Stability and Differential Privacy of Stochastic Gradient Descent for Pairwise Learning with Non-Smooth Loss

Zhenhuan Yang¹ Yunwen Lei² Siwei Lyu³ Yiming Ying^{1,*}
University at Albany, SUNY, USA¹ University of Birmingham, UK² University at Buffalo, SUNY, USA³

Abstract

Pairwise learning has recently received increasing attention since it subsumes many important machine learning tasks (e.g. AUC maximization and metric learning) into a unifying framework. In this paper, we give the first-ever-known stability and generalization analysis of stochastic gradient descent (SGD) for pairwise learning with non-smooth loss functions, which are widely used (e.g. Ranking SVM with the hinge loss). We introduce a novel decomposition in its stability analysis to decouple the pairwisely dependent random variables, and derive generalization bounds which are consistent with the setting of pointwise learning. Furthermore, we apply our stability analysis to develop differentially private SGD for pairwise learning, for which our utility bounds match with the state-ofthe-art output perturbation method (Huai et al., 2020) with smooth losses. Finally, we illustrate the results using specific examples of AUC maximization and similarity metric learning. As a byproduct, we provide an affirmative solution to an open question on the advantage of the nuclear-norm constraint over the Frobenius-norm constraint in similarity metric learning.

1 Introduction

Let the input space \mathcal{X} be a compact domain of \mathbb{R}^d , the output space $\mathcal{Y} \subseteq \mathbb{R}$, and the domain of model parameters $\mathcal{W} \subseteq \mathbb{R}^d$. In the standard supervised learning, one aims to learn the relation between the input and output variables from a training dataset $S = \{z_i =$

Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s). * Corresponding author.

 $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, 2, ..., n$ which is i.i.d. from an unknown distribution P on $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. In such cases, the quality of a model parameter \mathbf{w} is often measured by a *pointwise loss* function $\ell(\mathbf{w}, z)$.

In this paper, we are concerned with another important class of learning tasks called pairwise learning where the quality of a model parameter \mathbf{w} is measured by a pairwise loss $\ell(\mathbf{w},z,z')$ on pairs of examples (z,z') as opposed to the pointwise loss $\ell(\mathbf{w},z)$ in standard classification and regression. This pairwise learning framework instantiates many important learning tasks such as similarity and metric learning (Weinberger and Saul, 2009; Xing et al., 2003; Ying and Li, 2012), AUC maximization and bipartite ranking (Agarwal and Niyogi, 2009; Clémençon et al., 2008; Gao et al., 2013; Ying et al., 2016; Zhao et al., 2011), gradient learning (Mukherjee and Wu, 2006; Mukherjee and Zhou, 2006), and minimum error entropy principle (Hu et al., 2013).

Stochastic gradient descent (SGD) has become the workhorse behind many machine learning algorithms for large-scale data analysis. SGD and its variants have been widely studied in the pointwise learning case (Bach and Moulines, 2013; Bottou and Cun, 2004; Lacoste-Julien et al., 2012; Rakhlin et al., 2012; Shalev-Shwartz et al., 2009; Ying and Zhou, 2006) as well as the pairwise learning case (Kar et al., 2013) Lin et al., 2017; Wang et al., 2012; Ying and Zhou, 2016). In particular, Kar et al. (2013); Wang et al. (2012) studied the online-to-batch conversion bounds for online pairwise learning. The work of Shen et al. (2020) studied the stability and generalization of SGD in pairwise learning and derived lower bounds for their optimization error over a class of pairwise losses. This work used the uniform stability (Agarwal et al., 2010) which was largely motivated by Hardt et al. (2016) in the pointwise case. However, there are some fundamental limitations in the work by Shen et al. (2020): it requires the pairwise loss to be both Lipschitz continuous and strongly smooth, and the parameter domain W is assumed to be bounded. Such assumptions are very restrictive which are violated in many cases such as the least square loss for AUC maximization $(1 - \mathbf{w}^{\top}(\mathbf{x} - \mathbf{x}'))^2 \mathbb{I}_{[y=1 \wedge y'=-1]}$ with $\mathbf{w} \in \mathbb{R}^d$ (\mathbb{I} is the indicator function) and the hinge loss for metric learning $(1 + \tau(y, y')(\mathbf{x} - \mathbf{x}')^{\top}\mathbf{w}(\mathbf{x} - \mathbf{x}'))_+$ where \mathbf{w} is a positive semi-definite matrix, and $\tau(y, y') = 1$ if \mathbf{x}, \mathbf{x}' are from the same class and -1 otherwise.

On the other important front, the concept of stability is closely related to differential privacy (DP) (Dwork et al., 2006, 2014) which is a well accepted mathematical definition for privacy protection. While private SGD has been extensively studied (Bassily et al., 2020, 2019; Wu et al., 2017) in pointwise learning, there is litter work on differentially private SGD for pairwise learning except the very recent work of Huai et al. (2020). However, the study (Huai et al., 2020) again requires the loss to be both Lipschitz continuous and strongly smooth.

In this paper, we study the stability, generalization, and differential privacy of SGD for pairwise learning with non-smooth losses. Our contributions can be summarized as follows.

