^2}\}.$$

Here there shall be some confusion about the lower index, where k is equivalent to l as above, not the original k here. Therefore, with probability $1 - \beta$, $\|\xi_i\| \leq 4\sqrt{d\sigma_i \log(4/\beta)}$. Thus, due to the choice of η , we have

$$\begin{split} F(w_l) - F(\hat{w}_l) \leqslant & 4L_f \sqrt{d} \sigma_l \log(4/\beta) = 4L_f \sqrt{d} \log(4/\beta) \frac{8C_l \sqrt{\log(1/\delta)}}{n_l \lambda_l \epsilon} \\ &= & 32L_f \sqrt{d} \log(4/\beta) \frac{C_l \eta_l n_l^{p-1}}{\epsilon} \\ &= & 32L_f \sqrt{d} \log(4/\beta) \widetilde{r}_{2k,n_l} \left(\frac{\epsilon n_l}{\sqrt{d\log(n)}}\right)^{\frac{1}{k}} \frac{\eta}{4^l} \frac{n^{p-1}}{(2^l)^{p-1}} \frac{1}{\epsilon} \\ &\leqslant & \widetilde{r}_{2k,n_l} \cdot \frac{32DL \log(1/\beta)}{\sqrt{n} \log^{p+\frac{5}{2}} n}. \end{split}$$

Finally, we reach the upper bound for $F(w_l) - F(w^*)$:

Theorem 10. *Finally, we reach the upper bound for* $F(w_l) - F(w^*)$ *:*

$$F(w_l) - F(w^*) \lesssim \widetilde{R}_{2k,n} D(\frac{1}{\sqrt{n}} + (\frac{\sqrt{d\log n}}{\epsilon n})^{\frac{k-1}{k}}) + \frac{D\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n}}$$

Proof. Rewrite this term into summation of their differences,

$$F(w_l) - F(w^*) = \sum_{i=1}^{l} [f(\hat{w}_i) - f(\hat{w}_{i-1})] + [f(w_l) - f(\hat{w}_l)],$$

 $\leq \frac{D^2}{nn^{2p}} + \eta \left(L_f^2 + \widetilde{R}_{2k,n} n^p \left(\frac{d\log n}{\epsilon^2 n^2} \right)^{\frac{k-1}{k}} \right)$

which holds with probability at least $1 - \beta$.

to l, we have with probability at least $1 - \beta$, for some constant C_0

 $\leq C_0 \sum_{i=1}^{l} \left\{ \lambda_i \| \widehat{w}_{i-1} - w_{i-1} \|^2 + \frac{cL_f^2 \log n_i \log(1/\beta)}{\lambda_i n_i} + \frac{cL_f D_i \sqrt{\log(1/\beta)}}{\sqrt{n_i}} \right\}$

 $+L_{f}^{2}\eta n^{p-1}\log(1/\beta)\left(\frac{\log n}{2^{p+1}}+\frac{1}{2^{p+1}\log^{p}n}\right)+L_{f}^{2}\eta n^{p-\frac{1}{2}}\sqrt{\log(1/\beta)}\frac{1-\frac{1}{n^{p+\frac{3}{2}}}}{2^{p+\frac{3}{2}}-1}$

 $F(\hat{w}_i) - F(\hat{w}_{i-1}) \leqslant \frac{cL^2 \log n_i \log (2/\beta)}{\lambda n_i} + \frac{cLD_i \sqrt{\log(2/\beta)}}{\sqrt{n_i}} + \frac{\lambda_i}{2} \|w_{i-1} - \hat{w}_{i-1}\|^2.$

For $\|w_i - \widehat{w}_i\|^2 \leq \frac{\eta_i}{\lambda_i} L_f^2 + \frac{\widetilde{r}_{n_i}^{(2k)}}{\lambda_i^2 C_i^{2k-2} 4^i} \leq O\left(\frac{\eta^2 n}{16^i 4^i} (L_f^2 + \frac{n \widetilde{r}_{n_i}^{(2k)}}{C_i^{2k-2} 4^i})\right)$, then summing over i from 1

 $\leq \lambda_1 \|\widehat{w}_0 - w_0\|^2 + \sum_{i=1}^l \lambda_i \|\widehat{w}_{i-1} - w_{i-1}\|^2 + \sum_{i=1}^l \frac{L_f^2 \log n_i \log(1/\beta)}{\lambda_i n_i} + \sum_{i=1}^l \frac{L_f D_i \sqrt{\log(1/\beta)}}{\sqrt{n_i}} + \sum_{i=1}^l \frac{L_f D_i \sqrt{\log(1/\beta)}}{\sqrt{\log(1/\beta)}} + \sum_{i=1}^l \frac{L_f D_i \sqrt{\log(1/\beta)}}{\sqrt{\log($

 $\leq \frac{D^2}{\eta n^{2p}} + \sum_{i=1}^{l} \lambda_i \left[\eta_i^2 n_i^p L_f^2 + \frac{\eta_i^2 n_i^{2p} \widetilde{e}_{n_i}^{(2k)}}{C_i^{2k-2}} \right] + \sum_{i=1}^{l} \frac{L_f^2 (\log n - \log 2^i) \log(1/\beta)}{n_i} \eta_i n_i^p + \sum_{i=1}^{l} L_f^2 \eta_i n_i^{p-\frac{1}{2}} \sqrt{\log(1/\beta)}$

 $\leq \frac{D^2}{nn^{2p}} + \eta \left(L_f^2 + \widetilde{R}_{2k,n} n^p (\frac{d\log n}{\epsilon^2 n^2})^{\frac{k-1}{k}} \right) + L_f^2 \eta n^{p-1} \log(1/\beta) \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l (\frac{1}{2^{p+\frac{3}{2}}})^i \left(\frac{1}{2^{p+\frac{3}{2}}} \right)^{\frac{k-1}{2}} + L_f^2 \eta n^{p-1} \log(1/\beta) \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-1})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \sum_{i=1}^l \frac{(\log n-i)}{4^i \cdot (2^{p-\frac{1}{2}})^i} + L_f^2 \sqrt{\log(1/\beta)} + L_f^2 \sqrt{\log(1/\beta)} \sum_{i=1$

