# Sharper Utility Bounds for Differentially Private Models

**Yilin Kang, Yong Liu, Jian Li, Weiping Wang**

## Abstract

In this paper, by introducing Generalized Bernstein condition, we propose the first $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ high probability excess population risk bound for differentially private algorithms under the assumptions $G$-Lipschitz, $L$-smooth, and Polyak-Łojasiewicz condition, based on gradient perturbation method. If we replace the properties $G$-Lipschitz and $L$-smooth by $\alpha$-Hölder smoothness (which can be used in non-smooth setting), the high probability bound comes to $\mathcal{O}\left(n^{-\frac{\alpha}{1+2\alpha}}\right)$ w.r.t $n$, which cannot achieve $\mathcal{O}\left(1/n\right)$ when $\alpha \in (0, 1]$. To solve this problem, we propose a variant of gradient perturbation method, **max$\{1, g\}$-Normalized Gradient Perturbation** (m-NGP). We further show that by normalization, the high probability excess population risk bound under assumptions $\alpha$-Hölder smooth and Polyak-Łojasiewicz condition can achieve $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$, which is the first $\mathcal{O}\left(1/n\right)$ high probability excess population risk bound w.r.t $n$ for differentially private algorithms under non-smooth conditions. Moreover, we evaluate the performance of the new proposed algorithm m-NGP, the experimental results show that m-NGP improves the performance of the differentially private model over real datasets. It demonstrates that m-NGP improves the utility bound and the accuracy of the DP model on real datasets simultaneously.

## 1. Introduction

Machine learning has been widely used and found effective in many fields in recent years (Singha et al. 2021; Swapna and Soman 2021; Ponnusamy et al. 2021). When training machine learning models, tremendous data was collected, and the data often contains sensitive information of individuals, which may leakage personal privacy (Shokri et al. 2017; Carlini et al. 2019). Under these circumstances, the privacy of machine learning models is of great importance.

Differential Privacy (DP) (Dwork et al. 2006; Dwork, Roth, and others 2014) is a theoretically rigorous tool to prevent sensitive information leakage. It introduces random noise when training and blocks adversaries from inferring any single individual included in the dataset by observing the model. The mathematical definition of DP is well accepted and relative technologies are performed by Google (Erlingsson, Pihur, and Korolova 2014), Apple (McMillan 2016) and Microsoft (Ding, Kulkarni, and Yekhanin 2017).

As such, DP has attracted attentions from researchers and has been applied to numerous machine learning problems (Ullman and Sealfon 2019; Xu et al. 2019; Bernstein and Sheldon 2019; Wang and Xu 2019; Heikkilä et al. 2019; Kulkarni et al. 2021; Bun, Elias, and Kulkarni 2021; Nguyen and Vullikanti 2021).

There are mainly three approaches to guarantee differential privacy: output perturbation (Chaudhuri, Monteleoni, and Sarwate 2011), objective perturbation (Chaudhuri, Monteleoni, and Sarwate 2011), and gradient perturbation (Song, Chaudhuri, and Sarwate 2013). Considering that gradient descent is a widely used optimization method, the gradient perturbation method can be used for a wide range of applications. Besides, adding random noise to the gradient allows the model to escape local minima (Raginsky, Rakhlin, and Telgarsky 2017), so we focus on the gradient perturbation method to guarantee differential privacy in this paper.

In this paper, we aim to minimize the population risk, and measure the utility of the differentially private model by the excess population risk. To get the excess population risk, an important step is to analyze the generalization error (the reason is demonstrated in Section 3). Complexity theory (Bartlett, Bousquet, and Mendelson 2002) and algorithm stability theory (Bousquet and Elisseeff 2002) are popular tools to analyze the generalization error. On one hand, (Chaudhuri, Monteleoni, and Sarwate 2011) applied the complexity theory and achieved an $\mathcal{O}\left(\max\{\frac{1}{\sqrt{n}}, \sqrt[2/3]{\frac{p}{n\epsilon}}\}\right)$ high probability excess population risk bound under the assumption of strongly convex; (Kifer, Smith, and Thakurta 2012) achieved $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ expected excess population risk bound via complexity theory. On the other hand, the sharpest known high probability generalization bounds for DP algorithms analyzed via stability theory under different assumptions (Wu et al. 2017; Bassily et al. 2019; Feldman, Koren, and Talwar 2020; Bassily et al. 2020; Wang et al. 2021) are $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon} + \frac{1}{\sqrt{n}}\right)$ or $\mathcal{O}\left(\frac{\sqrt[4]{p}}{\sqrt{n\epsilon}}\right)$, containing an inevitable $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ term, which is a bottleneck on the utility analysis. Thus, we are focusing on the following question, which is still an open problem:

*Can we achieve the high probability excess risk bounds with rate $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ for DP models via uniform stability?*

By introducing *Generalized Bernstein condition*

(Koltchinskii 2006), this paper answers the question positively. We remove the $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ term in the generalization error and provide the first high probability excess population risk bound with order $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ in the setting of DP. Comparing with previous high probability bounds, the improvement is approximately up to $\mathcal{O}\left(\sqrt{n}\right)$. The contributions of this paper include: (1) We prove that by introducing Generalized Bernstein condition (Koltchinskii 2006), the high probability excess population risk can be improved to $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$, under Lipschitz and smooth assumptions. To our knowledge, this is the first $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ high probability excess population risk bound for DP model. (2) We relax the assumptions $G$-Lipschitz and $L$-smooth, by introducing $\alpha$-Hölder smooth. Under these assumptions, we prove that the high probability excess population risk bound comes to $\mathcal{O}\left(\frac{\sqrt[4]{p}}{\sqrt{\epsilon}}n^{\frac{-\alpha}{1+2\alpha}}\right)$. Considering that $\alpha \in (0,1]$, the result cannot achieve $\mathcal{O}\left(\frac{1}{n}\right)$ w.r.t $n$, but better than previous one w.r.t $p$ and $\epsilon$. (3) To overcome the bottleneck, we design a variant of gradient perturbation method, called **max $\{1, g\}$-Normalized Gradient Perturbation** (m-NGP) algorithm. Via this new proposed algorithm, we prove that under the assumptions $\alpha$-Hölder smooth and PL condition, the high probability excess population risk bound can be improved to $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$. To the best of our knowledge, this is the first $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ high probability excess population risk bound for non-smooth loss in the field of differential privacy. (4) To evaluate the performance of our proposed max $\{1, g\}$-Normalized Gradient Perturbation algorithm, we perform experiments on several real datasets. The experimental results show that m-NGP improves the accuracy and the convergence rate of the differentially private model on real datasets.

The rest of the paper is organized as follows. Related work is given in Section 2. Preliminaries are introduced in Section 3. In Section 4, we propose sharper utility bounds under different assumptions and design a variant of gradient perturbation method, **max $\{1, g\}$-Normalized Gradient Perturbation**. The experimental results are shown in Section 5. Finally, we conclude the paper in Section 6.

## 2. Related Work

(Dwork et al. 2006) proposed the mathematical definition of differential privacy for the first time. Then, it was developed to protect the privacy in the field of machine learning (e.g. Empirical Risk Minimization (ERM)) via output perturbation, objective perturbation, and gradient perturbation methods. For DP-ERM formulations, (Chaudhuri, Monteleoni, and Sarwate 2011) first proposed output perturbation and objective perturbation methods, and (Song, Chaudhuri, and Sarwate 2013) first proposed the gradient perturbation method. Based on these works, (Kifer, Smith, and Thakurta 2012; Bassily, Smith, and Thakurta 2014; Abadi et al. 2016; Wang, Ye, and Xu 2017; Zhang et al. 2017; Wu et al. 2017; Bassily et al. 2019; Feldman, Koren, and Talwar 2020; Bassily et al. 2020) further improved the results under different assumptions.

Among the works mentioned above, some of them analyzed the privacy guarantees (Song, Chaudhuri, and Sarwate 2013; Abadi et al. 2016), some of them discussed the excess empirical risk bound (Wang, Ye, and Xu 2017; Zhang et al. 2017; Wu et al. 2017). Some works discussed the excess population risk under expectation, from different points of view, such as complexity theory, optimization theory, and stability theory: (Kifer, Smith, and Thakurta 2012) achieved an $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ excess population risk bound via complexity theory under expectation condition; (Bassily, Smith, and Thakurta 2014) achieved similar expected excess population risk bound under convexity assumption, via optimization theory; (Wang, Chen, and Xu 2019) proposed an $\mathcal{O}\left(\frac{p}{\log(n)\epsilon^2}\right)$ excess population risk bound under non-convex condition in expectation, via Langevin Dynamics method (Gelfand and Mitter 1991) and the stability of Gibbs algorithm; and (Feldman, Koren, and Talwar 2020) gives expected population risk bound of the order $\mathcal{O}\left(\frac{1}{n} + \frac{p}{n^2}\right)$, under strongly convex condition.

However, the behavior of the algorithm within a single or few runs cannot be well captured by expectation bounds, which is related to the probabilistic nature. In addition, in practical applications such as deep learning, it is often the case that the algorithm runs only once since the training process may take a long time. Therefore, obtaining a high probability bound is essential to ensure the performance of the algorithm on a single or few runs. So we focus on the high probability bound in this paper. Meanwhile, we concentrate on stability theory. Among many notions of stability, uniform stability is arguably the most popular one, which yields exponential generalization bounds. Via uniform stability, the high probability excess population risk bounds under different assumptions given by previous works all contain an $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ term, details can be found in Table 1. The reason is that when analyzing the generalization error, the technical routes follow works (Bousquet and Elisseeff 2002; Hardt, Recht, and Singer 2016). Besides, when analyzing the stability, previous works always do not consider the injected noise (e.g. (Wang et al. 2021)), or assume the random noise injected into adjacent datasets is the same. However, this is not reasonable, because 'adjacent dataset' is also the basis of DP. With the same random noise, it is hard to say that 'DP is guaranteed'.

In this paper, we consider the random noise injected into adjacent datasets and analyze the stability under the noisy version. By introducing the *Generalized Bernstein condition* (Koltchinskii 2006), we remove the $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ term when combining the stability and the generalization error, and further improve the excess population risk bound of differentially private models. The improved convergence rate is up to $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$, which positively answers the question given in Section 1: Can the high probability excess population risk bound achieve $\mathcal{O}\left(1/n\right)$ w.r.t $n$. The improvements are shown in Table 1. For 'TYPE', E. means expectation bound and H.P. means the high probability bound.

Table 1 first shows that by adding more assumptions, we achieve a better high probability excess population risk

Table 1: Previous excess population risk bounds and ours under different assumptions

| | ASSUMPTIONS | METHOD | UTILITY BOUND | TYPE |
|---|---|---|---|---|
| (BASSILY ET AL. 2019) | LIPSCHITZ, SMOOTH, CONVEX | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon} + \frac{1}{\sqrt{n}}\right)$ | E. |
| (FELDMAN, KOREN, AND TALWAR 2020) | LIPSCHITZ, CONVEX | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon} + \frac{1}{\sqrt{n}}\right)$ | E. |
| (FELDMAN, KOREN, AND TALWAR 2020) | LIPSCHITZ, STRONGLY CONVEX | GRADIENT | $\mathcal{O}\left(\frac{d}{n^2\epsilon^2} + \frac{1}{n}\right)$ | E. |
| (BASSILY ET AL. 2020) | LIPSCHITZ, CONVEX | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon} + \frac{1}{\sqrt{n}}\right)$ | H.P. |
| (WANG ET AL. 2021) | $\alpha$-HÖLDER SMOOTH, CONVEX | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon} + \frac{1}{\sqrt{n}}\right)$ | H.P. |
| (WANG ET AL. 2021) | $\alpha$-HÖLDER SMOOTH, CONVEX | OUTPUT | $\mathcal{O}\left(\frac{\sqrt[4]{p}}{\sqrt{n}\epsilon}\right)$ | H.P. |
| OURS | LIPSCHITZ, SMOOTH, PL CONDITION | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ | H.P. |
| OURS | $\alpha$-HÖLDER SMOOTH, PL CONDITION | GRADIENT | $\mathcal{O}\left(\frac{\sqrt[4]{p}}{n^{\frac{\alpha}{1+2\alpha}}\epsilon^{\frac{1}{2}}}\right)$ | H.P. |
| OURS (M-NGP) | $\alpha$-HÖLDER SMOOTH, PL CONDITION | GRADIENT | $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ | H.P. |

bound, $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$, which is state-of-the-art to the best of our knowledge. Then, we relax the assumptions and achieve $\mathcal{O}\left(\frac{\sqrt[4]{p}}{n^{\frac{\alpha}{1+2\alpha}}\epsilon^{\frac{1}{2}}}\right)$ high probability bound, but it cannot achieve the same bound ($\mathcal{O}\left(1/n\right)$ w.r.t $n$) under the condition that the loss function is Lipschitz, smooth, and satisfies PL condition. To overcome this problem, we propose an algorithm called m-NGP, and achieve the $\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ result under the same assumptions: $\alpha$-Hölder smooth and PL condition.

Moreover, although it is hard to directly compare the PL condition with convexity, PL condition can be applied to many non-convex conditions (more information can be found in Section 4.2). Besides, PL condition is weaker compared with strongly convex condition, and one of the best population risks under strongly convex condition is $\mathcal{O}\left(\frac{1}{n} + \frac{p}{n^2}\right)$ (Feldman, Koren, and Talwar 2020) (line 2 in Table 1). However, the result is an expectation one, different from ours. In this paper, we analyze the excess population risk bound of DP algorithm under high probability and PL cases, different from previous scenarios.