- We establish the first-ever-known stability bounds of SGD for pairwise learning with non-smooth loss functions. Our results hold true for both bounded and unbounded parameter domains. The proof techniques are mainly motivated by the recent work (Bassily et al.) 2020; Lei and Ying, 2020) where stability of SGD was established in the pointwise case. The main challenge here is that pairs of examples involved in pairwise learning are not statistically independent. To overcome this hurdle, we develop a novel approach for decoupling such pairwisely dependent random variables in the analysis. We also derive the first generalization bound in high probability for SGD in pairwise learning using the stability approach.
- We study the differential privacy guarantee and utility bounds of private SGD for pairwise learning by output perturbation method. Our idea is to use our stability results to derive its sensitivity with high probability w.r.t. the randomness of algorithm, and hence guarantee its differential privacy with smaller added noise. The resulting utility bound matches with the output perturbation method in Huai et al. (2020) for private SGD in pairwise learning with smooth losses.
- We provide concrete examples of pairwise learning including AUC maximization and similarity metric learning to illustrate our stability and differential privacy results. In particular, we give an affirmative solution to the open question raised in Cao et al. (2016) that whether similarity metric learning with nuclear-norm constraint can yield milder de-

pendence on the dimensionality than the Frobenius-norm constraint.

Other Related Work. Generalization analysis for the ERM formulation in pairwise learning was studied using U-Statistics (e.g. De la Pena and Giné (2012)) for ranking Clémençon et al. (2008); Rejchel (2012) and metric learning (Cao et al., 2016; Verma and Branson, 2015). There are a considerable amount of work on studying SGD and online learning algorithms in pairwise learning. In particular, generalization bounds for online pairwise learning algorithms were established in Kar et al. (2013); Wang et al. (2012) using online-to-batch conversion techniques (Cesa-Bianchi et al., 2004) which involves the Rademacher complexity or the covering number. The convergence (optimization error) of SGD type algorithms for pairwise learning was obtained in Lin et al. (2017); Ying and Zhou (2016) where the algorithms there directly minimize the population risk. In this setting, there is no need to consider generalization (estimation error) i.e. the difference between the empirical risk and the true population risk.

Algorithmic stability and generalization bounds were established in Agarwal and Niyogi (2009) for ranking problems, and in Jin et al. (2009) for regularized metric learning with a strongly convex objective function, and both studies considered the ERM formulation with a strongly convex objective function. Recently, the uniform stability and its trade-off with optimization errors were studied in Shen et al. (2020) for SGD in pairwise learning, which is inspired by the recent work in pointwise learning (Charles and Papailiopoulos, 2018) Hardt et al., 2016; Kuzborskij and Lampert, 2018). However, the loss there is assumed to be Lipschitz and strongly smooth and the domain $\mathcal W$ needs to be bounded.

The concept of stability was recently used to study the generalization (utility) of differentially private SGD algorithms, particularly in pointwise learning. Specifically, the work of Wu et al. (2017) studied the output perturbation using sensitivity analysis which is very close to the concept of uniform stability. In Bassily et al. (2019), using stability approach, the optimal excess generalization bound $\tilde{\mathcal{O}}(\max\{1/\sqrt{n},\sqrt{d}/(n\epsilon)\})$ was established for (ϵ, δ) -DP algorithms which, however, requires the loss function to be Lipschitz and strongly smooth, and the domain W be bounded. For the non-smooth loss, it proposed to smooth the loss by its Moreau envelope function which is not an ideal solution as the Moreau envelope function is not easy to compute for a general loss. In Feldman et al. (2020), multi-phrased SGD were proposed with the optimal population risk in which, for the non-smooth case, their algorithm is significantly more involved than the noisy SGD algorithm. In regard to the differential private SGD in the pairwise case, the only work that we are aware of is Huai et al. (2020) which studied both gradient perturbation and output perturbation with Gaussian noise. They derive the rate $\tilde{\mathcal{O}}(\sqrt{d}/(\sqrt[4]{n}\epsilon))$ for gradient perturbation and $\tilde{\mathcal{O}}(\sqrt{d}/(\sqrt[4]{n}\epsilon))$ for output perturbation. Note that the loss function there needs to be both Lipschitz continuous and strongly smooth.

2 Main Results

Before stating our main results, we first introduce necessary materials and notations. Given a pairwise loss function $\ell: \mathcal{W} \times \mathcal{Z} \times \mathcal{Z} \to \mathbb{R}$, we aim to minimize the following population risk

$$R(\mathbf{w}) = \mathbb{E}_{\mathbf{z}, \mathbf{z}'}[\ell(\mathbf{w}, \mathbf{z}, \mathbf{z}')],$$

where \mathbf{z} and \mathbf{z}' are drawn independently from the population distribution P on \mathcal{Z} . The population distribution is often unknown and we only have access to a set of i.i.d. training data $S = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\} \in \mathcal{Z}^n$. The task then reduces to minimizing the empirical risk

$$\min_{\mathbf{w} \in \mathcal{W}} R_S(\mathbf{w}) := \frac{1}{n(n-1)} \sum_{i,j=1, i \neq j}^n \ell(\mathbf{w}, \mathbf{z}_i, \mathbf{z}_j).$$
 (1)