 $\leq \frac{D^2}{nn^{2p}} + \eta (L_f^2 + \widetilde{R}_{2k,n} n^p (\frac{d\log n}{\epsilon^2 n^2})^{\frac{k-1}{k}}) + L_f^2 \eta n^{p-1} \log(1/\beta) \log n \cdot 2^{-(p+1)} + L_f^2 \eta n^{p-\frac{1}{2}} \sqrt{\log(1/\beta)} \cdot 2^{-(p+\frac{3}{2})}.$

Assume that $\exists p \text{ s.t. } L_f \leqslant O\left(n^{p/2}\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d\log n}}{\epsilon n}\right)^{\frac{k-1}{k}}\right)\right)$ and take $\eta = 0$

 $f(w_l) - f(w^*) \leqslant O\left(\widetilde{R}_{2k,n} D\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d\log n}}{\epsilon n}\right)^{\frac{k-1}{k}}\right) + \frac{D\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n}}\right),$

 $\frac{D}{n^{\frac{D}{2}}}\min\{\frac{1}{L_f},\frac{1}{\tilde{R}_{2k-n}n^{\frac{p+1}{2}}}(\frac{\epsilon n}{\sqrt{d\log n}})^{\frac{k-1}{k}},\frac{1}{n^{\frac{p-1}{2}}L^2\sqrt{\log n\log(1/\beta)}}\},\text{ then the above can be reduced to}$

 $\leq \frac{D^2}{\eta n^{2p}} + \sum_{i=1}^{l} \left| \eta_i^2 L_f^2 + \frac{\eta_i n_i^p \widetilde{e}_{n_i}^{(2k)}}{C_i^{2k-2}} \right| \\ + \sum_{i=1}^{l} L_f^2 \eta_i n_i^{p-1} (\log n - \log 2^i) \log(1/\beta) \\ + \sum_{i=1}^{l} L_f^2 \eta_i n_i^{p-\frac{1}{2}} \sqrt{\log(1/\beta)}$

By lemma 4,

 $f(w_l) - f(w^*)$

Proof of Thorem 2.

Theorem 11. Assume that loss function $F(\cdot)$ is (θ, λ) -TNC and $f(\cdot, x)$ is convex, α -smooth and L_f -Lipschitz for each x. Then algorithm 4 is (ϵ, δ) -DP based on different stepsizes $\{\gamma_k\}_{k=1}^m$ and noises if $\gamma_k \leq \frac{1}{\alpha}$. Then for sufficiently large n and (ϵ, δ) -DP, with probability at least $1 - \beta$, we have

Proof. The guarantee of (ϵ, δ) -DP is just followed by Theorem 1. For simplicity, we denote $a(n) = O\left(\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d\log n}}{\epsilon n}\right)^{\frac{k-1}{k}}\right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n}}\right)$. We set $\mu_0 = 0$ $2R_0^{1-\theta}a(n_0), \mu_k = 2^{(\theta-1)k}\mu_0$ and $R_k = \frac{R_0}{2^k}$, where $k = 1, \dots, m$. Then we have $\mu_k \cdot R_k^{\theta} = 2^{-k} \mu_0 R_0^{\theta}$. We can also assume that $\lambda \leq \frac{L}{R_0^{\theta-1}}$, otherwise we can set $\lambda = \frac{L}{R_0^{\theta-1}}$, which makes TNC still hold. Recall that $m = \lfloor \frac{1}{2} \log_2 \frac{2n}{\log_2 n} \rfloor - 1$, when $n \ge 256$, it follows that $0 < \frac{1}{2}\log_2 \frac{2n}{\log_2 n} - 2 \le m \le \frac{1}{2}\log_2 \frac{2n}{\log_2 n} - 1 \le \frac{1}{2}\log_2 n.$ Thus, we have $2^m \ge \frac{1}{4}\sqrt{\frac{2n}{\log_2 n}}$ (if we pick specific m such that $2^m \ge \frac{1}{4}\sqrt{\frac{2n}{\log_2 n}} \cdot \frac{1}{\log n_0\sqrt{\log(1/\beta)}}$) Thus $\mu_m = 2^{(\theta - 1)m} \mu_0 \ge 2^m \mu_0$ $\geq \frac{1}{4}\sqrt{\frac{2n}{\log_2 n}}\frac{1}{\log n_0\sqrt{\log\left(1/\beta\right)}} \cdot 2 \cdot R_0^{1-\theta}a\left(n_0\right)$ $=\frac{5\cdot R_0^{1-\theta}}{\log n_0\sqrt{\log(1/\beta)}}\sqrt{\frac{2n}{\log_2 n}}\left(\widetilde{R}_{2k,n_0}(\frac{1}{\sqrt{n_0}}+(\frac{\sqrt{d\log n_0}}{\epsilon n_0})^{\frac{k-1}{k}})+\frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_0}}\right)$ $\geq 5 \cdot \widetilde{R}_{2k,n_0} R_0^{1-\theta} \sqrt{\frac{2n}{\log_2 n}} \left(\frac{1}{\sqrt{\frac{2n}{\log_2 n}}} \right)$ $= 5 \cdot \tilde{R}_{2k,n_0} R_0^{1-\theta} \sqrt{\frac{\log_2 2n - \log_2 \log_2 n - 4}{\log_2 n}} \cdot \log n_0 \sqrt{\log(1/\beta)}$ $\geq \widetilde{R}_{2k,n_0} R_0^{1-\theta} \left(\text{ Since } 5 \cdot \sqrt{\frac{\log_2 2n - \log_2 \log_2 n - 4}{\log_2 n}} \geq 1 \text{ when } n \geq 256 \right)$ $\geq \lambda$ (By assumption). where the third inequality is given by throwing away the $\left(\frac{\sqrt{d \log n_0}}{\epsilon n_0}\right)^{\frac{k-1}{k}}$ and $\frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_0}}$ term and substituting *m* in term $\frac{1}{\sqrt{\frac{n}{m}}}$ with $\frac{1}{2} \log_2 \frac{2n}{\log_2 n} - 2$. Below, we consider the following two cases. **Case 1** If $\lambda \ge \mu_0$, then $\mu_0 \le \lambda \le \mu_m$. We have the following lemma. **Lemma 5.** Let k^* satisfies $\mu_{k^*} \leq \lambda \leq 2^{\theta-1}\mu_{k^*}$, then for any $1 \leq k \leq k^*$, the points $\{\hat{w}_k\}_{k=1}^m$ generated by Algorithm 4 satisfy $\|\hat{w}_{k-1} - w^*\|_2 \leq R_{k-1} = 2^{-(k-1)} \cdot R_0.$ (6) $F(\hat{w}_k) - F(w^*) \leq \mu_k R_k^{\theta} = 2^{-k} \mu_0 R_0^{\theta}$ (7)*Moreover, for* $k \ge k^*$ *, we have* $F\left(\hat{w}_{k}\right) - F\left(\hat{w}_{k^{*}}\right) \leqslant \mu_{k^{*}} R_{k^{*}}^{\theta}.$ (8)