## 3. Preliminaries

In this paper, we assume that there are $n$ data instances in dataset $D$, i.e. $D = \{z_1, \cdots, z_n\}$ where $z = (x, y)$ with input $x \in \mathcal{X}$ and label $y \in \mathcal{Y}$, and $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. The data space is denoted by $\mathcal{D}$ and the parameter space is denoted by $\mathcal{C}$, the loss function $\ell$ is defined as $\ell(\cdot, \cdot) : \mathcal{D} \times \mathcal{C} \to \mathbb{R}$. Databases $D, D' \in \mathcal{D}^n$ differing by one data instance are denoted as $D \sim D'$, called *adjacent databases*. For a given vector $\mathbf{x} = [x_1, \cdots, x_d]^T$, its $\ell_p$-norm is $\|\mathbf{x}\|_p = (\sum_{i=1}^d |x_i|^p)^{\frac{1}{p}}$. And $A \lesssim B$ represents that there exists some constant $c > 0$, $A \leq cB$.

**Definition 1** (Differential Privacy (Dwork et al. 2006)). *A randomized algorithm: $\mathcal{A} : \mathcal{D}^n \to \mathbb{R}^p$ is $(\epsilon, \delta)$-differential privacy (DP) if for all $D \sim D'$ and events $S \in range(\mathcal{A})$:*

$$\mathbb{P}\left[\mathcal{A}(D) \in S\right] \leq e^\epsilon \mathbb{P}\left[\mathcal{A}(D') \in S\right] + \delta.$$

Definition 1 implies that the adversaries cannot infer whether an individual participates when training the machine learning model, because essentially the same distributions will be drawn over any adjacent datasets. Some kind of attacks, such as membership inference attack, attribute inference attack, and memorization attack, can be thwarted by DP (Backes et al. 2016; Jayaraman and Evans 2019; Carlini et al. 2019).

Throughout this paper, we focus on gradient perturbation method to guarantee $(\epsilon, \delta)$-DP, the paradigm is based on gradient descent: at iteration $t$,

$$\hat{\theta}_t \leftarrow \hat{\theta}_{t-1} - \eta_t \left(\nabla_\theta R_n(\hat{\theta}_{t-1}) + b_t\right), \quad (1)$$

where $\eta_t$ is the learning rate at iteration $t$, $b_t$ is the random noise injected into the gradient, $\hat{\theta}$ is corresponding model with privacy, and $R_n(\theta)$ is the empirical risk, defined as $R_n(\theta) := \frac{1}{n}\sum_{i=1}^n \ell(z_i, \theta)$.

In this paper, we focus on minimizing the population risk $R(\theta) = \mathbb{E}_{z\sim\mathcal{D}}\left[\ell(z, \theta)\right]$. In the setting of differential privacy, the excess population risk is defined by $R(\hat{\theta}) - \min_{\theta\in\mathcal{C}} R(\theta)$, which can be decomposed into:

$$R(\hat{\theta}_n) - \min_{\theta\in\mathcal{C}} R(\theta)$$
$$= \underbrace{R(\hat{\theta}_n) - R_n(\hat{\theta}_n)}_{A} + \underbrace{R_n(\hat{\theta}_n) - R_n(\theta^*)}_{B} \quad (2)$$
$$+ R_n(\theta^*) - R(\theta^*),$$

where $\theta^* = \arg\min_{\theta\in\mathcal{C}} R(\theta), \theta_n^* = \arg\min_{\theta\in\mathcal{C}} R_n(\theta)$. In (2), part $A$ is exactly the generalization. Via the definition of $\theta_n^*$, we have $R_n(\hat{\theta}_n) - R_n(\theta^*) \le R_n(\hat{\theta}_n) - R_n(\theta_n^*)$, which bounds part $B$ by the optimization error (also called the excess empirical risk). In this way, we answer the question mentioned in Section 1: Why generalization error is an important step towards the excess population risk.

To get the generalization error, algorithm stability theory is a popular tool, we introduce uniform stability and uniform argument stability here.

**Definition 2** (Uniform Stability (Bousquet and Elisseeff 2002)). *An algorithm $\theta_n$ is $\gamma$-uniformly stable if for any $S = \{z_1, \cdots, z_i, \cdots, z_n\}$ and $S' = \{z_1, \cdots, z_i', \cdots, z_n\}$, where $i = 1, \cdots, n$, it holds that*

$$|\ell(z, \theta_n(S)) - \ell(z, \theta_n(S'))| \le \gamma.$$

In this paper, we use notation $\theta_n(S)$ for both algorithm and model parameter. By Definition 2, it is easy to follow that the uniform stability measures the upper bound of the difference (on the loss function) between the models derived from adjacent datasets.

**Definition 3** (Uniform Argument Stability (Bassily et al. 2020)). *Algorithm $\theta_n$ is $\gamma$-uniformly argument stable if for any $S = \{z_1, \cdots, z_i, \cdots, z_n\}$ and $S' = \{z_1, \cdots, z_i', \cdots, z_n\}$, where $i = 1, \cdots, n$, it holds that*

$$\|\theta_n(S) - \theta_n(S')\|_2 \le \gamma.$$

Definition 3 shows that the uniform argument stability measures the upper bound of the difference (on the model parameter) between the models derived from adjacent datasets.

Furthermore, we introduce some assumptions.

**Assumption 1** ($G$-Lipschitz). *The loss function $\ell : \mathcal{D}\times\mathcal{C} \to \mathbb{R}$ is $G$-Lipschitz over $\theta$ if for any $z \in \mathcal{D}$ and $\theta_1, \theta_2 \in \mathcal{C}$, we have: $|\ell(z, \theta_1) - \ell(z, \theta_2)| \le G\|\theta_1 - \theta_2\|_2$.*

With Assumption 1, one can easily get that if the loss function is $G$-Lipchitz, then $\gamma$-uniformly argument stability implies $G\gamma$-uniformly stability.

**Assumption 2** ($L$-smooth). *The loss function $\ell : \mathcal{D}\times\mathcal{C} \to \mathbb{R}$ is $L$-smooth over $\theta$ if for any $z \in \mathcal{D}$ and $\theta_1, \theta_2 \in \mathcal{C}$, we have: $\|\nabla_\theta\ell(z, \theta_1) - \nabla_\theta\ell(z, \theta_2)\|_2 \le L\|\theta_1 - \theta_2\|_2$.*

If loss function $\ell(\cdot, \cdot)$ is differentiable, smoothness yields: $\ell(z, \theta_1) - \ell(z, \theta_2) \le \langle\nabla_\theta\ell(z, \theta_2), \theta_1 - \theta_2\rangle + \frac{L}{2}\|\theta_1 - \theta_2\|_2^2$.

Assumptions $G$-Lipschitz and $L$-smooth are commonly used in the utility analysis of DP machine learning (Chaudhuri, Monteleoni, and Sarwate 2011; Kifer, Smith, and Thakurta 2012; Abadi et al. 2016; Bassily et al. 2019; Feldman, Koren, and Talwar 2020; Bassily et al. 2020). To relax the Lipschitz and smoothness assumptions, we introduce the $\alpha$-Hölder smoothness of the loss function:

**Assumption 3** ($\alpha$-Hölder smooth). *Let $\alpha \in (0, 1]$. The loss function $\ell : \mathcal{D} \times \mathcal{C} \to \mathbb{R}$ is $\alpha$-Hölder smooth over $\theta$ with parameter $H$ if for any $z \in \mathcal{D}$ and $\theta_1, \theta_2 \in \mathcal{C}$, we have: $\|\nabla_\theta\ell(z, \theta_1) - \nabla_\theta\ell(z, \theta_2)\|_2 \le H\|\theta_1 - \theta_2\|_2^\alpha$.*

**Lemma 1** ((Ying and Zhou 2017)). *If the loss function $\ell(\cdot, \cdot)$ is differentiable, then Assumption 3 yields $\ell(z, \theta_1) - \ell(z, \theta_2) \le \langle\nabla_\theta\ell(z, \theta_2), \theta_1 - \theta_2\rangle + \frac{H}{\alpha+1}\|\theta_1 - \theta_2\|_2^{\alpha+1}$.*

By the definition, it is easy to follow that if $\alpha = 1$, it is equivalent to $H$-smooth; and if $\alpha \to 0$, it satisfies the Lipschitz property given in Assumption 1. Besides, with bounded parameter space, i.e. $\|\mathcal{C}\|_2 \le M_\mathcal{C}$, $\alpha$-Hölder smoothness immediately implies $\max\{2HM_\mathcal{C}, H\}$-Lipschitz. Moreover, Assumption 3 instantiates many nonsmooth loss functions. For example, the $q$-norm hinge loss $\ell(z, \theta) = (\max(0, 1 - y\langle\theta, z\rangle))^q$ for classification and the $q$-th power absolute distance loss $\ell(z, \theta) = |y - \langle\theta, z\rangle|^q$ for regression (Lei and Ying 2020a), whose $\ell$ are $(q-1)$-Hölder smooth if $q \in (1, 2]$ (Li and Liu 2021). Lemma 1 shows that Hölder smoothness shares similar property with smoothness defined in Assumption 2.

# 4. Sharper Utility Bounds for Differentially Private Models

## 4.1. Privacy Guarantees

Before analyzing the excess population risk bound, we first discuss the privacy guarantees in this section. (Abadi et al. 2016) proposed the moments accountant method to measure the privacy costs of DP model training by stochastic gradient descent (SGD), (Wang, Ye, and Xu 2017) further analyzed it under the setting of gradient descent (GD). In this paper, we focus more on the utility analysis, to improve the excess population risk, so we directly apply it to the gradient perturbation method.

**Lemma 2** ((Wang, Ye, and Xu 2017)). *In gradient perturbation method in (1), if Assumption 1 holds, then for $\epsilon, \delta > 0$, it is $(\epsilon, \delta)$-DP if the Gaussian random noise $b_t \sim \mathcal{N}(0, \sigma^2 I_p)$, and for some constant $c$,*

$$\sigma^2 = c\frac{G^2 T\log(1/\delta)}{n^2\epsilon^2}.$$

**Remark 1.** *Lemma 2 only assumes the loss function to be $G$-Lipschitz. If we assume that $\ell(\cdot, \cdot)$ is $\alpha$-Hölder smooth with parameter $H$, then $G$ can be replaced by $\max\{2HM_\mathcal{C}, H\}$ as discussed above.*

## 4.2. Analysis of the excess population risk

Before analyzing the excess population risk, we first introduce some assumptions.

Most of the previous works assumed that the loss function is convex (or strongly convex) when analyzing the empirical and population risks. In this paper, we use the Polyak-Łojasiewicz (PL) condition to replace convexity.

**Assumption 4** (Polyak-Łojasiewicz condition). *Function $f(\theta)$ satisfies the Polyak-Łojasiewicz (PL) condition if there exists $\mu > 0$ and for every $\theta$,*

$$\|\nabla_\theta f(\theta)\|_2^2 \ge 2\mu\left(f(\theta) - f(\theta^*)\right),$$

*where $\theta^* = \arg\min_{\theta\in\mathcal{C}} f(\theta)$.*

In this paper, we assume the empirical risk and the population risk both satisfy the PL condition[1].

---

[1] PL condition can be directly derived from strongly convex (Karimi, Nutini, and Schmidt 2016), so all the results given in this paper hold when it comes to the strongly convex conditions.

The Polyak-Łojasiewicz condition is one of the weakest curvature conditions (Karimi, Nutini, and Schmidt 2016; Li and Liu 2021), weaker than 'one-point convexity' (Kleinberg, Li, and Yuan 2018), 'star convexity' (Zhou et al. 2019), and 'quasar convexity' (Hinder, Sidford, and Sohoni 2020). It is widely used in the analysis of non-convex learning (Wang, Ye, and Xu 2017; Charles and Papailiopoulos 2018; Lei and Ying 2020b; Lei and Tang 2021) and many popular non-convex objective functions satisfy the PL condition, such as: matrix factorization (Liu, Wu, and So 2016), robust regression (Liu, Wu, and So 2016), neural networks with one hidden layer (Li and Yuan 2017), mixture of two Gaussians (Balakrishnan, Wainwright, and Yu 2017), ResNets with linear activations (Hardt and Ma 2017), linear dynamical systems (Hardt, Ma, and Recht 2018), phase retrieval (Sun, Qu, and Wright 2018), and blind deconvolution (Li et al. 2019).

**Remark 2.** *(Karimi, Nutini, and Schmidt 2016) shows that the Polyak-Łojasiewicz inequality directly implies the following inequality: $R_n(\theta) - R_n(\theta_n^*) \geq \frac{\mu}{2}\|\theta - \theta_n^*\|_2^2$, which is called the Quadratic Growth (QG) condition. In the following, we use QG to bound the argument stability of the gradient perturbation algorithm.*

Then, with Assmption 4, we discuss the uniform argument stability of the private model.

**Lemma 3.** *In gradient perturbation method (1), if Assumptions 1, 4 hold, then*

$$\left\|\hat{\theta}_n - \hat{\theta}_n'\right\|_2 \leq 2\sqrt{2}\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G}{\mu n},$$

*where $\hat{\theta}_n$ and $\hat{\theta}_n'$ denote the models derived from adjacent datasets $S$ and $S'$.*

The proof is given in Appendix A.1. By the QG condition (implied by PL inequality) discussed in Remark 2, Lemma 3 connects the uniform argument stability with the empirical risk. And as a result, we only need to consider the noise when analyzing the empirical risk. Besides, when it comes to the stability of DP models, previous results often assume that the noise added to the gradient in each iteration is the same for adjacent datasets $S$ and $S'$ (e.g. (Wang et al. 2021)). This is not reasonable because noise injection is an independent process, so we expand it in this paper.

And to remove the $\mathcal{O}(1/\sqrt{n})$ term in previous results, we further need the Generalized Bernstein condition when analyzing the excess population risk.