Randomized optimization algorithm $\mathcal{A}: \mathcal{Z}^n \to \mathcal{W}$ provides an efficient approach to find an approximate solution to problem (1), which takes S as input and produces an output $\mathcal{A}(S) \in \mathcal{W}$. The randomized algorithm \mathcal{A} here can be either SGD for pairwise learning or its noisy variant for differential privacy. The performance of \mathcal{A} is quantified by the excess population risk: $\epsilon_{\text{risk}}(\mathcal{A}(S)) = R(\mathcal{A}(S)) - \inf_{\mathbf{w} \in \mathcal{W}} R(\mathbf{w})$. We can decompose $\epsilon_{\text{risk}}(\mathcal{A}(S))$ as follows:

$$\epsilon_{\text{risk}}(\mathcal{A}(S)) = [R(\mathcal{A}(S)) - R_S(\mathcal{A}(S))] + [R_S(\mathbf{w}_*) - R(\mathbf{w}_*)] + [R_S(\mathcal{A}(S)) - R_S(\mathbf{w}_*)],$$
(2)

where $\mathbf{w}_* \in \arg\min_{\mathbf{w} \in \mathcal{W}} R(\mathbf{w})$. The first term on the right hand side of (2) is called the estimation error. Since \mathbf{w}_* is fixed, the term $R_S(\mathbf{w}_*) - R(\mathbf{w}_*)$ can be trivially handled by the standard Hoeffding inequality. As a comparison, the estimation of the term $R(\mathcal{A}(S)) - R_S(\mathcal{A}(S))$, also called the generalization error, is much more challenging since $\mathcal{A}(S)$ depends on S. We will develop novel stability analysis to handle this term. The last term $R_S(\mathcal{A}(S)) - R_S(\mathbf{w}_*)$ is called the optimization error and we can bound it by applying optimization theory.

We now introduce some necessary assumptions and definitions. Let $\|\cdot\|_2$ denote the Euclidean norm on \mathbb{R}^d and $\langle\cdot,\cdot\rangle$ denote the corresponding inner product. Given a function $f: \mathcal{W} \to \mathbb{R}$, let $\partial f(\mathbf{w})$ be a subgradient of f at \mathbf{w} . A function f is said to be convex if

Algorithm 1 SGD for Pairwise Learning

Input: Data set $S = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$, step size η , number of iterations T, initial point $\mathbf{w}_1 = 0$ and initial sample $i_1 \in [n]$ from uniform distribution

for t = 1 to T do

Select $i_{t+1} \in [n]$ by uniform distribution

$$\mathbf{w}_{t+1} = \Pi_{\mathcal{W}} \left(\mathbf{w}_t - \frac{\eta}{t} \sum_{k=1}^t \partial \ell(\mathbf{w}_t, \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k}) \right)$$

end for

Output: $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

for any $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$, there holds

$$f(\mathbf{w}') \ge f(\mathbf{w}) + \langle \partial f(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle.$$

A function f is said to be G-Lipschitz continuous if, for any $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$, there holds

$$|f(\mathbf{w}) - f(\mathbf{w}')| \le G||\mathbf{w} - \mathbf{w}'||_2.$$

Throughout this paper, we assume that the (possibly non-smooth) loss function $\ell(\mathbf{w}, \mathbf{z}, \mathbf{z}')$ is nonnegative, convex and G-Lipschitz continuous w.r.t \mathbf{w} .

2.1 Stability and Excess Risk Analysis

In this subsection, we consider the stability and generalization of the SGD algorithms for pairwise learning. The SGD algorithm is described in Algorithm I which has been widely discussed in Lin et al. (2017); Wang et al. (2012); Ying and Zhou (2016). Note that $\Pi_{\mathcal{W}}(\cdot)$ is the projection onto the parameter space \mathcal{W} and $[n] = \{1, \ldots, n\}$. In this subsection, the notation \mathcal{A} denotes Algorithm I

In particular, we will use the uniform argument stability (UAS) (Liu et al., 2017) where its original concept was stated in expectation w.r.t. the internal randomness of \mathcal{A} . We will use its probabilistic version here. Specifically, let $S = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ and $S' = \{\mathbf{z}_1', \dots, \mathbf{z}_n'\}$ be two neighborhood datasets that differ only in one single example. For any $\gamma \in (0, 1)$, \mathcal{A} is called ϵ_{stab} -UAS with probability $1 - \gamma$ if for any neighborhood datasets S and S',

$$\mathbb{P}_{\mathcal{A}}[\|\mathcal{A}(S) - \mathcal{A}(S')\|_{2} > \epsilon_{\text{stab}}] \leq \gamma.$$

We emphasize the probability here is taken over the internal randomness of \mathcal{A} , i.e. the uniform distribution of generating i_t 's.

The following theorem states a high-probability UAS result for Algorithm 1 with non-smooth losses. Here, \mathbf{w}_{t+1} and \mathbf{w}'_{t+1} denote the (t+1)-th iterate of Algorithm 1 based on samples S and S', respectively. And, the notation $\tilde{\mathcal{O}}(\cdot)$ indicates that the bound is up to a logarithmic term.