Proof. We prove (6), (7) by induction. Note that (6) holds for k = 1. Assume (6) is true for some k > 1, then we have

$$F(\hat{w}_{k}) - F(w^{*}) \leqslant R_{k-1} \cdot \left(\widetilde{R}_{2k,n_{0}} \left(\frac{1}{\sqrt{n_{0}}} + \left(\frac{\sqrt{d \log n_{0}}}{\epsilon n_{0}} \right)^{\frac{k-1}{k}} \right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{0}}} \right)$$
$$= R_{k-1}a(n_{0})$$
$$= \frac{1}{2}\mu_{k}2^{(1-\theta)k}R_{0}^{\theta-1}R_{k-1}$$

$$R_k^{ heta}$$

 $= \mu_k$

972 Which is (7). By the definition of TNC, we have

$$\begin{aligned} \|\hat{w}_{k} - w^{*}\|_{2}^{\theta} &\leq \frac{1}{\lambda} \left(F\left(\hat{w}_{k}\right) - F\left(w^{*}\right) \right) \\ &\leq \frac{F\left(\hat{w}_{k}\right) - F\left(w^{*}\right)}{\mu_{k^{*}}} \\ &\leq \frac{\mu_{k}R_{k}^{\theta}}{\mu_{k^{*}}} \leq R_{k}^{\theta} \end{aligned}$$

Thus (6) is true for k + 1. Now we prove (8). Referring to Theorem 1, we know that

$$F(\hat{w}_k) - F(\hat{w}_{k-1}) \leq R_{k-1} \cdot a(n_0)$$

= $2^{k^*-k} R_{k^*-1} a(n_0)$
= $2^{k^*-k} \mu_{k^*} R_{k^*}^{\theta}$
= $\mu_k R_k^{\theta}$

Thus, for $k > k^*$,

$$F(\hat{w}_{k}) - F(\hat{w}_{k^{*}}) = \sum_{j=k^{*}+1}^{k} \left(F(\hat{w}_{j}) - F(\hat{w}_{j-1})\right)$$
$$\leqslant \sum_{j=k^{*}+1}^{k} 2^{k^{*}-j} \mu_{k^{*}} R_{k^{*}}^{\theta}$$
$$= \left(1 - 2^{k^{*}-k}\right) \mu_{k^{*}} R_{k^{*}}^{\theta}$$
$$\leqslant \mu_{k^{*}} R_{k^{*}}^{\theta}$$

Here completes the proof of the lemma. Now we proceed to prove Theorem 1 in this case.

 $F(\hat{w}_m) - F(w^*) = (F(\hat{w}_m) - F(\hat{w}_{k^*})) + (F(\hat{w}_{k^*}) - F(w^*))$ $\leq 2\mu_{k^*}R_{k^*}^{\theta}$ $\leq 4 \left(\frac{\mu_{k^*}}{\lambda}\right)^{\frac{1}{\theta-1}} \mu_{k^*} R_{k^*}^{\theta} \left(\operatorname{Since} \left(\frac{\mu_{k^*}}{\lambda}\right)^{\frac{1}{\theta-1}} \geq \frac{1}{2} \right)$ $=4\left(\frac{2^{(\theta-1)k^*}\mu_0}{\lambda}\right)^{\frac{1}{\theta-1}}\mu_{k^*}R_{k^*}^{\theta}$ $= 4 \left(2^{k^*} \mu_{k^*} R_{k^*}^{\theta} \mu_0^{\frac{1}{\theta-1}} \left(\frac{1}{\lambda} \right)^{\frac{1}{\theta-1}} \right)$ $=4\left(\mu_0 R_0^{\theta} \mu_0^{\frac{1}{\theta-1}} \left(\frac{1}{\lambda}\right)^{\frac{1}{\theta-1}}\right)$ $=4\left(R_{0}^{\theta}\mu_{0}^{\frac{\theta}{\theta-1}}\left(\frac{1}{\lambda}\right)^{\frac{1}{\theta-1}}\right)$ $= 4 \cdot \left((2 \cdot a (n_0))^{\frac{\theta}{\theta-1}} \left(\frac{1}{\lambda}\right)^{\frac{1}{\theta-1}} \right)$ $=4\cdot\left(\frac{1}{\lambda}\right)^{\frac{1}{\theta-1}}\cdot 2\left(\widetilde{R}_{2k,n_0}(\frac{1}{\sqrt{n_0}}+(\frac{\sqrt{d\log n_0}}{\epsilon n_0})^{\frac{k-1}{k}})+\frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_0}}\right)^{\frac{\theta}{\theta-1}}$

where $m = O(\log_2 n)$ (Recall that $m \leq \frac{1}{2} \log_2 n$). **Case 2** If $\lambda < \mu_0$, then $F\left(\hat{w}_{1}\right) - F\left(w^{*}\right) \leqslant R_{0}a\left(n_{0}\right)$ $= \left(\frac{2}{\mu_0}\right)^{\frac{1}{\theta-1}} \cdot a \left(n_0\right)^{\frac{\theta}{\theta-1}}$ $<\left(\frac{2}{\lambda}\right)^{\frac{1}{\theta-1}} \cdot a\left(n_0\right)^{\frac{\theta}{\theta-1}}$ Also, we have $F(\hat{w}_m) - F(\hat{w}_1) = \sum_{j=2}^{m} (F(\hat{w}_j) - F(\hat{w}_{j-1}))$ $\leqslant \sum_{i=2}^{m} R_{j-1} \cdot a\left(n_0\right)$ $=\sum_{i=2}^{m} 2^{-(j-1)} R_0 \cdot a(n_0)$ $= (1 - (1/2)^{m-1}) R_0 \cdot a(n_0) < R_0 \cdot a(n_0)$ By a similar argument process as in Case 1, we have $F(\hat{w}_m) - F(w^*) = (F(\hat{w}_m) - F(\hat{w}_1)) + (F(\hat{w}_1) - F(w^*))$ $\leqslant 2R_0 a\left(n_0\right) \leqslant 2\left(\frac{2}{\lambda}\right)^{\frac{1}{\theta-1}} \cdot a\left(n_0\right)^{\frac{\theta}{\theta-1}}$

$$=2\cdot\left(\frac{2}{\lambda}\right)^{\frac{1}{\theta-1}}\cdot\left(\widetilde{R}_{2k,n_0}\left(\frac{1}{\sqrt{n_0}}+\left(\frac{\sqrt{d\log n_0}}{\epsilon n_0}\right)^{\frac{k-1}{k}}\right)+\frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_0}}\right)^{\frac{\theta}{\theta-1}}$$