**Assumption 5** (Generalized Bernstein condition (Koltchinskii 2006)). *We say the loss function $\ell$ satisfies the generalized Bernstein condition if for some $B > 0$ for any $\theta \in \mathcal{C}$, we have:*

$$\mathbb{E}\left[(\ell(z, \theta) - \ell(z, \theta^*))^2\right] \leq B\left(R(\theta) - R(\theta^*)\right).$$

**Remark 3.** *Here, we discuss the connections between Assumptions 1, 4 and 5. Via Assumption 1, we have $\mathbb{E}\left[(\ell(z, \theta) - \ell(z, \theta^*))^2\right] \leq G^2\|\theta - \theta^*\|_2^2$. And Remark 2 shows that if $R(\theta)$ satisfies the Polyak-Łojasiewicz condition, we have $\frac{\mu}{2}\|\theta - \theta^*\|_2^2 \leq R(\theta) - R(\theta^*)$. Combining*

*these inequalities together, we observe that if Assumptions 1 and 4 hold, then the loss function $\ell(\cdot)$ satisfies Assumption 5 with parameter $B = \frac{2G^2}{\mu}$.*

With the stability of the private model and Assumption 5, now we come to the excess population risk.

**Theorem 1.** *If Assumptions 1, 2, and 4 hold, the loss function is bounded, i.e. $0 \leq \ell(\cdot, \cdot) \leq M_\ell$, taking $\sigma$ given by Lemma 2, $T = \mathcal{O}(\log(n))$, $\eta_1 = \cdots = \eta_T = \frac{1}{L}$, if $\zeta \in (\exp(-p/8), 1)$, then with probability at least $1 - \zeta$:*

$$R(\hat{\theta}_n) - R(\theta^*)$$
$$\leq c_1 \frac{G\log^{1.5}(n)\sqrt{p\log(1/\delta)}}{n\epsilon}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)$$
$$+ c_2 \frac{G^2 p\log(n)\log(1/\delta)}{n^2\epsilon^2}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2$$
$$+ c_3 \frac{\log(n)}{n}.$$

*for some constants $c_1, c_2, c_3 > 0$.*

Detailed proof can be found in Appendix A.3, we give a proof sketch here. First, by Lemma 3 and Lipschitzness, we get the uniform stability of gradient perturbation (1). Then, we analyze the generalization error via stability theory. By novel decomposition method, via Assumption 5 and its moments bound, we couple term $R(\hat{\theta}_n) - R_n(\hat{\theta}_n)$ and term $R_n(\theta^*) - R(\theta^*)$ in (2) together, to remove the $\mathcal{O}(1/\sqrt{n})$ term in the generalization error. In this way, a better excess population risk bound is achieved. The proof is motivated by (Klochkov and Zhivotovskiy 2021) in the non-private case. The key difference is that in the setting of DP, random noise is injected into the algorithm. In (Klochkov and Zhivotovskiy 2021), a key step to analyze the generalization error is summing $X_i = \mathbb{E}'[\ell(z_i, \theta_n') - \ell(z_i, \theta^*)]$ for $i = 1, \cdots, n$, where $\theta_n'$ is derived from an independent copy of the original dataset and $\mathbb{E}'$ means the expectation over the independent copy. When summing, $X_i$ is required to be zero mean. However, in the cases of DP, if we replace $\theta_n'$ by $\hat{\theta}_n'$, then $X_i$ are not zero mean. Besides, for output perturbation, a common way to decompose the excess population risk is $R(\hat{\theta}_n) - R(\theta^*) \leq R(\hat{\theta}_n) - R(\theta_n) + R(\theta_n) - R_n(\theta_n) + R_n(\theta_n) - R_n(\theta_n^*) + R_n(\theta^*) - R(\theta^*)$, which naturally solves the problem mentioned above (the generalization error is discussed over the non-private model). However, when it comes to the gradient perturbation method, we cannot solve the problem in this way, because the random noise is coupled with the gradient. So, we decouple the noise terms and overcome the challenge by the moment Bernstein inequality.

**Remark 4.** *If omitting $\log(\cdot)$ and constant terms, the excess population risk bound comes to:*

$$\mathcal{O}\left(\frac{p}{n^2\epsilon^2} + \frac{\sqrt{p}}{n\epsilon} + \frac{1}{n}\right) = \mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right).$$

*The result is better than previous ones containing $\mathcal{O}(\frac{1}{\sqrt{n}})$ terms. And it is the first high probability excess population*

*risk bound over DP algorithm overcoming the $\mathcal{O}(n^{-1/2})$ bottleneck, to the best of our knowledge.*

Then, we replace Assumptions 1 ($G$-Lipschitz) and 2 ($L$-smooth) by Assumption 3 ($\alpha$-Hölder smooth).

**Theorem 2.** *If Assumptions 3 and 4 hold, the loss function and the parameter space are bounded, i.e. $0 \le \ell(\cdot, \cdot) \le M_\ell$, $\|\mathcal{C}\|_2 \le M_\mathcal{C}$. Taking $\sigma$ given by Lemma 2, $T = \mathcal{O}\left(n^{\frac{2}{1+2\alpha}}\right)$, and $\eta_t = \frac{2}{\mu(t+\kappa)}$, where $\kappa \ge \frac{2H^{1/\alpha}}{\mu}$, if $\zeta \in (\exp(-p/8), 1)$, then with probability at least $1 - \zeta$:*

$$R(\hat{\theta}_n) - R(\theta^*)$$

$$\le c_1 \frac{G' \log(n) \sqrt[4]{p \log(1/\delta)}}{n^{\frac{\alpha}{1+2\alpha}} \epsilon^{1/2}} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^{\frac{1}{2}}$$

$$+ c_2 \frac{G'^2 \sqrt{p \log(1/\delta)}}{n^{\frac{2\alpha}{1+2\alpha}} \epsilon} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)$$

$$+ c_3 \frac{\log(n)}{n},$$

*for some constants $c_1, c_2$ and $c_3$, where $G' = \max\{2HM_\mathcal{C}, H\}$.*

Detailed proof can be found in Appendix A.4. The proof is similar to Theorem 1. Considering that the stability parameter given by Lemma 3 is related to the optimization error, so the key is to obtain the optimization error. The challenge is that the properties $G$-Lipschitz and $L$-smooth are replaced by the assumption $\alpha$-Hölder smooth when analyzing it. To overcome the challenge, we use Lemma 1 to bound the optimization error and Young's inequality is used to normalize the exponential rate. By choosing proper learning rate, we get an acceptable excess population risk bound. Details are shown in the proof of Lemma 7.

By Theorem 2, it is easy to follow that with high probability, the excess population risk satisfies:

$$R(\hat{\theta}_n) - R(\theta^*) = \mathcal{O}\left(\frac{p^{1/4}}{\epsilon^{1/2}} n^{\frac{-\alpha}{1+2\alpha}}\right).$$

**Remark 5.** *By the definition of $\alpha$-Hölder smooth, we have $\alpha \in (0, 1]$, so our result is worse than the gradient based result given by (Wang et al. 2021) ($\mathcal{O}(1/\sqrt{n})$ w.r.t $n$). One of the reasons is that (Wang et al. 2021) does not consider the noise injected into the model when analyzing stability. However, as discussed before, the noise addition is independent to datasets, so we expand the stability to the noisy version in this paper, which generalize previous settings. Besides, our result is of $\mathcal{O}(\epsilon^{-1/2})$ w.r.t $\epsilon$, and previous results are of the order $\mathcal{O}(\epsilon^{-1})$. In practice, $\epsilon$ is always set less than 1 to guarantee meaningful privacy, so our result is more superior when it comes to conditions that privacy requirements are strict (low $\epsilon$ conditions).*

Via the discussion mentioned above, we observe that the result given in Theorem 2 is far worse than which given in Theorem 1, if we replace Lipschitzness and smoothness by $\alpha$-Hölder smoothness. The reason is that when applying

---

**Algorithm 1** $\max\{\mathbf{1}, \mathbf{g}\}$-Normalized Gradient Perturbation

1: **Input:** Dataset $D$, learning rate at iteration $t$: $\eta_t$, the variance of the Gaussian noise injected to the gradient: $\sigma$.
2: Initialize $\theta_0$.
3: **for** $t = 0$ **to** $T - 1$ **do**
4:     $G_t \leftarrow R_n(\hat{\theta}_t) + b_t, b_t \sim \mathcal{N}\left(0, \sigma^2 I_p\right)$.
5:     **if** $\|G_t\|_2 < 1$ **then**
6:         $G_t \leftarrow G_t / \|G_t\|_2$.
7:     **end if**
8:     $\hat{\theta}_{t+1} \leftarrow \hat{\theta}_t - \eta_t G_t$.
9: **end for**
10: Return $\hat{\theta}_n = \hat{\theta}_T$.

---

Young's inequality in the optimization error analysis, an additional term $\frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)}$ appears, leading a loose excess population risk bound.

Motivated by this, we design a variant of gradient perturbation method, called **$\max\{\mathbf{1}, \mathbf{g}\}$-Normalized Gradient Perturbation** DP algorithm, to overcome the loose excess population risk bound. Details are shown in Algorithm 1.

**Remark 6.** *The difference between Algorithm 1 and (1) is that in lines 5 and 6, we normalize the $\ell_2$-norm of the gradient to 1 if it is less than 1. In this way, we can 'bypass' the Young's inequality when scaling $\|\theta_t - \theta_n^*\|_2^{1+\alpha}$ (derived from Lemma 1), further remove term $\frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)}$ in the theoretical analysis. Details can be found in Appendix A.4.*

Then, we improve the excess population risk bound given in Theorem 2.

**Theorem 3.** *If Assumptions 3, 4 hold, the loss function and the parameter space are bounded, i.e. $0 \le \ell(\cdot, \cdot) \le M_\ell$, $\|\mathcal{C}\|_2 \le M_\mathcal{C}$. Taking $\sigma$ given by Lemma 2, $T = \mathcal{O}\left(\log(n)\right)$, and $\eta_1 = \cdots = \eta_T = \eta$, where $\eta = \left(\frac{1}{H}\right)^{1/\alpha}$, if $\zeta \in (\exp(-p/8), 1)$, then with probability at least $1 - \zeta$,*

$$R(\hat{\theta}_n) - R(\theta^*)$$

$$\le c_1 \frac{G' \log^{1.5}(n) \sqrt{p \log(1/\delta)}}{n\epsilon} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)$$

$$+ c_2 \frac{G'^2 \log(n) p \log(1/\delta)}{n^2 \epsilon^2} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2$$

$$+ c_3 \frac{\log(n)}{n},$$

*for some constants $c_1, c_2, c_3 > 0$, where $G' = \max\{2HM_\mathcal{C}, H\}$.*

Detailed proof is given in Appendix A.5. The proof is similar to Theorems 1 and 2, the key difference is that by *gradient normalization*, Young's inequality is abandoned in the theoretical analysis (as discussed in Remark 6), which implies a better excess population risk bound.

**Remark 7.** *We introduce normalization from theoretical view, and the experiments show that it works in practice (de-*

*tails can be found in Section 5). Here, we provide some intuitions on why normalization works. For gradient perturbation method, it is easy to observe that the random noise $b_t$ is sampled independently at each iteration $t$, so it is reasonable to suppose $b_t$ does not change too much when it comes to different iterations. As a result, if the $\ell_2$-norm of the gradient is too small, the random noise will play a more important role and the gradient property will be concealed beneath the noise. Besides, the random noise itself would help the generalization error (Li, Luo, and Qiao 2020). So in m-NGP, we apply normalization to scale the $\ell_2$-norm of the gradient to 1 if it is less than 1, strengthens the gradient along with the random noise.*

By Theorem 3, it is easy to follow that with high probability, $R(\hat{\theta}_n) - R(\theta^*) = \mathcal{O}\big(\frac{\sqrt{p}}{n\epsilon}\big)$. The bound is of the same order as the result given in Theorem 1. This is also the first $\mathcal{O}(1/n)$ high probability excess population risk bound over DP algorithm w.r.t $n$ without smoothness.

**Remark 8.** *Theorems 2 and 3 require Assumptions 3 and 4, in this remark, we give some examples satisfy these assumptions. As discussed in Section 3, the q-norm hinge loss and q-th power absolute distance loss can be seemed as squared piecewise-linear loss functions when $q = 2$, so they satisfy Hölder smoothness. For one-layer neural networks with squared error loss and leaky ReLU activations, the same phenomenon holds because the neural network can be seemed as a matrix multiplication. Beside, (Charles and Papailiopoulos 2018) shows that several interesting machine learning setups satisfy Assumption 4, such as 1-layer neural networks with a squared error loss and leaky ReLU activations, loss functions of least squares minimizations, squared piecewise-linear functions with regularized term, etc. As a result, loss functions sataisfy those assumptions including but not limited to: (1) logistic regression and least squared minimization; (2) some of the squared piecewise linear functions; (3) some of the neural networks such as one-layer neural networks with squared error loss and leaky ReLU activations. The examples listed above are only part of the loss functions who satisfy those assumptions. We do not only focus on least squared minimization and logistic regression models, but extend the condition from smoothness to Hölder smoothness, and from strongly convex (convex) to the PL condition (some of the non-convex cases).*

## 5. Experiments

In this section, we perform experiments on real datasets to evaluate our proposed our proposed m-NGP algorithm.