Theorem 1. Suppose that we run Algorithm $\boxed{1}$ under random selection with replacement for t iterations based on S and S'. Then, with probability $1 - \gamma$ w.r.t. the internal randomness of A, we have, for any S and S', that

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_{2}^{2} \le 4e\eta^{2}G^{2}\left[t + \ln^{2}(et)\right] \times \left(\frac{t}{n} + \ln(1/\gamma) + \sqrt{\frac{t\ln(1/\gamma)}{n}}\right)^{2}$$
(3)

In particular, if $T \geq n$, then the output of Algorithm I is ϵ_{stab} -UAS with high probability where

$$\epsilon_{stab} = \tilde{\mathcal{O}}\Big(\eta\sqrt{T} + \frac{\eta T \ln(T)}{n}\Big).$$

The proof of Theorem $\boxed{1}$ is given in Section $\boxed{3.1}$ This bound matches the result in the pointwise learning with non-smooth losses (Bassily et al., 2020) Lei and Ying, 2020) up to a logarithmic term of T. The proof is motivated by Lei and Ying (2020) in the pointwise case but more involved in pairwise learning. Indeed, the key challenge, in comparison with pointwise learning, is that the sub-gradient estimator at the t-th step depends not only on the current example $\mathbf{z}_{i_{t+1}}$ but also on previous examples $\{\mathbf{z}_{i_k}: k=1,\cdots,t\}$.

To our best knowledge, Shen et al. (2020) is the only available work which considered the stability of SGD in pairwise learning. However, their work required the loss to be Lipschitz continuous and strongly smooth to ensure the non-expansiveness of the gradient update, which is very critical for the proof of the main results there. The non-smoothness assumption in our paper makes the corresponding gradient update no longer non-expansive, and therefore the arguments in Shen et al. (2020) no longer apply. We bypass this obstacle by a refined control of the expansiveness between adjacent steps. To address this dependence issue, the work of Shen et al. (2020) counts the number m of different examples $\mathbf{z}_i \neq \mathbf{z}_i'$ encountered by SGD until iteration t, which obeys a binomial distribution. In contrast, high-probability analysis here for non-smooth loss is more challenging and involved because directly applying concentration inequality to similar binomial distribution yields an undesired estimation. We overcome this hurdle by decomposing the sub-gradients into sum of t pairs of dependent random variables first, and then upper bound this sum by two sums of independent random variables. From this new decomposition, we can apply the Chernoff-type tail bounds to these two sums of independent random variables to get the desired estimation. One can see Section 3.1 for more details.

Based on Theorem 1 and the error decomposition 2, we derive the excess risk bounds for bounded (Theorem 2) and unbounded domains (Theorem 3). To

bound the optimization error, we need the following variant of Rademacher average (Bartlett and Mendelson, 2002)

$$\mathcal{R}_t(\ell \circ \mathcal{W}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \Big[\sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{t} \sum_{k=1}^t \sigma_k \ell(\mathbf{w}, \mathbf{z}_i, \mathbf{z}_{i_k}) \Big].$$
(4)

Here σ_k are Rademacher random variables taking values in $\{\pm 1\}$ with equal probability 1/2, and the expectation is taken over \mathbf{z}_i , \mathbf{z}_{i_k} and σ_k .

Theorem 2. Suppose W is bounded with diameter D. Denote $M = \sup_{\mathbf{z},\mathbf{z}'} \ell(0,\mathbf{z},\mathbf{z}')$. Assume we run Algorithm \mathbb{I} for $T \geq n$ iterations under random selection with replacement rule. Then for any $\gamma \in (0,1)$, with probability at least $1 - \gamma$ w.r.t. the sample S and the internal randomness of A, we have

$$\begin{split} \epsilon_{risk}(\bar{\mathbf{w}}_T) \leq & \frac{4}{T} \sum_{t=1}^{T} \mathcal{R}_t(\ell \circ \mathcal{W}) + \frac{D^2}{2T\eta} + \frac{\eta G^2}{2} + c_2 \sqrt{\frac{\ln(6T/\gamma)}{n}} \\ & + c_1 \eta \lceil \ln(n) \rceil \Big(\sqrt{T} + \frac{\sqrt{3}T \ln(eT) \ln(6/\gamma)}{n} \Big), \end{split}$$

where $c_1 = 100\sqrt{6}e^{3/2}G \max\{1, G\} \ln(6e/\gamma)$ and $c_2 = (6+19e)(M+GD)$.

In particular, if $\mathcal{R}_t(\ell \circ \mathcal{W}) = \mathcal{O}(1/\sqrt{t})$ and we choose $T = n^2$ and $\eta = \mathcal{O}(n^{-3/2})$ then with high probability we have

$$\epsilon_{risk}(\bar{\mathbf{w}}_T) = \tilde{\mathcal{O}}\Big(\frac{\ln^2(n)}{\sqrt{n}}\Big).$$

Theorem 2 is proved in Appendix A.2. Using standard technique (Bartlett and Mendelson, 2002), the Rademacher complexity estimation of $\mathcal{R}_t(\ell \circ \mathcal{W}) = \mathcal{O}(1/\sqrt{t})$ holds true in many cases when \mathcal{X} and \mathcal{W} are bounded (e.g. see Section 4 for concrete examples of AUC maximization and similarity metric learning). It is worthy of mentioning that the choice of $T=n^2$ is consistent with pointwise learning with non-smooth loss (Bassily et al.) 2020; Lei and Ying, 2020).

We can also derive excess generalization bounds for Algorithm \blacksquare even when \mathcal{W} is unbounded. Specifically, let $D = \|\mathbf{w}_*\|_2$ and $\mathcal{W}_D = \{\mathbf{w} \in \mathcal{W} | \|\mathbf{w}\|_2 \leq D\}$. The main idea is to show that the iterate \mathbf{w}_t from Algorithm \blacksquare has an adaptive bound, i.e. $\mathbf{w}_t \in \mathcal{W}_t = \{\mathbf{w} \in \mathcal{W} | \|\mathbf{w}\|_2^2 \leq (G^2 + M)\eta t\}$.