Combining the two cases, we conclude that with probability at least $1 - \beta$,

$$\begin{array}{l} 1059\\ 1060\\ 1061\\ 1062\\ 1062\\ 1063\\ 1064\\ 1065 \end{array} F\left(\hat{w}_{m}\right) - F\left(w^{*}\right) \leqslant O\left(\frac{1}{\lambda^{\frac{1}{\theta-1}}} \cdot \left(\widetilde{R}_{2k,n}(\frac{\sqrt{\log n}}{\sqrt{n}} + \left(\frac{\sqrt{d\log^{3} n}}{\epsilon n}\right)^{\frac{k-1}{k}}) + \frac{\sqrt{\log n\log(1/\beta)}}{2^{p+1}\sqrt{n}}\right)^{\frac{\theta}{\theta-1}}\right) \right)$$

Proof of Theorem 3. *Proof.* The guarantee of (ϵ, δ) -DP is just followed by Theorem 1 and the parallel theorem of Differential Privacy. In the following we focus on the utility.

Since $k = \lfloor (\log \log_{\bar{\theta}} 2) \cdot \log \log n \rfloor$, then $k \leq (\log_{\bar{\theta}} 2) \cdot \log \log n$, namely $2^k \leq (\log n)^{\log 2}$ and $\frac{2^k-1}{(\log n)^{\log_{\theta^2}}} \leq 1$. Observe that the total sample number used in the algorithm is $\sum_{i=1}^k n_i \leq 1$ $\textstyle{\textstyle \sum_{i=1}^k \frac{2^{i-1}n}{(\log n)^{\log_{\bar{\theta}}2}} = \frac{\left(2^k-1\right)n}{(\log n)^{\log_{\bar{\theta}}2}} \leqslant n}.$

For the output of phase i, denote $\Delta_i = F(w_i) - F(w^*)$, and let $D_i^{\theta} = ||w_i - w^*||_2^{\theta}$. The assumption of TNC implies that $F(w_i) - F(w^*) \ge \lambda \|w_i - w^*\|_2^{\theta}$, which is $F(w_i) - F(w^*) \ge \lambda \|w_i - w^*\|_2^{\theta}$ when we take expectations at both sides, namely

$$\Delta_i \ge \lambda D_i^{\theta}.\tag{9}$$

1080 Thus, we have

1082

1086

$$\Delta_{i} \leqslant c\widetilde{R}_{2k,n}D_{i-1}\left(\frac{1}{\sqrt{n_{i}}} + \left(\frac{\sqrt{d\log n_{i}}}{\epsilon n_{i}}\right)^{\frac{k-1}{k}}\right) + \frac{cD_{i-1}\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{i}}}$$

$$\stackrel{(9)}{\leqslant} \left(\frac{\Delta_{i-1}}{\lambda}\right)^{\frac{1}{\theta}} \left(c\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n_{i}}} + \left(\frac{\sqrt{d\log n_{i}}}{\epsilon n_{i}}\right)^{\frac{k-1}{k}}\right) + \frac{e\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{i}}}\right),$$
(10)

where the first inequality comes from Theorem 1 and the second inequality uses (9). Denote $E_{i} = \frac{c^{\theta}}{\lambda} \left(\widetilde{R}_{2k,n} \left(\frac{1}{\sqrt{n_{i}}} + \left(\frac{\sqrt{d \log n_{i}}}{\epsilon n_{i}} \right)^{\frac{k-1}{k}} \right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{i}}} \right)^{\theta}$ Then (10) can be simplified as $\Delta_{i} \leq \left(\Delta_{i-1} E_{i} \right)^{\frac{1}{\theta}}.$ (11)

1092

1093 1094 Notice that $n_i/n_{i-1} = 2$, then $\frac{E_{i-1}}{E_i} \leq \left(\frac{n_i}{n_{i-1}}\right)^{\theta} = 2^{\theta}$, namely: 1095 1096 $E_i > 2^{-\theta}E_{i-1}$.

1096 1097

1133

1098 Then we can rearrange the above inequality as

$$\frac{\Delta_i}{E_i^{\frac{1}{\theta-1}}} \leqslant \frac{\left(\Delta_{i-1}E_i\right)^{\frac{1}{\theta}}}{E_i^{\frac{1}{\theta-1}}} \leqslant 2^{\frac{1}{\theta-1}} \left(\frac{\Delta_{i-1}}{E_{i-1}^{\frac{1}{\theta-1}}}\right)^{\frac{1}{\theta}},\tag{13}$$

(12)

¹¹⁰³ where the first inequality uses (11) and the second inequality applies (12).