The experiments are performed on classification task over datasets Iris (Dua and Graff 2017), Breast Cancer (Mangasarian and Wolberg 1990), Credit Card Fraud (Bontempi and Worldline 2018), Bank (Moro, Cortez, and Rita 2014), and Adult (Dua and Graff 2017), the number of total data instances are 150, 699, 984, 41188, and 45222, respectively. We split the training and testing sets randomly and evaluate the accuracy on the testing set and the convengence rate on the training set. In all the experiments, the privacy budget $\delta$ is set $\frac{1}{n}$ and we choose $\epsilon = 0.1$ to 1.0.
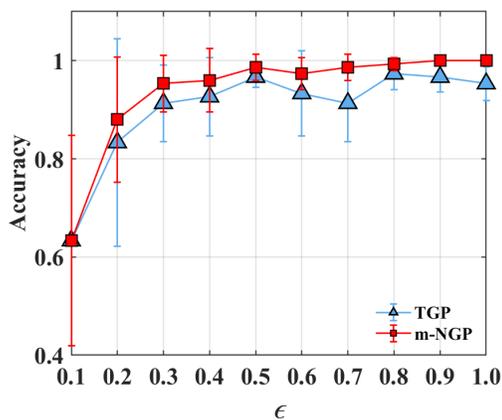
We apply regularized logistic regression to the classification task, which satisfies the assumptions mentioned before, the experimental results are shown in Figure 1. We show the results over datasets Iris and Adult in this section, experiments on other datasets are shown in Appendices B.1 and B.2. For convergence rate, the shadow area represents the maximum and minimum loss over mutiple experiments, reflecting the variance. The shadow area in part (d) of Figure 1 is not obvious, the reason is that the variances are small. Over most datasets, the accuracy and the convergence rate of m-NGP is better than traditional gradient perturbation method, which is in line with the theoretical analysis. Moreover, in Appendix B.3, we perform experiments to demonstrate the effects brought by the dimension parameter $p$, the experimental results follow the theoretical results: with increasing $p$, the accuracy becomes worse in general.
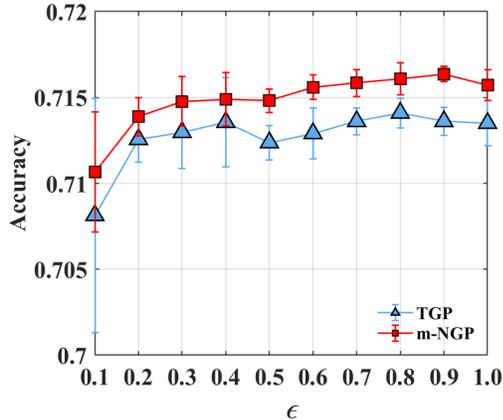
## 6. Conclusions

In this paper, we first propose a state-of-the-art $\mathcal{O}\big(\frac{\sqrt{p}}{n\epsilon}\big)$ high probability excess population risk bound for gradient perturbation based DP algorithms, under the assumptions of $G$-Lipschitz, $L$-smooth, and Polyak-Łojasiewicz condition. The result positively answers the open problem: *Can we achieve high probability excess risk bound with rate $\mathcal{O}(1/n)$ w.r.t $n$ for DP models via uniform stability?* Then, we extend the result to a more general case, requiring $\alpha$-Hölder smoothness and Polyak-Łojasiewicz condition. However, the result is not as satisfactory as before, we achieve an $\mathcal{O}\big(n^{\frac{-\alpha}{1+2\alpha}}\big)$ high probability excess population risk bound, which cannot achieve an $\mathcal{O}(1/n)$ bound. To get a better result, we further propose a new algorithm: max$\{1, g\}$-Normalized Gradient Perturbation (m-NGP). Detailed theoretical analysis shows that m-NGP can achieve $\mathcal{O}\big(\frac{\sqrt{p}}{n\epsilon}\big)$ high probability excess population risk bound, under the assumptions of $\alpha$-Hölder smoothness and Polyak-Łojasiewicz condition, which is the first $\mathcal{O}(1/n)$ high probability bound w.r.t $n$ under non-smoothness differentially private cases. Experimental results show that the accuracy of m-NGP algorithm is better than traditional gradient perturbation method. Thus, our proposed max$\{1, g\}$-Normalized Gradient Perturbation method improves the excess population risk bound and the accuracy of the DP model over real datasets, simultaneously.
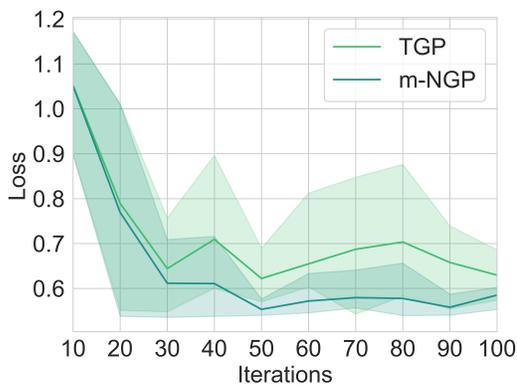
## References

[Abadi et al. 2016] Abadi, M.; Chu, A.; Goodfellow, I.; McMahan, H. B.; Mironov, I.; Talwar, K.; and Zhang, L. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 308–318.

[Backes et al. 2016] Backes, M.; Berrang, P.; Humbert, M.; and Manoharan, P. 2016. Membership privacy in micrornabased studies. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 319–330.

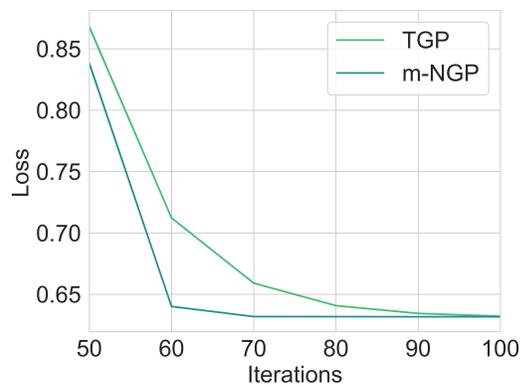[Balakrishnan, Wainwright, and Yu 2017] Balakrishnan, S.; Wainwright, M. J.; and Yu, B. 2017. Statistical guaran-

(a) Iris

(b) Adult

(c) Iris

(d) Adult

Figure 1: Comparisons between Traditional Gradient Perturbation (TGP) and max$\{\mathbf{1}, \mathbf{g}\}$-Normalized Gradient Perturbation (m-NGP).

tees for the em algorithm: From population to sample-based analysis. *The Annals of Statistics* 77 – 120.

[Bartlett, Bousquet, and Mendelson 2002] Bartlett, P. L.; Bousquet, O.; and Mendelson, S. 2002. Localized rademacher complexities. In *Computational Learning Theory, 15th Annual Conference on Computational Learning Theory, COLT 2002*, 44–58.

[Bassily et al. 2019] Bassily, R.; Feldman, V.; Talwar, K.; and Guha Thakurta, A. 2019. Private stochastic convex optimization with optimal rates. In *Advances in Neural Information Processing Systems*, 11279–11288.

[Bassily et al. 2020] Bassily, R.; Feldman, V.; Guzmán, C.; and Talwar, K. 2020. Stability of stochastic gradient descent on nonsmooth convex losses. In *Advances in Neural Information Processing Systems*, 4381–4391.

[Bassily, Smith, and Thakurta 2014] Bassily, R.; Smith, A.; and Thakurta, A. 2014. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, 464–473.

[Bernstein and Sheldon 2019] Bernstein, G., and Sheldon, D. R. 2019. Differentially private bayesian linear regression. In *Advances in Neural Information Processing Systems*, 523–533.

[Bontempi and Worldline 2018] Bontempi, G., and Worldline. 2018. ULB the machine learning group.

[Boucheron, Lugosi, and Massart 2013] Boucheron, S.; Lugosi, G.; and Massart, P. 2013. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.

[Bousquet and Elisseeff 2002] Bousquet, O., and Elisseeff, A. 2002. Stability and generalization. *The Journal of Machine Learning Research* 499–526.

[Bousquet, Klochkov, and Zhivotovskiy 2020] Bousquet, O.; Klochkov, Y.; and Zhivotovskiy, N. 2020. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory*, 610–626.

[Bun, Elias, and Kulkarni 2021] Bun, M.; Elias, M.; and

Kulkarni, J. 2021. Differentially private correlation clustering. In *Proceedings of the 38th International Conference on Machine Learning*, 1136–1146.

[Carlini et al. 2019] Carlini, N.; Liu, C.; Erlingsson, Ú.; Kos, J.; and Song, D. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)*, 267–284.

[Charles and Papailiopoulos 2018] Charles, Z., and Papailiopoulos, D. 2018. Stability and generalization of learning algorithms that converge to global optima. In *Proceedings of the 35th International Conference on Machine Learning*, 745–754.

[Chaudhuri, Monteleoni, and Sarwate 2011] Chaudhuri, K.; Monteleoni, C.; and Sarwate, A. D. 2011. Differentially private empirical risk minimization. *Journal of Machine Learning Research* 1069–1109.

[Ding, Kulkarni, and Yekhanin 2017] Ding, B.; Kulkarni, J.; and Yekhanin, S. 2017. Collecting telemetry data privately. In *Advances in Neural Information Processing Systems*, 3571–3580.

[Dua and Graff 2017] Dua, D., and Graff, C. 2017. UCI machine learning repository.

[Dwork et al. 2006] Dwork, C.; McSherry, F.; Nissim, K.; and Smith, A. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, 265–284.

[Dwork, Roth, and others 2014] Dwork, C.; Roth, A.; et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 211–407.

[Erlingsson, Pihur, and Korolova 2014] Erlingsson, Ú.; Pihur, V.; and Korolova, A. 2014. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*, 1054–1067.

[Feldman, Koren, and Talwar 2020] Feldman, V.; Koren, T.; and Talwar, K. 2020. Private stochastic convex optimization: Optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, 439–449.

[Gelfand and Mitter 1991] Gelfand, S. B., and Mitter, S. K. 1991. Recursive stochastic algorithms for global optimization in rˆd. *SIAM Journal on Control and Optimization* 999–1018.

[Hardt and Ma 2017] Hardt, M., and Ma, T. 2017. Identity matters in deep learning. In *5th International Conference on Learning Representations, 2017*.

[Hardt, Ma, and Recht 2018] Hardt, M.; Ma, T.; and Recht, B. 2018. Gradient descent learns linear dynamical systems. *Journal of Machine Learning Research* 29:1–29:44.

[Hardt, Recht, and Singer 2016] Hardt, M.; Recht, B.; and Singer, Y. 2016. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of The 33rd International Conference on Machine Learning*, 1225–1234.

[Heikkilä et al. 2019] Heikkilä, M.; Jälkö, J.; Dikmen, O.; and Honkela, A. 2019. Differentially private markov chain monte carlo. In *Advances in Neural Information Processing Systems 32*. 4115–4125.

[Hinder, Sidford, and Sohoni 2020] Hinder, O.; Sidford, A.; and Sohoni, N. 2020. Near-optimal methods for minimizing star-convex functions and beyond. In *Proceedings of Thirty Third Conference on Learning Theory*, 1894–1938.

[Jayaraman and Evans 2019] Jayaraman, B., and Evans, D. 2019. Evaluating differentially private machine learning in practice. In *Proceedings of the 28th USENIX Conference on Security Symposium*, 1895–1912.

[Karimi, Nutini, and Schmidt 2016] Karimi, H.; Nutini, J.; and Schmidt, M. 2016. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 795–811.

[Kifer, Smith, and Thakurta 2012] Kifer, D.; Smith, A.; and Thakurta, A. 2012. Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, 25–1.

[Kleinberg, Li, and Yuan 2018] Kleinberg, B.; Li, Y.; and Yuan, Y. 2018. An alternative view: When does SGD escape local minima? In *Proceedings of the 35th International Conference on Machine Learning*, 2698–2707.

[Klochkov and Zhivotovskiy 2021] Klochkov, Y., and Zhivotovskiy, N. 2021. Stability and deviation optimal risk bounds with convergence rate $o(1/n)$. *arXiv preprint arXiv:2103.12024*.

[Koltchinskii 2006] Koltchinskii, V. 2006. Local rademacher complexities and oracle inequalities in risk minimization. *The Annals of Statistics* 2593–2656.

[Kulkarni et al. 2021] Kulkarni, T.; Jälkö, J.; Koskela, A.; Kaski, S.; and Honkela, A. 2021. Differentially private bayesian inference for generalized linear models. In *Proceedings of the 38th International Conference on Machine Learning*, 5838–5849.

[Lei and Tang 2021] Lei, Y., and Tang, K. 2021. Learning rates for stochastic gradient descent with nonconvex objectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

[Lei and Ying 2020a] Lei, Y., and Ying, Y. 2020a. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *Proceedings of the 37th International Conference on Machine Learning*, 5809–5819.

[Lei and Ying 2020b] Lei, Y., and Ying, Y. 2020b. Sharper generalization bounds for learning with gradient-dominated objective functions. In *International Conference on Learning Representations*.

[Li and Liu 2021] Li, S., and Liu, Y. 2021. Improved learning rates for stochastic optimization: Two theoretical viewpoints.

[Li and Yuan 2017] Li, Y., and Yuan, Y. 2017. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems 30, 2017*, 597–607.

[Li et al. 2019] Li, X.; Ling, S.; Strohmer, T.; and Wei, K. 2019. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Applied and computational harmonic analysis* 893–934.

[Li, Luo, and Qiao 2020] Li, J.; Luo, X.; and Qiao, M. 2020. On generalization error bounds of noisy gradient methods for non-convex learning. In *8th International Conference on Learning Representations*. OpenReview.net.

[Liu, Wu, and So 2016] Liu, H.; Wu, W.; and So, A. M. 2016. Quadratic optimization with orthogonality constraints: Explicit lojasiewicz exponent and linear convergence of line-search methods. In *Proceedings of the 33nd International Conference on Machine Learning, 2016*, 1158–1167.

[Mangasarian and Wolberg 1990] Mangasarian, O. L., and Wolberg, W. H. 1990. Cancer diagnosis via linear programming. Technical report, University of Wisconsin-Madison Department of Computer Sciences.