Theorem 3. Denote $M = \sup_{\mathbf{z}, \mathbf{z}'} \ell(0, \mathbf{z}, \mathbf{z}')$ and $D = \|\mathbf{w}_*\|_2$. Suppose we run Algorithm $[\![]\!]$ for $T \geq n$ iterations. For any $\gamma \in (0,1)$, with probability at least $1-\gamma$ w.r.t. the sample S and the internal randomness of A,

we have

$$\begin{split} &\epsilon_{risk}(\bar{\mathbf{w}}_T) \leq \frac{2}{T} \sum_{t=1}^T \left(\mathcal{R}_t(\ell \circ \mathcal{W}_t) + \mathcal{R}_t(\ell \circ \mathcal{W}_D) \right) + \frac{D^2}{2T\eta} + \frac{\eta G^2}{2} \\ &+ c_4 \sqrt{\eta \ln(6T/\gamma)} + c_5 \sqrt{\frac{\eta T \ln(6e/\gamma)}{n}} + c_3 \sqrt{\frac{\ln(6T/\gamma)}{n}} \\ &+ c_1 \eta \lceil \ln(n) \rceil \left(\sqrt{T} + \frac{4T \ln(eT) \sqrt{\ln(6n/\gamma)}}{n} \right), \end{split}$$

where $c_1 = 100\sqrt{6}e^{3/2}G \max\{1,G\} \ln(6e/\gamma)$, $c_3 = (7 + 12\sqrt{2}e)M + 4GD + 16eG$, $c_4 = 3G\sqrt{G^2 + 2M}$ and $c_5 = 12\sqrt{2}eG\sqrt{G^2 + 2M}$.

In particular, if $\mathcal{R}_t(\ell \circ \mathcal{W}_t) = \mathcal{O}(\eta \sqrt{t})$ and $\mathcal{R}_t(\ell \circ \mathcal{W}_D) = \mathcal{O}(1/\sqrt{t})$ and we choose $T = n^{4/3}$ and $\eta = \mathcal{O}(n^{-1})$, then with high probability we have

$$\epsilon_{risk}(\bar{\mathbf{w}}_T) = \tilde{\mathcal{O}}\Big(\frac{\ln^2(n)}{n^{1/3}}\Big).$$

Theorem 3 is proved in Appendix A.3. In particular, one can show that the Rademacher complexity can be estimated using standard technique (Bartlett and Mendelson, 2002) such that $\mathcal{R}_t(\ell \circ \mathcal{W}_D) = \mathcal{O}(D/\sqrt{t})$ when \mathcal{X} is a bounded domain. Therefore by the definition of \mathcal{W}_t one can similarly show that $\mathcal{R}_t(\ell \circ \mathcal{W}_t) =$ $\mathcal{O}(\eta t/\sqrt{t}) = \mathcal{O}(\eta\sqrt{t})$. One can see more discussion on such estimation in Section 4. Therefore, Theorem 3 mainly differs from Theorem 2 in the additional $\tilde{\mathcal{O}}(\sqrt{\eta T/n})$ term where $T \geq n$. This is due to the unboundedness of W. Our excess risk bound is consistent with the results in Lin et al. (2016) in the pointwise setting (up to a logarithmic term), where the authors studied SGD for non-smooth loss functions in the pointwise setting using uniform convergence. However, the bound there is given in expectation while we have provided a high-probability bound.

2.2 Differentially Private Pairwise Learning

We show the implication of stability analysis in analyzing differentially private SGD in pairwise learning. We start by introducing the notion of differential privacy.

Definition 1 (Differential Privacy (Dwork et al., 2006)). A (randomized) algorithm \mathcal{A} is called (ϵ, δ) -differentially private (DP) if, for all neighboring datasets S, S' differing by only one example and for all events O in the output space of \mathcal{A} , the following holds

$$\mathbb{P}[\mathcal{A}(S) \in O] \le e^{\epsilon} \mathbb{P}[\mathcal{A}(S') \in O] + \delta.$$

There are other forms of differential privacy such as Gaussian differential privacy (Bu et al.) 2020; Dong et al., 2019). In this paper we restrict our attention

Algorithm 2 Private SGD for Pairwise Learning with Output Perturbation

Input: Private dataset $S = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$, privacy parameter ϵ, δ , stepsize η , number of iterations T, initial point $\mathbf{w}_1 = 0$ and initial sample $i_1 \in [n]$ from uniform distribution

for t = 1 to T do

Select $i_{t+1} \in [n]$ from uniform distribution

$$\mathbf{w}_{t+1} = \Pi_{\mathcal{W}} \left(\mathbf{w}_t - \frac{\eta}{t} \sum_{k=1}^t \partial \ell(\mathbf{w}_t, \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k}) \right)$$

end for

 $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Sample $\mathbf{u} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ with σ^2 being given by (5)

Output: $\mathbf{w}_{\mathrm{priv}} = \Pi_{\mathcal{W}}(\bar{\mathbf{w}}_T + \mathbf{u})$

to the standard DP mentioned above. In particular, we consider Gaussian mechanism (Dwork et al.) [2006), i.e. given any query function $q: S^n \to \mathbb{R}^d$, let $\mathcal{A}(S) = q(S) + \mathbf{u}$ where $\mathbf{u} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ with \mathbf{I}_d being the identical matrix. For all neighborhood datasets S, S' that differ by one example, the ℓ_2 -sensitivity Δ of the query function q is defined as $\Delta(q) = \sup_{S,S'} \|q(S) - q(S')\|_2$.