1104 1105 It can be verified that (13) is equivalent to

$$\frac{\Delta_i}{2^{\frac{\theta}{(\theta-1)^2}}E_i^{\frac{1}{\theta-1}}} \leqslant \left(\frac{\Delta_{i-1}}{2^{\frac{\theta}{(\theta-1)^2}}E_{i-1}^{\frac{1}{\theta-1}}}\right)^{\frac{1}{\theta}} \leqslant \left(\frac{\Delta_1}{2^{\frac{\theta}{(\theta-1)^2}}E_1^{\frac{1}{\theta-1}}}\right)^{\frac{1}{\theta^{i-1}}}$$

According to Lemma 1, $\Delta_1 \leq (L^{\theta} \lambda^{-1})^{\frac{1}{\theta-1}}$. Also observe that

$$E_{1} = \frac{c^{\theta}}{\lambda} \left(\widetilde{R}_{2k,n} \left(\frac{1}{\sqrt{n_{1}}} + \left(\frac{\sqrt{d \log n_{1}}}{\epsilon n_{1}} \right)^{\frac{k-1}{k}} \right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{1}}} \right)^{\theta} \ge \frac{c^{\theta} \widetilde{R}_{2k,n}^{\theta}}{\lambda} \frac{1}{\left(\sqrt{n_{1}}\right)^{\theta}} \ge \frac{c^{\theta} \widetilde{R}_{2k,n}^{\theta}}{\lambda} \frac{1}{n^{\theta}}.$$

$$E_{1} = \frac{c^{\theta}}{\lambda} \left(\widetilde{R}_{2k,n} \left(\frac{1}{\sqrt{n_{1}}} + \left(\frac{\sqrt{d \log n_{1}}}{\epsilon n_{1}} \right)^{\frac{k-1}{k}} \right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_{1}}} \right)^{\theta} \ge \frac{c^{\theta} \widetilde{R}_{2k,n}^{\theta}}{\lambda} \frac{1}{\left(\sqrt{n_{1}}\right)^{\theta}} \ge \frac{c^{\theta} \widetilde{R}_{2k,n}^{\theta}}{\lambda} \frac{1}{n^{\theta}}.$$

Let $c_1 = c^{\frac{\theta}{\theta-1}} 2^{\frac{\theta}{(\theta-1)^2}}$, then $\frac{\Delta_1}{2^{\frac{\theta}{(\theta-1)^2}} E_1^{\frac{1}{\theta-1}}} \leq \frac{n^{\frac{\theta}{\theta-1}}}{c_1}$, which implies that for $l = \lfloor (\log_{\bar{\theta}} 2) \cdot \log\log n \rfloor$, 1117

$$\frac{\Delta_l}{2^{\frac{\theta}{(\theta-1)^2}} E_l^{\frac{1}{\theta-1}}} \leqslant \left(\frac{n^{\frac{\theta}{\theta-1}}}{c_1}\right)^{\frac{1}{\theta^{l-1}}}$$

Let $C_1 = 2^{\frac{\theta^3}{\theta-1} + \theta^2 |\log c_1|}$. In the following we prove that

1124
1125
1126
1127

$$\left(\frac{n^{\frac{\theta}{\theta-1}}}{c_1}\right)^{\frac{1}{\theta^{l-1}}} \leqslant C_1.$$

1128 Since $l + 1 \ge (\log_{\bar{\theta}} 2) \log \log n \ge (\log_{\theta} 2) \log \log n$, it follows that

1129
1130
$$(l-1)\log\theta + \log\log C_1 \ge \log\left(\frac{\theta}{\theta-1} + |\log c_1|\right) + \log\log n$$
1131

1132 which indicates

$$\left(\frac{\theta}{\theta-1} + |\log c_1|\right)\log n \leqslant \theta^{l-1}\log C_1.$$

1135 Thus we have $\frac{\theta}{\theta-1}\log n - \log c_1 \leq \theta^{l-1}\log C_1$, which is equivalent to our object $\left(\frac{n^{\frac{\theta}{\theta-1}}}{c_1}\right)^{\frac{1}{\theta^{k-1}}} \leq C_1$. Now we know

$$\frac{\Delta_l}{2^{\frac{\theta^2}{(\theta-1)^2}} E_l^{\frac{1}{\theta-1}}} \leqslant \left(\frac{n^{\frac{\theta}{\theta-1}}}{c_1}\right)^{\frac{1}{\theta^{l-1}}} \leqslant C_1,$$

which indicates that $\frac{\Delta_l}{E_l^{\frac{\theta}{\theta-1}}} \leq 2^{\frac{\theta}{(\theta-1)^2}}C_1 = 2^{\theta^2 \left(\frac{\theta^2-\theta+1}{(\theta-1)^2}+|\log c_1|\right)} := C$. As a result, we hold a solution with error:

$$F(w_l) - F(w^*) \leqslant CE_l^{\frac{1}{\theta-1}} = C\left(\frac{c^{\theta}}{\lambda}\right)^{\frac{1}{\theta-1}} \left(\widetilde{R}_{2k,n}\left(\frac{1}{\sqrt{n_l}} + \left(\frac{\sqrt{d\log n_l}}{\epsilon n_l}\right)^{\frac{k-1}{k}}\right) + \frac{\sqrt{\log(1/\beta)}}{2^{p+1}\sqrt{n_l}}\right)^{\frac{\theta}{\theta-1}}$$

Proof of Theorem 4. We first define the set of distributions $\{Q_v\}_{v \in \mathcal{V}}$. Specifically, by the standard 1152 Gilbert-Varshamov bound, there exists a set $\mathcal{V} \subset \{\pm\}^d$ such that: (1) $|\mathcal{V}| \ge 2^{\frac{d}{20}}$, (2) for all $v, v' \in \mathcal{V}$, 1153 $d_{ham}(v, v') \ge \frac{d}{8}$ (1). For each $v \in \mathcal{V}$, we define Q_v as

$$X_{v} = \begin{cases} 0, \text{ with probability } 1 - p \\ p^{-\frac{1}{k}} \frac{\tilde{r}_{k}}{2\sqrt{d}} v, \text{ with probability } p \end{cases}$$
(14)

1158 We can see that for each $X_v \sim Q_v$, we always have $\|\mu_v = \mathbb{E}[X_v]\|_2 = p^{\frac{k-1}{k}} \frac{\tilde{r}_k}{2} = \mu$.

1160 We then consider the loss function $f(w, x) = -\langle w, x \rangle + \frac{1}{\theta} ||w||_{2}^{\theta}$, i.e., $F_{P}(w) = -\langle w, \mathbb{E}_{P}[x] \rangle + \frac{1}{\theta} ||w||_{2}^{\theta}$ for distribution P. By (30) we know it satisfies $(\theta, 1)$ -TNC when $\theta \ge 2$. Moreover, for each Q_{v} we have

$$\mathbb{E}[\sup_{w \in \mathcal{W}} \|\nabla f(w, x)\|_{2}^{k}] = \mathbb{E}[\sup_{w \in \mathcal{W}} \|\|w\|_{2}^{\theta - 2}w - x\|_{2}^{k}] \leqslant \mathbb{E}[\|2x\|_{2}^{k}] = \tilde{r}_{k}^{k} = \tilde{r}^{(k)},$$
(15)

where the first inequality is due to the radius of \mathcal{W} is $(\frac{p^{-\frac{1}{k}}\tilde{r}_k}{2})^{\frac{1}{\theta-1}}$. Thus we can see $F_P(w)$ satisfies Assumption 1. For convenience we denote $F_{Q_v}(w) = F_v(w)$.