[McMillan 2016] McMillan, R. 2016. Apple tries to peek at user habits without violating privacy. *The Wall Street Journal*.

[Moro, Cortez, and Rita 2014] Moro, S.; Cortez, P.; and Rita, P. 2014. A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems* 22–31.

[Nguyen and Vullikanti 2021] Nguyen, D., and Vullikanti, A. 2021. Differentially private densest subgraph detection. In *Proceedings of the 38th International Conference on Machine Learning*, 8140–8151.

[Ponnusamy et al. 2021] Ponnusamy, V.; Christopher Clement, J.; Sriharipriya, K. C.; and Natarajan, S. 2021. *Smart Healthcare Technologies for Massive Internet of Medical Things*. Springer International Publishing. 71–101.

[Raginsky, Rakhlin, and Telgarsky 2017] Raginsky, M.; Rakhlin, A.; and Telgarsky, M. 2017. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. In *Proceedings of the 2017 Conference on Learning Theory*, 1674–1703.

[Shokri et al. 2017] Shokri, R.; Stronati, M.; Song, C.; and Shmatikov, V. 2017. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy, 2017*, 3–18.

[Singha et al. 2021] Singha, M. K.; Dwivedi, P.; Sankhe, G.; Patra, A.; and Rojwal, V. 2021. *Role of Sensors, Devices and Technology for Detection of COVID-19 Virus*. Springer International Publishing. 293–312.

[Song, Chaudhuri, and Sarwate 2013] Song, S.; Chaudhuri, K.; and Sarwate, A. D. 2013. Stochastic gradient descent with differentially private updates. In *2013 IEEE Global Conference on Signal and Information Processing*, 245–248.

[Sun, Qu, and Wright 2018] Sun, J.; Qu, Q.; and Wright, J. 2018. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics* 1131–1198.

[Swapna and Soman 2021] Swapna, G., and Soman, K. P. 2021. *Diabetes Detection and Sensor-Based Continuous Glucose Monitoring – A Deep Learning Approach*. Springer International Publishing. 245–268.

[Ullman and Sealfon 2019] Ullman, J., and Sealfon, A. 2019. Efficiently estimating erdos-renyi graphs with node differential privacy. In *Advances in Neural Information Processing Systems*, 3765–3775.

[Wang and Xu 2019] Wang, D., and Xu, J. 2019. Principal component analysis in the local differential privacy model. *Theoretical Computer Science*.

[Wang et al. 2021] Wang, P.; Lei, Y.; Ying, Y.; and Zhang, H. 2021. Differentially private sgd with non-smooth loss. *arXiv preprint arXiv:2101.08925*.

[Wang, Chen, and Xu 2019] Wang, D.; Chen, C.; and Xu, J. 2019. Differentially private empirical risk minimization with non-convex loss functions. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, 6526–6535.

[Wang, Ye, and Xu 2017] Wang, D.; Ye, M.; and Xu, J. 2017. Differentially private empirical risk minimization revisited: Faster and more general. In *Advances in Neural Information Processing Systems*, 2722–2731.

[Wu et al. 2017] Wu, X.; Li, F.; Kumar, A.; Chaudhuri, K.; Jha, S.; and Naughton, J. 2017. Bolt-on differential privacy for scalable stochastic gradient descent-based analytics. In *Proceedings of the 2017 ACM International Conference on Management of Data*, 1307–1322.

[Xu et al. 2019] Xu, C.; Ren, J.; Zhang, D.; Zhang, Y.; Qin, Z.; and Ren, K. 2019. Ganobfuscator: Mitigating information leakage under gan via differential privacy. *IEEE Transactions on Information Forensics and Security* 2358–2371.

[Yang et al. 2021] Yang, Z.; Lei, Y.; Lyu, S.; and Ying, Y. 2021. Stability and differential privacy of stochastic gradient descent for pairwise learning with non-smooth loss. In *International Conference on Artificial Intelligence and Statistics*, 2026–2034.

[Ying and Zhou 2017] Ying, Y., and Zhou, D.-X. 2017. Unregularized online learning algorithms with general loss functions. *Applied and Computational Harmonic Analysis* 224–244.

[Zhang et al. 2017] Zhang, J.; Zheng, K.; Mou, W.; and Wang, L. 2017. Efficient private erm for smooth objectives. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, 3922–3928.

[Zhou et al. 2019] Zhou, Y.; Yang, J.; Zhang, H.; Liang, Y.; and Tarokh, V. 2019. SGD converges to global minimum in deep learning via star-convex path. In *7th International Conference on Learning Representations, 2019*.

# A. Details of Proofs

## A.1. Proof of Lemma 3

**Lemma 4.** *In gradient perturbation method (1), if Assumptions 1, 4 hold, then*

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq 2\sqrt{2}\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G}{\mu n},$$

*where $\hat{\theta}_n$ and $\hat{\theta}'_n$ denote the models derived from adjacent datasets $S$ and $S'$.*

*Proof.* By triangle inequality, we have:

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq \left\|\hat{\theta}_n - \theta_n^*(S)\right\|_2 + \|\theta_n^*(S) - \theta_n^*(S')\|_2 + \left\|\theta_n^*(S') - \hat{\theta}'_n\right\|_2. \tag{3}$$

Recalling that PL condition implies QG condition:

$$R_n(\theta) - R_n(\theta_n^*) \geq \frac{\mu}{2}\|\theta - \theta_n^*\|_2^2. \tag{4}$$

So we have:

$$\left\|\hat{\theta}_n - \theta_n^*(S)\right\|_2 \leq \sqrt{\frac{2}{\mu}\left(R_n(\hat{\theta}_n) - R_n(\theta_n^*)\right)},$$

$$\left\|\theta_n^*(S') - \hat{\theta}'_n\right\|_2 \leq \sqrt{\frac{2}{\mu}\left(R_n(\hat{\theta}'_n) - R_n(\theta_n^*(S'))\right)}.$$

Due to the symmetry, the two inequalities can be integrated into one, i.e.

$$\left\|\hat{\theta}_n - \theta_n^*(S)\right\|_2 + \left\|\theta_n^*(S') - \hat{\theta}'_n\right\|_2 \leq 2\sqrt{2}\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}}. \tag{5}$$

Now we turn to term $\|\theta_n^*(S) - \theta_n^*(S')\|_2$.

For $R_n(\theta_n^*(S')) - R_n(\theta_n^*(S))$, if assuming the different data instance is $z_j$, we have:

$$
\begin{aligned}
R_n(\theta_n^*(S')) - R_n(\theta_n^*(S)) &= \frac{1}{n}\left(\ell\left(z_j, \theta_n^*(S')\right) - \ell\left(z_j, \theta_n^*(S)\right)\right) + \frac{1}{n}\sum_{i \neq j}\ell\left(z_i, \theta_n^*(S')\right) - \ell\left(z_i, \theta_n^*(S)\right)\\
&= \frac{1}{n}\left(\ell\left(z_j, \theta_n^*(S')\right) - \ell\left(z_j, \theta_n^*(S)\right)\right) + \frac{1}{n}\left(\ell\left(z_j', \theta_n^*(S')\right) - \ell\left(z_j', \theta_n^*(S)\right)\right)\\
&\quad + R_{n'}(\theta_n^*(S')) - R_{n'}(\theta_n^*(S))\\
&\leq \frac{2G}{n}\|\theta_n^*(S') - \theta_n^*(S)\|_2 + R_{n'}(\theta_n^*(S')) - R_{n'}(\theta_n^*(S))\\
&\leq \frac{2G}{n}\|\theta_n^*(S') - \theta_n^*(S)\|_2,
\end{aligned}
\tag{6}
$$

where $R_{n'}(\theta)$ is the empirical risk over dataset $S'$, the first inequality holds because $G$-Lipschitzness and the last inequality holds because $R_{n'}(\theta_n^*(S')) - R_{n'}(\theta_n^*(S)) \leq 0$ due to the definition of $\theta_n^*(S')$.

By inequality (4),

$$R_n(\theta_n^*(S')) - R_n(\theta_n^*(S)) \geq \frac{\mu}{2}\|\theta_n^*(S') - \theta_n^*(S)\|_2^2. \tag{7}$$

Combining inequalities (6) and (7), we have:

$$\frac{\mu}{2}\|\theta_n^*(S') - \theta_n^*(S)\|_2^2 \leq \frac{2G}{n}\|\theta_n^*(S') - \theta_n^*(S)\|_2,$$

which implies

$$\|\theta_n^*(S') - \theta_n^*(S)\|_2 \leq \frac{4G}{\mu n}. \tag{8}$$

Plugging inequalities (5) and (8) back into (3), we have:

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq 2\sqrt{2}\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G}{\mu n},$$

which completes the proof.

$\square$

## A.2. The Optimization Error

As discussed before, the excess population risk can be decomposed into:

$$R(\hat{\theta}_n) - R(\theta^*) = R(\hat{\theta}_n) - R_n(\hat{\theta}_n) + R_n(\hat{\theta}_n) - R_n(\theta^*) + R_n(\theta^*) - R(\theta^*)$$
$$\leq R(\hat{\theta}_n) - R_n(\hat{\theta}_n) + R_n(\hat{\theta}_n) - R_n(\theta_n^*) + R_n(\theta^*) - R(\theta^*). \tag{9}$$

In Lemma 3, the stability is also related to the optimization error (excess empirical risk), so we first discuss the optimization error of the private model $\hat{\theta}_n$, i.e. $R_n(\hat{\theta}_n) - R_n(\theta_n^*)$, under different assumptions.

To get the optimization error bound, we need the following lemma given in (Yang et al. 2021).

**Lemma 5** ((Yang et al. 2021)). *If Gaussian random noise $b \sim \mathcal{N}(0, \sigma^2 I_p)$, then for $\zeta \in (\exp(-p/8), 1)$, we have with probability $1 - \zeta$,*

$$\|b\|_2 \leq \sigma\sqrt{p}\left(1 + \left(\frac{8\log(1/\zeta)}{p}\right)^{1/4}\right).$$

**Lemma 6.** *If the Assumptions 1, 2, 4 hold and the DP model is trained by $T$-iterations gradient perturbation method (1), then taking $T = \mathcal{O}(\log(n))$, $\eta_1 = \cdots = \eta_T = \frac{1}{L}$, if $\zeta \in (\exp(-p/8), 1)$, with probability at least $1 - \zeta$,*

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \lesssim \frac{G^2 p \log(n) \log(1/\delta)}{n^2 \epsilon^2}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2.$$

*Proof.* Note that we assume the loss function is $L$-smooth (Assumption 2, denoted by $L$) and satisfies the PL condition (Assumption 4, denoted by $PL$), at iteration $t$, taking $\eta_t = \frac{1}{L}$, we have:

$$
\begin{aligned}
R_n(\hat{\theta}_{t+1}) - R_n(\hat{\theta}_t) &\overset{(L)}{\leq} \langle \nabla_\theta R_n(\hat{\theta}_t), \hat{\theta}_{t+1} - \hat{\theta}_t \rangle + \frac{L}{2}\left\|\hat{\theta}_{t+1} - \hat{\theta}_t\right\|_2^2 \\
&= -\eta_t \langle \nabla_\theta R_n(\hat{\theta}_t), \nabla_\theta R_n(\hat{\theta}_t) + b_{t+1} \rangle + \frac{L\eta_t^2}{2}\left\|\nabla_\theta R_n(\hat{\theta}_t) + b_{t+1}\right\|_2^2 \\
&= -\frac{1}{L}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 + \frac{1}{2L}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 + \frac{1}{2L}\|b\|_2^2 \\
&= -\frac{1}{2L}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 + \frac{1}{2L}\|b_{t+1}\|_2^2 \\
&\overset{(PL)}{\leq} -\frac{\mu}{L}\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{1}{2L}\|b_{t+1}\|_2^2.
\end{aligned}
\tag{10}
$$

Adding $R_n(\hat{\theta}_t) - R_n(\theta_n^*)$ to both sides, we have:

$$R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*) \leq \left(1 - \frac{\mu}{L}\right)\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{1}{2L}\|b_{t+1}\|_2^2.$$

Summing over $T$ iterations, we have:

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \leq \left(1 - \frac{\mu}{L}\right)^T\left(R_n(\hat{\theta}_0) - R_n(\theta_n^*)\right) + \frac{1}{2L}\sum_{t=0}^{T-1}\left(1 - \frac{\mu}{L}\right)^t\|b_{t+1}\|_2^2. \tag{11}$$

With Lemma 5, with probability at least $1 - \xi$, we have:

$$\|b_t\|_2^2 \leq \sigma^2 p\left(1 + \left(\frac{8\log(T/\xi)}{p}\right)^{1/4}\right)^2,$$

for all $t = 1, \cdots, T$.

Then, with the high probability upper bound of $b_t$, inequality (11) comes to:

$$
\begin{aligned}
R_n(\hat{\theta}_n) - R_n(\theta_n^*) &\leq \left(1 - \frac{\mu}{L}\right)^T\left(R_n(\hat{\theta}_0) - R_n(\theta_n^*)\right) + \frac{1}{2L}\sum_{t=0}^{T-1}\left(1 - \frac{\mu}{L}\right)^t \sigma^2 p\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2 \\
&\leq \left(1 - \frac{\mu}{L}\right)^T M_\ell + \frac{\left(1 - \left(1 - \frac{\mu}{L}\right)^{T-1}\right)\sigma^2 p}{2\mu}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2 \\
&\leq \left(1 - \frac{\mu}{L}\right)^T M_\ell + \frac{\sigma^2 p}{2\mu}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2,
\end{aligned}
\tag{12}
$$

where the second inequality holds because $0 \le \ell(\cdot, \cdot) \le M_\ell$, and $0 < \frac{\mu}{L} < 1$ because of the definitions of $\mu$ and $L$ (Wang, Ye, and Xu 2017).