We develop a private version of SGD for pairwise learning. In this subsection, the notation \mathcal{A} denotes Algorithm $\boxed{2}$. The idea is to add Gaussian noise to the output of the non-private Algorithm $\boxed{1}$. In return, Algorithm $\boxed{2}$ is guaranteed to be (ϵ, δ) -DP by properly choosing σ as shown below.

Theorem 4. Given the total number of iterations T, for any privacy budget $\epsilon > 0$ and $\delta > 0$, Algorithm 2 satisfies (ϵ, δ) -differential privacy with

$$\sigma^2 = \frac{8e\eta^2 G^2 \ln(2.5/\delta)}{\epsilon^2} \left(T + \frac{3T^2 \ln^2(eT) \ln^2(2/\delta)}{n^2} \right). (5)$$

The proof of Theorem 4 is given in Section 3.2. The goal here is to guarantee privacy with the added noise being as small as possible. The key observation is the UAS of the non-private output $\bar{\mathbf{w}}_T$ can be used to quantify the high-probability sensitivity of the query function $q(S) = \bar{\mathbf{w}}_T$. Specifically, subsampling forms an event of probability measure $1 - \delta/2$ under which a small sensitivity $\tilde{\mathcal{O}}(\eta\sqrt{T} + \eta T \ln(T)/n)$ holds true. Hence, under this event, we only need to add noise with $\sigma = \tilde{\mathcal{O}}((\eta\sqrt{T} + \eta T \ln(T)) \ln(2/\delta)/(n\epsilon))$ to guarantee a slightly restrictive $(\epsilon, \delta/2)$ -DP. Therefore the algorithm is $(\epsilon; \delta)$ -DP over the whole event space. Wu et al. (2017) studied differential private SGD by output perturbation method in the pointwise learning setting and they also utilized the idea of bounding sensitivity by UAS. However, they considered the stability and sensitivity regardless of the randomness of the algorithm, which is not suitable for high probability analysis of utility bound later. In contrast, our technique

can also be applied to derive privacy guarantee and high probability utility in pointwise learning. Huai et al. (2020) also studied the sensitivity of SGD for Pairwise learning. However, they focused on the online setting where the data arrives in a streaming manner, and hence the different example between S and S' will only appear once in the algorithm. While in our stochastic setting the different example can be used more than once by subsampling, it is more challenging to measure the sensitivity. Moreover, their analysis depends on the strong smoothness of the loss function while we allow the loss function to be non-smooth.

In order to derive the utility bound of Algorithm 2, we need a new error decomposition scheme as follow

$$\epsilon_{\text{risk}}(\mathbf{w}_{\text{priv}}) = R(\mathbf{w}_{\text{priv}}) - R(\mathbf{w}_{*})
= R(\mathbf{w}_{\text{priv}}) - R(\bar{\mathbf{w}}_{T}) + R(\bar{\mathbf{w}}_{T}) - R(\mathbf{w}_{*}),$$
(6)

where $R(\bar{\mathbf{w}}_T) - R(\mathbf{w}_*)$ measures the excess risk incurred by the non-private output $\bar{\mathbf{w}}_T$ (Algorithm 1) and $R(\mathbf{w}_{\text{priv}}) - R(\bar{\mathbf{w}}_T)$ measures the effect of perturbation by adding random noises. The utility bound is given as follow.

Theorem 5. Suppose W is bounded with diameter D. Consider Algorithm 2 for T iterations under random selection with replacement rule. For any privacy budget $\epsilon > 0$, $\delta > 0$, and for any $\gamma \in (\max\{4\delta, \exp(-d/8)\}, 1)$, with probability at least $1-\gamma$, we have

$$\epsilon_{risk}(\mathbf{w}_{priv}) \leq \frac{4}{T} \sum_{t=1}^{T} \mathcal{R}_{t}(\ell \circ \mathcal{W}) + \frac{D^{2}}{2T\eta} + \frac{\eta G^{2}}{2} + c_{2} \sqrt{\frac{\ln(6T/\gamma)}{n}} + 2G\sigma\sqrt{d}\ln^{1/4}(4/\gamma). + c_{1}\eta\lceil\ln(n)\rceil\left(\sqrt{T} + \frac{\sqrt{3}T\ln(eT)\ln(2/\delta)}{n}\right),$$

where $c_1 = 100\sqrt{6}e^{3/2}G \max\{1, G\} \ln(6e/\gamma)$ and $c_2 = (6+19e)(M+GD)$.