¹¹⁶⁸ By the form the $F_v(w)$ we can also see that

$$\nabla F_v(w^*) = 0 \equiv \|w^*\|_2^{\theta-2} w^* = \mu_v.$$
(16)

Thus the optimal solution $w_v^* = \frac{\mu_v}{\mu^{\frac{\theta}{\theta-1}}} \in \mathcal{W}$ by our assumption on n and thus $p \leq 1$. In total we have

$$\mathcal{M}(\mathcal{W}, \mathcal{P}, \mathcal{F}_{k}^{\theta}(\mathcal{P}, \tilde{r}_{k}), \rho) \geq \inf_{\mathcal{A} \in \mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{\mathcal{A}, D \in Q_{v}^{n}} [F_{v}(\mathcal{A}(D)) - \min_{w \in \mathcal{W}} F_{v}(w)],$$
(17)

$$\geq \inf_{\mathcal{A}\in\mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v\in\mathcal{V}} \mathbb{E}_{\mathcal{A},D\in Q_v^n} \|\mathcal{A}(D) - w_v^*\|_2^{\theta} = \inf_{\mathcal{A}\in\mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v\in\mathcal{V}} \mathbb{E}_{\mathcal{A},D\in Q_v^n} \|\mathcal{A}(D) - \frac{\mu_v}{\mu^{\frac{\theta-2}{\theta-1}}}\|_2^{\theta}.$$
(18)

1180 Next, we recall the following private Fano's lemma:

Lemma 6. [Theorem 1.4 in (20)] Let \mathcal{P} be a class of distributions over a data universe \mathcal{X} . For each **distribution** $p \in \mathcal{T}$, there is a deterministic function $\theta(p) \in \mathcal{T}$, where \mathcal{T} is the parameter space. Let $\rho: \mathcal{T} \times \mathcal{T} :\mapsto \mathbb{R}_+$ be a semi-metric function on the space \mathcal{T} and $\Phi: \mathbb{R}_+ \mapsto \mathbb{R}_+$ be a non-decreasing function with $\Phi(0) = 0$. We further assume that $X = \{X_i\}_{i=1}^n$ are n i.i.d observations drawn according to some distribution $p \in \mathcal{P}$, and $Q: \mathcal{X}^n \mapsto \Theta$ be some algorithm whose output Q(X) is an estimator. Consider a set of distributions $\mathcal{V} = \{p_1, p_2, \cdots, p_M\} \subseteq \mathcal{P}$ such that for all $i \neq j$,

•
$$\Phi(\rho(\theta(p_i), \theta(p_j)) \ge \alpha,$$

1188 •
$$D_{KL}(p_i, p_j) \leq \beta$$
, where D_{KL} is the KL-divergence,
1189

1190 • D_{TV}

$$D_{TV}(p_i, p_j) \leqslant \gamma$$
,

1192 then we have for any ρ -zCDP mechanism Q.

$$\frac{1}{M} \sum_{i \in [M]} \mathbb{E}_{X \sim p_i^n, Q} [\Phi(\rho(Q(X), \theta(p_i))] \ge \frac{\alpha}{2} \max\{1 - \frac{n\beta + \log 2}{\log M}, 1 - \frac{\rho(n^2\gamma^2 + n\gamma(1-\gamma)) + \log 2}{\log M}\}$$

1197 Now we will leverage the above lemma to lower bound equation 18. We can see in our set of 1198 probabilities $\{Q_v\}_{v \in \mathcal{V}}$, for any $v, v' \in \mathcal{V}$ we have $D_{TV}(Q_v, Q_{v'}) \leq p$. And

$$\|\frac{\mu_{v}}{\mu^{\frac{\theta-2}{\theta-1}}} - \frac{\mu_{v'}}{\mu^{\frac{\theta-2}{\theta-1}}}\|_{2}^{\theta} = \frac{1}{\mu^{\frac{\theta(\theta-2)}{\theta-1}}} \|p^{\frac{k-1}{k}} \frac{\tilde{r}_{k}}{2\sqrt{d}} (v-v')\|_{2}^{\theta} \ge C \frac{p^{\frac{\theta(k-1)}{k}}}{\mu^{\frac{\theta(\theta-2)}{\theta-1}}} \tilde{r}_{k}^{\theta} = \Omega(\tilde{r}_{k}^{\frac{\theta}{\theta-1}} p^{\frac{k-1}{k}} \frac{\theta}{\theta-1}).$$
(19)

Taking $p = \frac{\sqrt{d}}{n\sqrt{\rho}}$ and by Lemma 6 we have

$$\inf_{\mathcal{A}\in\mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v\in\mathcal{V}} \mathbb{E}_{\mathcal{A},D\in Q_v^n} \|\mathcal{A}(D) - \frac{\mu_v}{\mu^{\frac{\theta-2}{\theta-1}}} \|_2^{\theta} \ge \Omega\left(\left(\tilde{r}_k \left(\frac{\sqrt{d}}{n\sqrt{\rho}} \right)^{\frac{k-1}{k}} \right)^{\frac{\theta}{\theta-1}} \right).$$
(20)

Proof of Theorem 5. The lower bound for non-private case follows the proof in (2). Here we extend to the heavy-tailed case. For the index set \mathcal{V} we consider the same one as in the proof of Theorem 4. For each $v \in \mathcal{V}$ we define $X \sim P_v$ as

For
$$j \in [d], X_j = \begin{cases} v_j e_j \frac{\tilde{r}_k}{2\sqrt{d}}, \text{ with probability } \frac{1+\delta}{2}, \\ -v_j e_j \frac{\tilde{r}_k}{2\sqrt{d}}, \text{ with probability } \frac{1-\delta}{2}. \end{cases}$$
 (21)

1216 We can see that for each $X_v \sim Q_v$, we always have $\|\mu_v = \mathbb{E}[X_v]\|_2 = \delta \frac{\tilde{r}_k}{2} = \mu$.