Taking $\sigma = c\frac{G\sqrt{T \log(1/\delta)}}{n\epsilon}$ given in Lemma 2, and taking $T = \mathcal{O}(\log(n))$, then if $\zeta \in (\exp(-p/8), 1)$, with probability at least $1 - \zeta$, we have:

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \lesssim \frac{G^2 p \log(n) \log(1/\delta)}{2\mu n^2 \epsilon^2} \left( 1 + \left( \frac{8 \log(T/\zeta)}{p} \right)^{1/4} \right)^2 .$$

The result follows.

$\square$

**Lemma 7.** *If the loss function is $\alpha$-Hölder smooth with parameter $H$, satisfies the PL inequality with parameter $2\mu$ and the DP model is trained by $T$-iterations gradient perturbation method (1), then taking $T = \mathcal{O}\left( n^{\frac{2}{1+2\alpha}} \right)$, $\eta_t = \frac{2}{\mu(t+\kappa)}$, where $\kappa \ge \frac{2H^{1/\alpha}}{\mu}$, if $\zeta \in (\exp(-p/8), 1)$, with probability at least $1 - \zeta$,*

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \lesssim \frac{G'^2 \sqrt{p \log(1/\delta)}}{n^{\frac{2\alpha}{1+2\alpha}} \epsilon} \left( 1 + \left( \frac{8 \log(T/\zeta)}{p} \right)^{1/4} \right).$$

*Proof.* The proof is motivated by (Li and Liu 2021).

Like the proof of Lemma 6, by assuming that the loss function is $\alpha$-Hölder smooth (Assumption 3, denoted by $\alpha$), via Lemma 1, at iteration $t$,

$$R_n(\hat{\theta}_{t+1}) - R_n(\hat{\theta}_t) \overset{(\alpha)}{\le} \langle \nabla_\theta R_n(\hat{\theta}_t), \hat{\theta}_{t+1} - \hat{\theta}_t \rangle + \frac{H}{2} \left\| \hat{\theta}_{t+1} - \hat{\theta}_t \right\|_2^{\alpha+1}$$

$$\le \langle \nabla_\theta R_n(\hat{\theta}_t), \hat{\theta}_{t+1} - \hat{\theta}_t \rangle + \frac{H}{\alpha+1} \left\| \hat{\theta}_{t+1} - \hat{\theta}_t \right\|_2^{\alpha+1}$$

$$= -\eta_t \langle \nabla_\theta R_n(\hat{\theta}_t), \nabla_\theta R_n(\hat{\theta}_t) + b_{t+1} \rangle + \frac{H\eta_t^{\alpha+1}}{\alpha+1} \left\| \nabla_\theta R_n(\hat{\theta}_t) + b_{t+1} \right\|_2^{\alpha+1}$$

$$\le -\eta_t \left( \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 + \langle \nabla_\theta R_n(\hat{\theta}_t), b_{t+1} \rangle \right)$$

$$+ \frac{H\eta_t^{\alpha+1}}{\alpha+1} \left( \frac{1-\alpha}{2} + \frac{\alpha+1}{2} \left( \left\| \nabla_\theta R_n(\hat{\theta}_t) + b_{t+1} \right\|_2^{\alpha+1} \right)^{\frac{2}{\alpha+1}} \right)$$

$$\le -\eta_t \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 + \left( H\eta_t^{\alpha+1} - \eta_t \right) \langle \nabla_\theta R_n(\hat{\theta}_t), b_{t+1} \rangle$$

$$+ \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)} + \frac{H\eta_t^{\alpha+1}}{2} \left( \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 + \|b_{t+1}\|_2^2 \right),$$

where the third inequality is because of Young's inequality: if $p^{-1} + q^{-1} = 1$ and $p > 0$, then $uv \le p^{-1}|u|^p + q^{-1}|v|^q$. Here we set $p^{-1} = (1-\alpha)/2$, $q^{-1} = (\alpha+1)/2$. And the last inequality holds because of Cauthy-Schwarz inequality.

Noting that $\eta_t = \frac{2}{\mu(t+\kappa)}$ and $\kappa \ge \frac{2H^{1/\alpha}}{\mu}$, so we have: $\eta_t \le \left( \frac{1}{H} \right)^{1/\alpha}$.

As a result, we have:

$$H\eta_t^{\alpha+1} \le H \left[ \left( \frac{1}{H} \right)^{1/\alpha} \right]^\alpha \eta_t \le \eta_t. \tag{13}$$

As a result,

$$R_n(\hat{\theta}_{t+1}) - R_n(\hat{\theta}_t) \le -\eta_t \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 + \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)} + \frac{H\eta_t^{\alpha+1}}{2} \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2$$

$$+ \eta_t \|\nabla_\theta R_n(\hat{\theta}_t)\|_2 \|b_{t+1}\|_2 + \frac{H\eta_t^{\alpha+1}}{2} \|b_{t+1}\|_2^2$$

$$\le -\frac{\eta_t}{2} \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 + \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)} + 2\eta_t \|\nabla_\theta R_n(\hat{\theta}_t)\|_2 \|b_{t+1}\|_2 + \frac{\eta_t}{2} \|b_{t+1}\|_2^2 \tag{14}$$

$$\overset{(PL)}{\le} -\frac{\eta_t}{4} \left\| \nabla_\theta R_n(\hat{\theta}_t) \right\|_2^2 - \mu\eta_t \left( R_n(\hat{\theta}_t) - R_n(\theta_n^*) \right) + \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)}$$

$$+ \eta_t G' \|b_{t+1}\|_2 + \frac{\eta_t}{2} \|b_{t+1}\|_2^2,$$

where the second inequality is because of (13) and the last inequality holds because we assume that the loss function satisfies the PL condition with parameter $2\mu$, and $G' = \max\{2HM_{\mathcal{C}}, H\}$, as discussed in Lemma 2.

Adding $\frac{\eta_t}{4}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 - R_n(\theta_n^*)$ to both sides of (14), we have

$$R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*) + \frac{\eta_t}{4}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 \leq (1 - \mu\eta_t)\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)}$$
$$+ \eta_t G'\|b_{t+1}\|_2 + \frac{\eta_t}{2}\|b_{t+1}\|_2^2.$$

Taking $\eta_t = \frac{2}{\mu(t+\kappa)}$,

$$R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*) + \frac{1}{2\mu(t+\kappa)}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 \leq \frac{t+\kappa-2}{t+\kappa}\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{H\eta_t^{\alpha+1}(1-\alpha)}{2(\alpha+1)}$$
$$+ \eta_t G'\|b_{t+1}\|_2 + \frac{\eta_t}{2}\|b_{t+1}\|_2^2.$$

Multiply both side by $(t+\kappa)(t+\kappa-1)$,

$$(t+\kappa)(t+\kappa-1)\left(R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*)\right) + \frac{t+\kappa-1}{2\mu}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2$$

$$\leq (t+\kappa-1)(t+\kappa-2)\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + (t+\kappa)^{-\alpha}(t+\kappa-1)\frac{H(1-\alpha)}{2(\alpha+1)}\left(\frac{2}{\mu}\right)^{1+\alpha} \tag{15}$$

$$+ \frac{2G'(t+\kappa-1)}{\mu}\|b_{t+1}\|_2 + \frac{t+\kappa-1}{\mu}\|b_{t+1}\|_2^2.$$

With Lemma 5, for each $t = 1, \cdots, T$, with probability at least $1 - \xi$, we have:

$$\|b_t\|_2 \leq \sigma\sqrt{p}\left(1 + \left(\frac{8\log(1/\xi)}{p}\right)^{1/4}\right),$$

$$\|b_t\|_2^2 \leq \sigma^2 p\left(1 + \left(\frac{8\log(1/\xi)}{p}\right)^{1/4}\right)^2.$$

So with probability at least $1 - \xi$:

$$(t+\kappa)(t+\kappa-1)\left(R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*)\right) + \frac{t+\kappa-1}{2\mu}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2$$

$$\leq (t+\kappa-1)(t+\kappa-2)\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + (t+\kappa)^{-\alpha}(t+\kappa-1)\frac{H(1-\alpha)}{2(\alpha+1)}\left(\frac{2}{\mu}\right)^{1+\alpha}$$

$$+ \frac{2G'(t+\kappa-1)\sigma\sqrt{p}}{\mu}\left(1 + \left(\frac{8\log(1/\xi)}{p}\right)^{1/4}\right) + \frac{(t+\kappa-1)\sigma^2 p}{\mu}\left(1 + \left(\frac{8\log(1/\xi)}{p}\right)^{1/4}\right)^2.$$

By summing over $T$ iterations and taking $\xi = \zeta/T$, with probability at least $1 - \zeta$, we have:

$$(T+\kappa)(T+\kappa-1)\left(R_n(\hat{\theta}_{T+1}) - R_n(\theta_n^*)\right) + \sum_{t=1}^{T}\frac{t+\kappa-1}{2\mu}\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2$$

$$\leq \kappa(\kappa-1)\left(R_n(\hat{\theta}_1) - R_n(\theta_n^*)\right) + \sum_{t=1}^{T}(t+\kappa)^{-\alpha}(t+\kappa-1)\frac{H(1-\alpha)}{2(\alpha+1)}\left(\frac{2}{\mu}\right)^{1+\alpha}$$

$$+ \sum_{t=1}^{T}\frac{2G'(t+\kappa-1)\sigma\sqrt{p}}{\mu}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right) \tag{16}$$

$$+ \sum_{t=1}^{T}\frac{(t+\kappa-1)\sigma^2 p}{\mu}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2.$$

Here, for simplicity, we represent $t = 0, \cdots, T-1$ by $t = 1, \cdots, T$.

Then we bound term $\sum_{t=1}^{T} (t+\kappa)^{-\alpha} (t+\kappa-1) \frac{H(1-\alpha)}{2(\alpha+1)} \left(\frac{2}{\mu}\right)^{1+\alpha}$.

Note that:

$$\sum_{t=1}^{T} (t+\kappa)^{-\alpha} (t+\kappa-1) \leq \sum_{t=1}^{T} (t+\kappa)^{1-\alpha} \leq \int_{1}^{T} (t+\kappa)^{1-\alpha} dt \leq \frac{(T+\kappa)^{2-\alpha}}{2-\alpha}.$$

Plugging the result above back into (16), and note that $0 \leq \ell(\cdot,\cdot) \leq M_\ell$, we have:

$$(T+\kappa)(T+\kappa-1)\left(R_n(\hat{\theta}_{T+1}) - R_n(\theta_n^*)\right) \leq \kappa(\kappa-1)M_\ell + \frac{(T+\kappa)^{2-\alpha}}{2-\alpha} \frac{H(1-\alpha)}{2(\alpha+1)} \left(\frac{2}{\mu}\right)^{1+\alpha}$$

$$+ \left(T\kappa + \frac{T(T-1)}{2}\right) \frac{2G'\sigma\sqrt{p}}{\mu} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)$$

$$+ \left(T\kappa + \frac{T(T-1)}{2}\right) \frac{\sigma^2 p}{\mu} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2.$$

As a result, taking $\sigma$ given in Lemma 2, with probability at least $1-\zeta$, we have:

$$R_n(\hat{\theta}_{T+1}) - R_n(\theta_n^*) \lesssim T^{-\alpha} + \frac{G'^2 \sqrt{Tp\log(1/\delta)}}{n\epsilon} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right).$$

Taking $T = \mathcal{O}\left(n^{\frac{2}{1+2\alpha}}\right)$, with probability at least $1-\zeta$, we have:

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \lesssim \frac{G'^2 \sqrt{p\log(1/\delta)}}{n^{\frac{2\alpha}{1+2\alpha}}\epsilon} \left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right).$$

The result follows.

$\square$

## A.3. Proof of Theorem 1

Before the detailed proof, we first prove the following lemma 10. To get Lemma 10, we need the following lemmas given in (Bousquet, Klochkov, and Zhivotovskiy 2020).

**Lemma 8** ((Bousquet, Klochkov, and Zhivotovskiy 2020)). *Assume that $z_1, \cdots, z_n$ are independent variables and the function $g_i : \mathcal{Z}^n \to \mathbb{R}$ satisfy the following properties for $i = 1, \cdots, n$,*

- $\mathbb{E}_{z_i} g_i(z_1, \cdots, z_n) = 0$ *almost surely;*
- $|\mathbb{E}[g_i(z_1, \cdots, z_n)|z_i]| \leq K$ *almost surely;*
- $|g_i(z_1, \cdots, z_n) - g_i(z_1, \cdots, z_{j-1}, z_j', z_{j+1}, \cdots, z_n)| \leq \beta.$

*Then the following inequality holds for all $q \geq 2$,*

$$\left\|\sum_{i=1}^{n} g_i\right\|_q \leq 12\sqrt{2}\beta qn \log(n) + 4K\sqrt{qn}.$$

**Lemma 9** ((Bousquet, Klochkov, and Zhivotovskiy 2020)). *Under the uniform stability condition with parameter $\gamma$ and uniformly bounded loss function $\ell(\cdot,\cdot) \leq M_\ell$, we have for $g_i = \mathbb{E}_{z_i'}\left(\ell(z_i, \theta_n^{(i)}) - \mathbb{E}_z \ell(z, \theta_n^{(i)})\right)$,*

$$\left|n\left(R_n(\theta_n) - R(\theta_n)\right) - \sum_{i=1}^{n} g_i\right| \leq 2\gamma n.$$

**Lemma 10.** *Defining the DP algorithm (model) training by $T$-iterations gradient perturbation method (like (1)) $\hat{\theta}_n = \hat{\theta}(z_1, \cdots, z_n)$ and its independent copy $\hat{\theta}_n' = \hat{\theta}(z_1', \cdots, z_n')$. Then for all $q \geq 2$,*

$$\left\|R_n(\hat{\theta}_n) - R(\hat{\theta}_n) - \frac{1}{n}\sum_{i=1}^{n} \mathbb{E}\left[\ell(z_i, \hat{\theta}_n')|z_i\right] + \mathbb{E}R(\hat{\theta}_n)\right\|_q \lesssim Gq\log(n)\left(\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{G}{\mu n}\right).$$

*Proof.* Via Lemma 4,

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq 2\sqrt{2}\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G}{\mu n}.$$

Recalling the definition of $\gamma$-uniformly stability: If for any $z, z', z_1, \cdots, z_n \in \mathcal{Z}$ and $i = 1, \cdots, n$, it holds that

$$|\ell(z, \theta_n(z_1, \cdots, z_n)) - \ell(z, \theta_n(z_1, \cdots, z_{i-1}, z', z_{i+1}, \cdots, z_n))| \leq \gamma.$$

Then with $G$-Lipschitzness, we have:

$$\left|\ell(z, \hat{\theta}_n) - \ell(z, \hat{\theta}'_n)\right| \leq 2\sqrt{2}G\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G^2}{\mu n},$$

where $\hat{\theta}_n$ and $\hat{\theta}'_n$ are private models derived from any adjacent datasets. In the following, we use $\hat{\theta}_n^{(i)}$ to represent $\hat{\theta}'_n$, which means that the single different data instance is the $i^{th}$ one.