In particular, letting σ satisfy (5) and choosing $T = n^2$ and $\eta = \mathcal{O}(n^{-3/2})$, then with high probability we have

$$\epsilon_{risk}(\mathbf{w}_{priv}) = \tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\sqrt{n}\epsilon}\right).$$

Theorem 5 is proved in Appendix A.4. The difference compared to Theorem 2 is the additional $\tilde{\mathcal{O}}(\sigma\sqrt{d})$ term caused by $R(\mathbf{w}_{\text{priv}})-R(\bar{\mathbf{w}}_T)$ in 6. The utility bound $\tilde{\mathcal{O}}(\sqrt{d}/(\sqrt{n}\epsilon))$ matches that of the output perturbation for pairwise learning studied in 6 Huai et al. 6 Which, however, requires the loss to be both strongly smooth and Lipschitz continuous. Our analysis only needs the loss to be Lipschitz continuous.

3 Main Proofs for Theorems 1 and 4

In this section, we provide technical proofs for Theorems and and Proofs of other Theorems can be found in the Appendix. Throughout this section, we let $\hat{L}_{t+1}(\mathbf{w}_t)$ denote the accumulated loss until $\mathbf{z}_{i_{t+1}}$ is revealed. i.e. $\hat{L}_{t+1}(\mathbf{w}_t) = \frac{1}{t} \sum_{k=1}^t \ell(\mathbf{w}_t, \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k})$.

3.1 Proof of Theorem 1

To prove Theorem we need the following Chernoff's bound for a summation of independent Bernoulli random variables (Wainwright, 2019).

Lemma 1 (Chernoff bound for Bernoulli vector). Let X_1, \ldots, X_t be independent random variables taking values in $\{0,1\}$. Let $X = \sum_{j=1}^t X_j$ and $\mu = \mathbb{E}[X]$. Then for any $\tilde{\gamma} > 0$, with probability at least $1 - \exp(-\mu \tilde{\gamma}^2/(2 + \tilde{\gamma}))$ we have $X \leq (1 + \tilde{\gamma})\mu$.

Proof of Theorem [7]. Without loss of generality, assume that S and S' differs in n-th position. Denote $\delta_{t+1,k} = \partial \ell(\mathbf{w}_t, \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k}) - \partial \ell(\mathbf{w}_t', \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k})$ and $\delta'_{t+1,k} = \partial \ell(\mathbf{w}_t, \mathbf{z}_{i_{t+1}}, \mathbf{z}_{i_k}) - \partial \ell(\mathbf{w}_t', \mathbf{z}_{i_{t+1}}', \mathbf{z}_{i_k}')$. The following recursive inequality holds

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_{2}^{2} = \|\mathbf{w}_{t} - \eta \partial \hat{L}_{t+1}(\mathbf{w}_{t}) - \mathbf{w}'_{t} + \eta \partial \hat{L}'_{t+1}(\mathbf{w}'_{t})\|_{2}^{2}$$

$$= \|\mathbf{w}_{t} - \mathbf{w}'_{t} - \frac{\eta}{t} \sum_{k=1}^{t} \delta'_{t+1,k} \|_{2}^{2}$$

$$\leq \frac{1}{t} \sum_{k=1}^{t} \|\mathbf{w}_{t} - \mathbf{w}'_{t} - \eta \delta'_{t+1,k} \|_{2}^{2}.$$
(7)

Now we estimate the term on the right hand side of (7) by considering two cases. For the case $i_{t+1} \neq n$ and $i_k \neq n$, we have $\mathbf{z}_{i_{t+1}} = \mathbf{z}'_{i_{t+1}}$ and $\mathbf{z}_{i_k} = \mathbf{z}'_{i_k}$. Then

$$\|\mathbf{w}_{t} - \mathbf{w}'_{t} - \eta \delta_{t+1,k}\|^{2}$$

$$= \|\mathbf{w}_{t} - \mathbf{w}'_{t}\|_{2}^{2} + \eta^{2} \|\delta_{t+1,k}\|_{2}^{2} - 2\eta \langle \mathbf{w}_{t} - \mathbf{w}'_{t}, \delta_{t+1,k} \rangle$$

$$\leq \|\mathbf{w}_{t} - \mathbf{w}'_{t}\|_{2}^{2} + 4\eta^{2} G^{2},$$

where the last inequality holds because ℓ is G-Lipschitz and convex. If $i_{t+1} = n$ or $i_k = n$, then $\mathbf{z}_{i_{t+1}} \neq \mathbf{z}'_{i_{t+1}}$ or $\mathbf{z}_{i_k} \neq \mathbf{z}'_{i_k}$. It follows from the Young's inequality that for any p > 0

$$\begin{aligned} &\|\mathbf{w}_{t} - \mathbf{w}_{t}' - \eta \delta_{t+1,k}''\|^{2} \\ &\leq (1+p) \|\mathbf{w}_{t} - \mathbf{w}_{t}'\|_{2}^{2} + (1+1/p)\eta^{2} \|\delta_{t+1,k}'\|_{2}^{2} \\ &\leq (1+p) \|\mathbf{w}_{t} - \mathbf{w}_{t}'\|_{2}^{2} + 4(1+1/p)\eta^{2}G^{2}. \end{aligned}$$

Combining the above two inequalities together and let $Y_t = \frac{1}{t} \sum_{k=1}^{t} \mathbb{I}_{[i_{t+1}=n \vee i_k=n]}$, we have

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2^2 \le (1 + pY_t) \|\mathbf{w}_t - \mathbf{w}'_t\|_2^2 + 4(1 + Y_t/p)\eta^2 G^2.$$