1217 We then consider the loss function $f(w, x) = -\langle w, x \rangle + \frac{1}{\theta} ||w||_2^{\theta}$, i.e., $F_P(w) = -\langle w, \mathbb{E}_P[x] \rangle + \frac{1}{\theta} ||w||_2^{\theta}$ for distribution *P*. By (30) we know it satisfies $(\theta, 1)$ -TNC when $\theta \ge 2$. Moreover, for each Q_v we have

$$\mathbb{E}[\sup_{w\in\mathcal{W}} \|\nabla f(w,x)\|_{2}^{k}] = \mathbb{E}[\sup_{w\in\mathcal{W}} \|\|w\|_{2}^{\theta-2}w - x\|_{2}^{k}] \leqslant \mathbb{E}[\|2x\|_{2}^{k}] = \tilde{r}_{k}^{k} = \tilde{r}^{(k)},$$
(22)

where the first inequality is due to the radius of \mathcal{W} is $(\frac{\tilde{r}_k}{2})^{\frac{1}{\theta-1}}$. Thus we can see $F_P(w)$ satisfies Assumption 1. For convenience we denote $F_{Q_v}(w) = F_v(w)$.

1225 By the form the $F_v(w)$ we can also see that

$$\nabla F_v(w^*) = 0 \equiv \|w^*\|_2^{\theta-2} w^* = \mu_v.$$
(23)

Thus the optimal solution $w_v^* = \frac{\mu_v}{\mu^{\frac{\theta}{\theta-1}}} \in \mathcal{W}$ by our assumption on n and thus $p \leq 1$. In total we have

$$\mathcal{M}(\mathcal{W}, \mathcal{P}, \mathcal{F}_{k}^{\theta}(\mathcal{P}, \tilde{r}_{k}), \rho) \geq \inf_{\mathcal{A} \in \mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \mathbb{E}_{\mathcal{A}, D \in Q_{v}^{n}} [F_{v}(\mathcal{A}(D)) - \min_{w \in \mathcal{W}} F_{v}(w)],$$
(24)

$$\geq \inf_{\mathcal{A}\in\mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v\in\mathcal{V}} \mathbb{E}_{\mathcal{A},D\in Q_v^n} \|\mathcal{A}(D) - w_v^*\|_2^{\theta} = \inf_{\mathcal{A}\in\mathcal{Q}(\rho)} \frac{1}{|\mathcal{V}|} \sum_{v\in\mathcal{V}} \mathbb{E}_{\mathcal{A},D\in Q_v^n} \|\mathcal{A}(D) - \frac{\mu_v}{\mu^{\frac{\theta-2}{\theta-1}}} \|_2^{\theta}.$$
(25)

1237 We can see in our set of probabilities $\{Q_v\}_{v \in \mathcal{V}}$, for any $v, v' \in \mathcal{V}$ we have $D_{KL}(Q_v, Q_{v'}) \leq \delta^2$. 1238 And

$$\left\|\frac{\mu_{v}}{\mu^{\frac{\theta-2}{\theta-1}}} - \frac{\mu_{v'}}{\mu^{\frac{\theta-2}{\theta-1}}}\right\|_{2}^{\theta} = \frac{1}{\mu^{\frac{\theta(\theta-2)}{\theta-1}}} \left\|\frac{\delta \tilde{r}_{k}}{2\sqrt{d}}(v-v')\right\|_{2}^{\theta} \ge C \frac{\delta^{\theta}}{\mu^{\frac{\theta(\theta-2)}{\theta-1}}} \tilde{r}_{k}^{\theta} = \Omega(\tilde{r}_{k}^{\frac{\theta}{\theta-1}}\delta^{\frac{\theta}{\theta-1}}).$$
(26)

Thus by Fano's lemma or Lemma 6, taking $\delta = \sqrt{\frac{d}{n}}$ we have the result.

Proof of Theorem 6. Proof of Privacy. We first recall the following lemma:

We know that for each $x \in \mathcal{B}_t$, we have $\mathcal{R}(\Pi_C(\nabla f(w, x))) = \Pi_C(\nabla f(w, x)) + \zeta_x$, with $\zeta_x \sim \mathcal{N}(0, \sigma_1^2)$ and $\sigma_1^2 = \frac{8C^2 \log \frac{1}{\delta_0}}{\epsilon_0^2}$) is an (ϵ_0, δ_0) -LDP randomizer. As we randomly shuffled the data in the beginning, thus, the algorithm will be $(\hat{\epsilon}, \hat{\delta} + O(e^{\hat{\epsilon}}\delta_0 n))$ -DP where $\hat{\epsilon} = O\left((1 - e^{-\epsilon_0}) \cdot (\frac{\sqrt{e^{\epsilon_0} \log(1/\hat{\delta})}}{\sqrt{n}} + \frac{e^{\epsilon_0}}{n})\right)$.

Now, assume that $\epsilon_0 \leq \frac{1}{2}$, then $\exists c_1 > 0$, s.t.,

$$\hat{\epsilon} \leqslant c_1 (1 - e^{-\epsilon_0}) \cdot \left(\frac{\sqrt{e^{\epsilon_0} \log(1/\hat{\delta})}}{\sqrt{n}} + \frac{e^{\epsilon_0}}{n}\right)$$

$$\leqslant c_1 \cdot \left((e^{\epsilon_0/2} - e^{-\epsilon_0/2}) \cdot \sqrt{\frac{\log(1/\hat{\delta})}{n}} + \frac{e^{\epsilon_0} - 1}{n} \right)$$

$$\leqslant c_1 \cdot \left(\left(\left(1 + \epsilon_0 \right) - \left(1 - \frac{\epsilon_0}{2} \right) \right) \cdot \sqrt{\frac{\log(1/\hat{\delta})}{n} + \frac{(1 + 2\epsilon_0) - 1}{n}} \right)$$

$$= c_1 \cdot \epsilon_0 \cdot \left(\frac{3}{2}\sqrt{\frac{\log(1/\hat{\delta})}{n}} + \frac{2}{n}\right)$$