Considering the function $g_i(z_1, \cdots, z_n) = \mathbb{E}_{z'_i}[\ell(z_i, \hat{\theta}_n^{(i)})] - \mathbb{E}_{z'_i}[R(\hat{\theta}_n^{(i)})]$, via the definition of $R(\hat{\theta}^{(i)})$, we have: $\mathbb{E}_{z_i}g_i(z_1, \cdots, z_n) = 0$.

With the stability of the DP model, we have:

$$\left|g_i(z_1, \cdots, z_n) - g_i(z_1, \cdots, z_{j-1}, z'_j, z_{j+1}, \cdots, z_n)\right| \leq \beta := 2\left(2\sqrt{2}G\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G^2}{\mu n}\right).$$

If considering $h_i(z_1, \cdots, z_n) = g_i(z_1, \cdots, z_n) - \mathbb{E}[g_i(z_1, \cdots, z_n)|z_i]$, we have:

$$\mathbb{E}_{z_i} h_i(z_1, \cdots, z_n) = 0$$

almost surely, and

$$\left|h_i(z_1, \cdots, z_n) - h_i(z_1, \cdots, z_{j-1}, z'_j, z_{j+1}, \cdots, z_n)\right| \leq 2\beta = 4\left(2\sqrt{2}G\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G^2}{\mu n}\right).$$

Via the definition of $h_i$, we observe that $\mathbb{E}[h_i|z_i] = 0$ almost surely, which further implies $K = 0$ in Lemma 8, so we have for $q \geq 2$:

$$\left\|\sum_{i=1}^{n} h_i\right\|_q = \left\|\sum_{i=1}^{n}(g_i - \mathbb{E}[g_i|z_i])\right\|_q \leq 192Gqn\log(n)\left(\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{2\mu}} + \frac{G}{\mu n}\right).$$

Via Lemma 9, we have:

$$\left|n\left(R_n(\hat{\theta}_n) - R(\hat{\theta}_n)\right) - \sum_{i=1}^{n} g_i\right| \leq 2n\left(2\sqrt{2}G\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{4G^2}{\mu n}\right).$$

Noting that

$$\mathbb{E}[g_i|z_i] = \mathbb{E}\left[\ell(z_i, \hat{\theta}'_n)|z_i\right] - \mathbb{E}R(\hat{\theta}'_n),$$

we have:

$$\left\|R_n(\hat{\theta}_n) - R(\hat{\theta}_n) - \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left[\ell(z_i, \hat{\theta}'_n)|z_i\right] + \mathbb{E}R(\hat{\theta}_n)\right\|_q \lesssim Gq\log(n)\left(\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{G}{\mu n}\right).$$

The result follows.

$\square$

To get Theorem 1, we further need the following lemma given in (Boucheron, Lugosi, and Massart 2013).

**Lemma 11** ((Boucheron, Lugosi, and Massart 2013))**.** *If $X_1, \cdots, X_n$ are zero mean, independent and bounded $|X_i| \leq M$ almost surely, then for $q \geq 2$,*

$$\|X_1 + \cdots X_n\|_q \leq 6\sqrt{\left(\sum_{i=1}^{n}\mathbb{E}[X_i^2]\right)q} + 4qM.$$

Then, we can start our proof.

**Theorem 4.** *If Assumptions 1, 2, and 4 hold, the loss function is bounded, i.e.* $0 \leq \ell(\cdot, \cdot) \leq M_\ell$, *taking $\sigma$ given by Lemma 2,* $T = \mathcal{O}(\log(n))$, $\eta_1 = \cdots = \eta_T = \frac{1}{L}$, *if $\zeta \in (\exp(-p/8), 1)$, then with probability at least $1 - \zeta$:*

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \frac{G^2 p \log(n) \log(1/\delta)}{n^2 \epsilon^2} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)^2$$

$$+ c_2 \frac{G \log^{1.5}(n) \sqrt{p \log(1/\delta)}}{n\epsilon} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right) + c_3 \frac{\log(n)}{n}.$$

*for some constants $c_1, c_2, c_3 > 0$.*

*Proof.* Via Lemma 10, we have:

$$R_n(\hat{\theta}_n) - R(\hat{\theta}_n) = \rho + \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \mathbb{E} R(\hat{\theta}_n),$$

where $\|\rho\|_q \lesssim Gq \log(n) \left(\sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{G}{\mu n}\right)$ for $q \geq 2$ and $\mathbb{E}'$ denotes the expectation taken over the independent copy.

Plugging this back to (9), we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq \left(R_n(\hat{\theta}_n) - R_n(\theta_n^*)\right) + (R_n(\theta^*) - R(\theta^*)) - \rho - \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}' \ell(z_i, \hat{\theta}'_n) + \mathbb{E} R(\hat{\theta}_n).$$

Noting that $R_n(\theta^*) = \frac{1}{n} \sum_{i=1}^{n} \ell(z_i, \theta^*)$, we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq \left(R_n(\hat{\theta}_n) - R_n(\theta_n^*)\right) + \left(\mathbb{E} R(\hat{\theta}_n) - R(\theta^*)\right) - \rho - \frac{1}{n} \sum_{i=1}^{n} \left(\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \ell(z_i, \theta^*)\right). \tag{17}$$

Based on the definition of $R(\theta)$, Assumption 5 is equivalent to:

$$\mathbb{E}\left[(\ell(z, \theta) - \ell(z, \theta^*))^2\right] \leq B \left(\mathbb{E}\ell(z, \theta) - \mathbb{E}\ell(z, \theta^*)\right). \tag{18}$$

With $G$-Lipschitz and PL inequality with parameter $\mu$, we have $B = 2G^2/\mu$.

So, via (18),

$$\mathbb{E}\left[\left(\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \ell(z_i, \theta^*)\right)^2\right] \leq \frac{2G^2}{\mu} \left(\mathbb{E}\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \mathbb{E}\ell(z_i, \theta^*)\right)$$

$$= \frac{2G^2}{\mu} \left(\mathbb{E}[R(\hat{\theta}'_n)] - R(\theta^*)\right), \tag{19}$$

where the last equation holds because $\mathbb{E}\mathbb{E}' \ell(z_i, \hat{\theta}'_n) = \mathbb{E}[R(\hat{\theta}'_n)]$.

Note that term $\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \ell(z_i, \theta^*)$ can be decomposed as the following:

$$\underbrace{\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \ell(z_i, \theta^*)}_{X_i} = \underbrace{\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \mathbb{E}' \ell(z_i, \theta'_n)}_{X'_i} + \underbrace{\mathbb{E}' \ell(z_i, \theta'_n) - \ell(z_i, \theta^*)}_{X''_i}.$$

Via triangle inequality,

$$\|X_i\|_q \leq \|X'_i\|_q + \|X''_i\|_q.$$

Recalling the definition of $R_n(\hat{\theta}_n) - R_n(\theta_n^*)$, we have:

$$\left\|\frac{1}{n} \sum_{i=1}^{n} X'_i\right\|_q = R_n(\hat{\theta}_n) - R_n(\theta_n^*). \tag{20}$$

Via Lemma 11, since $\mathbb{E}[R(\theta'_n)] - R(\theta^*)$ is exactly the expectation of each $X''_i$, we have for $q \geq 2$,

$$\left\|\frac{1}{n} \sum_{i=1}^{n} \mathbb{E}'\left[\ell(z_i, \theta'_n)\right] - \ell(z_i, \theta^*) - \mathbb{E}[R(\theta'_n)] + R(\theta^*)\right\|_q \lesssim \sqrt{\mathbb{E}\left[\left(\mathbb{E}' \ell(z_i, \hat{\theta}'_n) - \ell(z_i, \theta^*)\right)^2\right] \frac{q}{n}}$$

$$\leq \sqrt{\frac{2G^2}{\mu} \left(\mathbb{E}[R(\theta_n)] - R(\theta^*)\right) \frac{q}{n}} + \frac{qM_\ell}{n}, \tag{21}$$

where the last inequality holds because of (19) and $\mathbb{E}[R(\theta_n)] = \mathbb{E}[R(\hat{\theta}'_n)]$.

Plugging (20) and (21) back into (17), we obtain for each $q \geq 2$ and some constant $C > 0$,

$$\left\| R(\hat{\theta}_n) - R(\theta^*) - \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) \right\|_q$$

$$\leq C \left( Gq \log(n) \left( \sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \frac{G}{\mu n} \right) + \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) + \sqrt{\frac{2G^2}{\mu} \left( \mathbb{E}[R(\theta_n)] - R(\theta^*) \right) \frac{q}{n}} + \frac{qM_\ell}{n} \right)$$

$$\leq \varphi C \left( \mathbb{E}[R(\theta_n)] - R(\theta^*) \right) + C \left( Gq \log(n) \sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) \right)$$

$$+ C \left( \left( \frac{G^2 \log(n)}{\mu} + \frac{2G^2}{\mu\varphi} + M_\ell \right) \frac{q}{n} \right), \tag{22}$$

where the last inequality holds because for $a, b, \varphi > 0$, $\sqrt{ab} \leq \varphi a + b/\varphi$.

Taking $q = 2$, and via Cauchy-Schwarz inequality,

$$\mathbb{E}[R(\hat{\theta}_n)] - R(\theta^*) - \mathbb{E}\left[ R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right]$$

$$\leq \left\| R(\hat{\theta}_n) - R^* - \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) \right\|_2$$

$$\leq \varphi C \left( \mathbb{E}[R(\hat{\theta}_n)] - R(\theta^*) \right) + C \left( 2G \log(n) \sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}} + \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) \right)$$

$$+ \frac{2C}{n} \left( \frac{G^2 \log(n)}{\mu} + \frac{2G^2}{\mu\varphi} + M_\ell \right).$$

The inequality above can be rewritten as:

$$\mathbb{E}[R(\hat{\theta}_n)] - R(\theta^*) \leq \frac{1}{1 - \varphi C} \mathbb{E}[R_n(\hat{\theta}_n) - R_n(\theta_n^*)] + \frac{C}{1 - \varphi C} \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) + \frac{2CG \log(n)}{1 - \varphi C} \sqrt{\frac{R_n(\hat{\theta}_n) - R_n(\theta_n^*)}{\mu}}$$

$$+ \frac{2C}{(1 - \varphi C)n} \left( \frac{G^2 \log(n)}{\mu} + \frac{2G^2}{\mu\varphi} + M_\ell \right).$$

Taking this back to (22), we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \left( R_n(\hat{\theta}_n) - R_n(\theta_n^*) \right) + c_2 \log(n) \sqrt{R_n(\hat{\theta}_n) - R_n(\theta_n^*)} + c_3 \frac{\log(n)}{n}, \tag{23}$$

for some constants $c_1, c_2$ and $c_3$, where we combine $R_n(\hat{\theta}_n) - R_n(\theta_n^*)$ and its expectation together.

Then via Lemma 6, with probability at least $1 - \zeta$, we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \frac{G^2 p \log(n) \log(1/\delta)}{n^2 \epsilon^2} \left( 1 + \left( \frac{8 \log(T/\zeta)}{p} \right)^{1/4} \right)^2$$

$$+ c_2 \frac{G \log^{1.5}(n) \sqrt{p \log(1/\delta)}}{n\epsilon} \left( 1 + \left( \frac{8 \log(T/\zeta)}{p} \right)^{1/4} \right) + c_3 \frac{\log(n)}{n}.$$

for some constants $c_1, c_2, c_3 > 0$.

The result follows. $\square$

## A.4. Proof of Theorem 2

**Theorem 5.** *If the loss function is $\alpha$-Hölder smooth (Assumption 3) with parameter $H$, and satisfies the PL condition with parameter $2\mu$ (Assumption 4), the loss function and the parameter space are bounded, i.e. $0 \leq \ell(\cdot, \cdot) \leq M_\ell$, $\|\mathcal{C}\|_2 \leq M_\mathcal{C}$. Taking $\sigma$ given by Lemma 2, $T = \mathcal{O}\left( n^{\frac{2}{1+2\alpha}} \right)$, and $\eta_t = \frac{2}{\mu(t+\kappa)}$, where $\kappa \geq \frac{2H^{1/\alpha}}{\mu}$, if $\zeta \in (\exp(-p/8), 1)$, then with*

*probability at least $1 - \zeta$:*

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \frac{G'^2 \sqrt{p \log(1/\delta)}}{n^{\frac{2\alpha}{1+2\alpha}} \epsilon} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)$$

$$+ c_2 \frac{G' \log(n) \sqrt[4]{p \log(1/\delta)}}{n^{\frac{\alpha}{1+2\alpha}} \epsilon^{1/2}} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)^{1/2} + c_3 \frac{\log(n)}{n},$$

*for some constants $c_1, c_2$ and $c_3$, where $G' = \max\{2HM_{\mathcal{C}}, H\}$.*

*Proof.* Like inequality (23) in the proof of Theorem 1 (Appendix A.3), we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \left(R_n(\hat{\theta}_n) - R_n(\theta_n^*)\right) + c_2 \log(n) \sqrt{R_n(\hat{\theta}_n) - R_n(\theta_n^*)} + c_3 \frac{\log(n)}{n}, \tag{24}$$

for some constants $c_1, c_2$ and $c_3$.

The differences between inequalities (23) and (24) are in constants, for example, Lipschitz constant $G$ discussed in Appendix A.3 comes to $G' = \max\{2HM_{\mathcal{C}}, H\}$ here, as discussed before; and in Appendix A.3, the PL condition is with parameter $\mu$, rather than $2\mu$ here.

Combining the result obtained by Lemma 7, taking $T = \mathcal{O}\left(n^{\frac{2}{1+2\alpha}}\right)$ and $\eta_t = \frac{2}{\mu(t+\kappa)}$ with $\kappa \geq \frac{2H^{1/\alpha}}{\mu}$, then with probability at least $1 - \zeta$,

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \frac{G'^2 \sqrt{p \log(1/\delta)}}{n^{\frac{2\alpha}{1+2\alpha}} \epsilon} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)$$

$$+ c_2 \frac{G' \log(n) \sqrt[4]{p \log(1/\delta)}}{n^{\frac{\alpha}{1+2\alpha}} \epsilon^{1/2}} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)^{1/2} + c_3 \frac{\log(n)}{n},$$

for some constants $c_1, c_2$ and $c_3$, where $G' = \max\{2HM_{\mathcal{C}}, H\}$.

The result holds. $\qquad\square$

## A.5. Proof of Theorem 3

**Theorem 6.** *If Assumptions 3, 4 hold, the loss function and the parameter space are bounded, i.e. $0 \leq \ell(\cdot, \cdot) \leq M_\ell$, $\|\mathcal{C}\|_2 \leq M_{\mathcal{C}}$. Taking $\sigma$ given by Lemma 2, $T = \mathcal{O}(\log(n))$, and $\eta_1 = \cdots = \eta_T = \eta$, where $\eta = \left(\frac{1}{H}\right)^{1/\alpha}$, if $\zeta \in (\exp(-p/8), 1)$, then with probability at least $1 - \zeta$,*

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1 \frac{G'^2 \log(n) p \log(1/\delta)}{n^2 \epsilon^2} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right)^2$$

$$+ c_2 \frac{G' \log^{1.5}(n) \sqrt{p \log(1/\delta)}}{n\epsilon} \left(1 + \left(\frac{8 \log(T/\zeta)}{p}\right)^{1/4}\right) + c_3 \frac{\log(n)}{n},$$

*for some constants $c_1, c_2, c_3 > 0$, where $G' = \max\{2HM_{\mathcal{C}}, H\}$.*

*Proof.* The proof is similar to Theorems 1 and 2, we first analyze the optimization error $R_n(\hat{\theta}_n) - R_n(\theta_n^*)$.

For algorithm 1, with normalization, if taking $\eta_t = H^{-1/\alpha}$ we have:

$$R_n(\hat{\theta}_{t+1}) - R_n(\hat{\theta}_t) \overset{(\alpha)}{\leq} \langle \nabla_\theta R_n(\hat{\theta}_t), \hat{\theta}_{t+1} - \hat{\theta}_t \rangle + \frac{H}{2} \left\|\hat{\theta}_{t+1} - \hat{\theta}_t\right\|_2^{\alpha+1}$$

$$= -\eta_t \langle \nabla_\theta R_n(\hat{\theta}_t), \nabla_\theta R_n(\hat{\theta}_t) + b_t \rangle + \frac{H \eta_t^{\alpha+1}}{2} \left(\left\|\nabla_\theta R_n(\hat{\theta}_t) + b_t\right\|_2\right)^{\alpha+1}$$

$$\leq -\eta_t \left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 + \frac{H \eta_t^{\alpha+1}}{2} \left(\left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2 + \|b\|_2\right)^2 + (H \eta_t^{\alpha+1} - \eta_t) \langle \nabla_\theta R_n(\hat{\theta}_t), b_t \rangle$$

$$\leq \frac{-\eta_t}{2} \left\|\nabla_\theta R_n(\hat{\theta}_t)\right\|_2^2 + \frac{H \eta_t^{\alpha+1}}{2} \|b\|_2^2$$

$$\overset{(PL)}{\leq} -\mu \eta_t \left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{\eta_t}{2} \|b\|_2^2,$$

where the second inequality holds because by normalization, and the third inequality holds because $\eta_t = \left(\frac{1}{H}\right)^{1/\alpha}$.

Summing $R_n(\theta_t) - R_n(\theta_n^*)$ to both sides, we have:

$$R_n(\hat{\theta}_{t+1}) - R_n(\theta_n^*) \leq (1 - \mu\eta_t)\left(R_n(\hat{\theta}_t) - R_n(\theta_n^*)\right) + \frac{\eta_t}{2}\|b\|_2^2.$$

Summing over $T$ iterations,

$$R_n(\hat{\theta}_T) - R_n(\theta_n^*) \leq (1 - \mu\eta)^T \left(R_n(\hat{\theta}_0) - R_n(\theta_n^*)\right) + \frac{\eta}{2}\sum_{t=0}^{T-1}(1 - \mu\eta)^t \|b_t\|_2^2,$$

where $\eta = \left(\frac{1}{H}\right)^{1/\alpha}$.

With Lemma 5, with probability at least $1 - \xi$,

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \leq (1 - \mu\eta)^T M_\ell + \frac{\sigma^2 p}{2\mu}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2.$$

Noting that by definition, $\mu\eta \leq 1$, so if taking $T = \mathcal{O}(\log(n))$, with probability at least $1 - \zeta$, we have:

$$R_n(\hat{\theta}_n) - R_n(\theta_n^*) \leq c\frac{G'^2\log(n)p\log(1/\delta)}{n^2\epsilon^2}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2.$$

Then, like in (24), we have:

$$R(\hat{\theta}_n) - R(\theta^*) \leq c_1\left(R_n(\hat{\theta}_n) - R_n(\theta_n^*)\right) + c_2\log(n)\sqrt{R_n(\hat{\theta}_n) - R_n(\theta_n^*)} + c_3\frac{\log(n)}{n}$$

$$\leq c_1\frac{G'^2\log(n)p\log(1/\delta)}{n^2\epsilon^2}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right)^2$$

$$+ c_2\frac{G'\log^{1.5}(n)\sqrt{p\log(1/\delta)}}{n\epsilon}\left(1 + \left(\frac{8\log(T/\zeta)}{p}\right)^{1/4}\right) + c_3\frac{\log(n)}{n},$$

for some constants $c_1, c_2, c_3 > 0$.

The result follows.

$\square$

# B. More Experimental Results

## B.1. Accuracies on More Datasets

In this section, we show the experimental results on datasets Breast Cancer, Credit Card Fraud, and Bank. Details are shown in Figure 2.

The results are similar to which given by Figure 1 in Section 5: although there are some fluctuations over some datasets (such as Bank), the performance of our proposed m-NGP method is similar to or better than traditional method on most datasets.

## B.2. Convergence Rate and Normalization

In this section, we perform experiments to demonstrate the effects on the convergence rate caused by normalization when applying m-NGP. The privacy budget $\epsilon$ is set 0.5. Detailed results are shown in Figure 3.

In Figure 3, the lines with dark color and light color correspond to m-NGP and TGP, respectively, and the shadow area represents the maximum and minimum loss over mutiple experiments, reflecting the variance. And the horizontal axis is iterations and the ordinate is the loss. The experimental results show that over most datasets, m-NGP (normalization) achieves faster convergence rate, comparing with TGP, which is in line with the theoretical analysis.

## B.3. Accuracy and Dimension $p$

In this section, we perform experiments to demonstrate the effects on the accuracy brought by the dimensions of data instances. The experiments are performed on datasets Credit Card Fraud, Bank, and Adult, whose dimensions are 29, 48, and 104, respectively. And the privacy budget $\epsilon$ is set 0.5. The results are shown in Figure 4.

For abscissa, the first dimensions of parts (a), (b), and (c) are set $p = 29, 48, 104$, they are original features given by the datasets. And the dimensions more than them are all set 0, to evaluate the effects brought by the magnitude of $p$, without introducing new information.

Experimental results show that although there may exist some fluctuations caused by the injected random noise, the accuracy decreases with the increasing of $p$ overall, which is in line with the theoretical analysis given in Section 4.
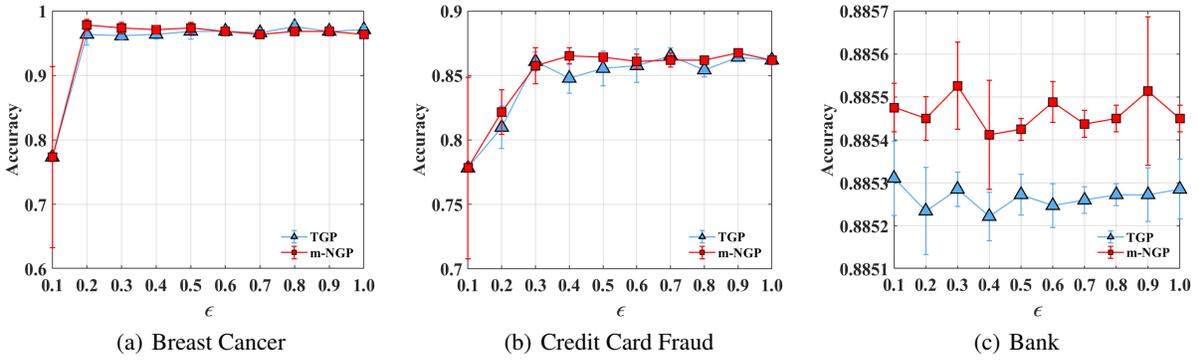
(a) Breast Cancer     (b) Credit Card Fraud     (c) Bank

Figure 2: Comparisons between Traditional Gradient Perturbation (TGP) method and $\max\{\mathbf{1}, \mathbf{g}\}$-Normalized Gradient Perturbation (m-NGP) method.
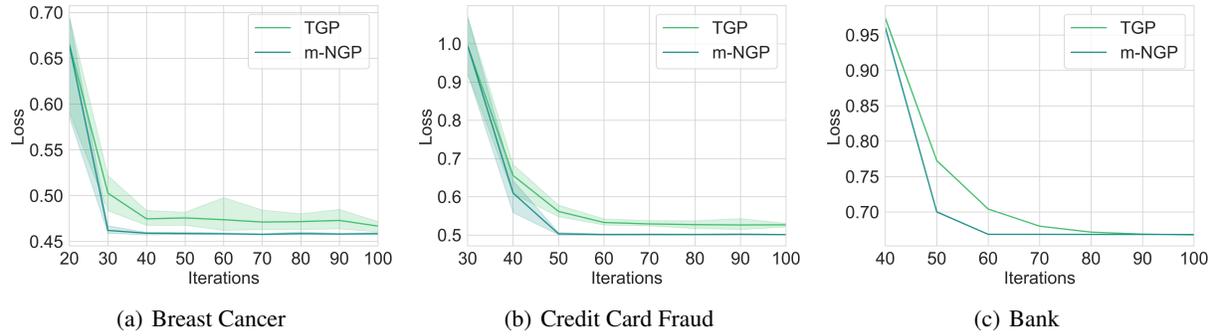


(a) Breast Cancer     (b) Credit Card Fraud     (c) Bank

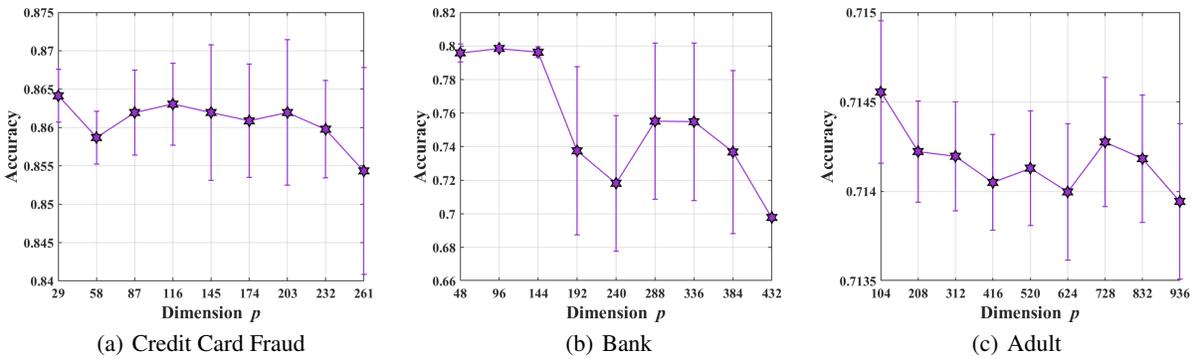Figure 3: Convergence Rates of TGP and m-NGP.



(a) Credit Card Fraud     (b) Bank     (c) Adult

Figure 4: Effects of dimension $p$ on m-NGP.