1277 Set $\hat{\delta} = \frac{\delta}{2}, \delta_0 = c_2 \cdot \frac{\delta}{e^{\epsilon_n}}$ for some constant $c_2 > 0$ and replace $\epsilon_0 = \frac{2\sqrt{2}C\sqrt{\log \frac{1}{\delta_0}}}{\sigma_1}$): 1282 $\hat{\epsilon} \leqslant c_1 \cdot \frac{2\sqrt{2}C\sqrt{\log \frac{1}{\delta_0}}}{\sigma_1}) \cdot \left(\frac{3}{2}\sqrt{\frac{\log(1/\hat{\delta})}{n}} + \frac{2}{n}\right)$

$$\leq O\left(\frac{C \cdot \sqrt{\log(1/\delta)\log(e^{\hat{\epsilon}}n/\delta)}}{\sigma_1\sqrt{n}}\right).$$

For any $\epsilon \leq 1$, if we set $\sigma = O\left(\frac{C \cdot \sqrt{\log(1/\delta) \log(e^{\hat{\epsilon}}n/\delta)}}{\epsilon \sqrt{n}}\right)$, then we have $\hat{\epsilon} \leq \epsilon$. Furthermore, we need $\epsilon_0 = \frac{2\sqrt{2}C\sqrt{\log\frac{1}{\delta_0}}}{\sigma_1} \leqslant \frac{1}{2}$, which would be ensured if we set $\epsilon = O\left(\sqrt{\frac{\log(n/\delta)}{n}}\right)$. This implies that for $\sigma_1 = O\left(\frac{C \cdot \sqrt{\log(1/\delta) \log(e^{\epsilon}n/\delta)}}{\epsilon \sqrt{n}}\right)$, algorithm 6 satisfies (ϵ, δ) -DP as long as $\epsilon = 0$ $O\left(\sqrt{\frac{\log(n/\delta)}{n}}\right) \text{ if releasing } \mathcal{R}(\Pi_C(\nabla f(w, x))) \text{ for all } x. \text{ Thus in step 6 we can see } \widetilde{\nabla} F_t\left(w_t^{md}\right) = 0$ $\frac{T}{n}\sum_{x\in\mathcal{B}_t} (\mathcal{R}(\Pi_C(\nabla f(w_t^{md}, x))) \text{ is } (\epsilon, \delta) \text{-DP for each } t. \text{ And since } \{B_t\} \text{ are disjoint, Algorithm 6}$ is (ϵ, δ) -DP.

Lemma 8. (4) Let $\{z_i\}_{i=1}^s \sim \mathcal{D}^s$ be \mathbb{R}^d -valued random vectors with $\mathbb{E}z_i = \nu$ and $\mathbb{E} ||z_i||^k \leq r^{(k)}$ for some $k \ge 2$. Denote the noiseless average of clipped samples by $\hat{\nu} := \frac{1}{s} \sum_{i=1}^{s} \prod_{C} (z_i)$ and $\tilde{\nu} := \hat{\nu} + N$. Then, $\|\mathbb{E}\widetilde{\nu}-\nu\| = \|\mathbb{E}\widehat{\nu}-\nu\| \leqslant \mathbb{E}\|\widehat{\nu}-\nu\| \leqslant \frac{r^{(k)}}{(k-1)C^{k-1}}$, and $\mathbb{E}\|\widetilde{\nu}-\mathbb{E}\widetilde{\nu}\|^2 = \mathbb{E}\|\widetilde{\nu}-\mathbb{E}\widehat{\nu}\|^2 \leqslant d\sigma^2 + \frac{r^{(2)}}{s}$. claim: we can improve the noise to $\Sigma^2 := \sup_{t \in [T]} \mathbb{E}[\|N_t\|^2] \leqslant d\sigma^2 + \frac{r^2T}{n} \approx \frac{dC^2T}{\epsilon^2 n^2} + \frac{r^2T}{n}$. Excess risk: Consider round $t \in [T]$ of Algorithm 6, where Algorithm 1 is run on input data $\{\nabla f(w_t, x_i^t)\}_{i=1}^{n/T}$. Denote the bias of Algorithm 1 by $b_t := \mathbb{E}\widetilde{\nabla}F_t(w_t) - \nabla F(w_t)$, where $\widetilde{\nabla}F_t(w_t) = \widetilde{\nu}$ in the notation of Algorithm 1. Also let $\widehat{\nabla}F_t(w_t) := \widehat{\mu}$ (in the notation of Lemma 8) and denote the noise by $N_t = \widetilde{\nabla} F_t(w_t) - \nabla F(w_t) - b_t = \widetilde{\nabla} F_t(w_t) - \mathbb{E} \widetilde{\nabla} F_t(w_t)$. Then we have $B := \sup_{t \in [T]} \|b_t\| \leqslant \frac{r^{(k)}}{(k-1)C^{k-1}}$ and $\Sigma^2 := \sup_{t \in [T]} \mathbb{E}\left[\|N_t\|^2\right] \leqslant d\sigma^2 + \frac{r^2T}{n} \leqslant d\sigma^2$ $O\left(\frac{dC^2T}{r^2n^2} + \frac{r^2T}{n}\right)$, by Lemma 5. Plugging these estimates for B and Σ^2 into Proposition 40 of (29) and setting $C = r \left(\frac{\epsilon n}{\sqrt{d \log(1/\delta)}}\right)^{1/k}$, we get $\mathbb{E}F\left(w_{T}^{ag}\right) - F^{*} \leqslant O\left(\frac{\beta D^{2}}{T^{2}} + \frac{D(\Sigma + B)}{\sqrt{T}} + BD\right)$ $\leqslant O\left(\frac{\beta D^2}{T^2} + \frac{CD\sqrt{d\log(1/\delta)}}{\epsilon n} + \frac{rD}{\sqrt{n}} + \frac{r^{(k)}D}{C^{k-1}}\right)$ $\leqslant O\left(\frac{\beta D^2}{T^2} + rD\left[\frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d\log(1/\delta)}}{\epsilon n}\right)^{(k-1)/k}\right]\right).$ Now, our choice of T $T = \min\{\sqrt{\frac{\beta D}{r}} \cdot \left(\frac{\epsilon n}{\sqrt{d\log(1/\delta)}}\right)^{\frac{\kappa-1}{2k}}, \sqrt{\frac{\beta D}{r}} \cdot n^{1/4}\},\$ implies that $\frac{\beta D^2}{T^2} \leqslant rD \left| \frac{1}{\sqrt{n}} + \left(\frac{\sqrt{d \log(1/\delta)}}{\epsilon n} \right)^{(k-1)/k} \right|$ and we get the result upon plugging in T. Proof of Theorem 7. Similar to the proof of Theorem 3.