Applying the above inequality recursively we have

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_{2}^{2} \stackrel{(a)}{\leq} \sum_{j=1}^{t} \prod_{l=j+1}^{t} (1+pY_{l})(4+4Y_{j}/p)\eta_{j}^{2}G^{2}$$

$$\stackrel{(b)}{\leq} \sum_{j=1}^{t} \prod_{l=j+1}^{t} (1+p)^{Y_{l}} (4+4Y_{j}/p)\eta_{j}^{2}G^{2}$$

$$\stackrel{(c)}{\leq} (1+p)^{\sum_{l=1}^{t} Y_{l}} \eta^{2}G^{2} (4t+4\sum_{l=1}^{t} Y_{l}/p), \tag{8}$$

where (a) is due to the recursive relation, (b) is due to $1 + ax \leq (1 + a)^x$ for a > 0 and $x \geq 0$ and (c) inequality is due to $\prod_{i=a}^b x^i \leq x^{\sum_{i=1}^b i}$ for $a \geq 1$. We note that Y_1, \dots, Y_t are dependent variables, but the sum of Y_i 's has the following decomposition:

$$\begin{split} \sum_{l=1}^{t} Y_{l} &= \sum_{l=1}^{t} \frac{1}{l} \sum_{k=1}^{l} \mathbb{I}_{[i_{l+1} = n \vee i_{k} = n]} \leq \sum_{l=1}^{t} \frac{1}{l} \sum_{k=1}^{l} \left(\mathbb{I}_{[i_{l+1} = n]} + \mathbb{I}_{[i_{k} = n]} \right) \\ &= \sum_{l=1}^{t} \mathbb{I}_{[i_{l+1} = n]} + \sum_{l=1}^{t} \frac{1}{l} \sum_{k=1}^{l} \mathbb{I}_{[i_{k} = n]} \\ &\leq \sum_{l=1}^{t} \mathbb{I}_{[i_{l+1} = n]} + \ln(t) \sum_{k=1}^{l} \mathbb{I}_{[i_{k} = n]} \leq \ln(et) \sum_{k=1}^{t+1} \mathbb{I}_{[i_{k} = n]}. \end{split}$$

Applying Lemma [1] with $X_k = \mathbb{I}_{[i_k=n]}$ and $X = \sum_{k=1}^t X_k$, with probability at least $1 - \gamma$, we have

$$\sum_{k=1}^{t+1} \mathbb{I}_{[i_k=n]} \le \frac{t+1}{n} + \ln(1/\gamma) + \sqrt{\frac{(t+1)\ln(1/\gamma)}{n}}.$$

For the simplicity of notation, let $c_{\gamma,t} = \ln(et)((t+1)/n + \ln(1/\gamma) + \sqrt{(t+1)\ln(1/\gamma)/n})$. Plugging the above inequality back into (8), we derive the following inequality with probability $1-\gamma$

$$\|\mathbf{w}_{t+1} - \mathbf{w}'_{t+1}\|_2^2 \le 4\eta^2 G^2 (1+p)^{c_{\gamma,t}} (t + c_{\gamma,t}/p).$$

By selecting $p = 1/c_{\gamma,t}$ in the above equality, we have $(1+p)^{c_{\gamma,t}} \leq e$. Therefore we have proved (3) in Theorem 1. Now, since the bound on left hand side of (3) is monotonically increasing, with probability $1-\gamma$, we have

$$\|\bar{\mathbf{w}}_{T} - \bar{\mathbf{w}}_{T}'\|_{2}^{2} \leq \frac{1}{T} \sum_{t=1}^{T} \|\mathbf{w}_{T} - \mathbf{w}_{T}'\|_{2}^{2}$$

$$\leq 4e\eta^{2} G^{2} \left(T + \frac{3T^{2} \ln^{2}(eT) \ln^{2}(1/\gamma)}{n^{2}}\right), \tag{9}$$

where we have used the fact that $T \geq n$. Therefore the ϵ_{stab} -UAS bound holds by calling the convexity of ℓ_2 -norm.

3.2 Proof of Theorem 4

In order to establish the privacy guarantee of Algorithm 2 we need the following lemmas. The first lemma characterizes the necessary scale of σ of Gaussian mechanism (Dwork et al., 2014).

Lemma 2 (Gaussian mechanism). For a Gaussian mechanism $\mathcal{A}(S) = q(S) + \mathbf{u}$ with $\mathbf{u} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$, if q has ℓ_2 -sensitivity $\Delta(q)$ and assume that $\sigma \geq \sqrt{2 \ln(1.25/\delta)} \Delta(q)/\epsilon$, then \mathcal{A} yields (ϵ, δ) -DP.

The next lemma indicates that differential privacy is immune to post-processing (Dwork et al., 2014).

Lemma 3 (Post-processing). Let $A : \mathbb{Z}^n \to W$ be a (randomized) algorithm that is (ϵ, δ) -DP. Let $f : W \to W$ be an arbitrary randomized mapping. Then $f \circ A : \mathbb{Z}^n \to W$ is (ϵ, δ) -DP.

Proof of Theorem \Box Consider the mechanism $\mathcal{A}'_T = \bar{\mathbf{w}}_T + \mathbf{u}$ and for any S, S', consider the ℓ_2 -sensitivity $\Delta_T = ||\bar{\mathbf{w}}_T - \bar{\mathbf{w}}'_T||_2$. Let $I = \{i_1, \dots, i_T\}$ be the sequence of sampling after T iterations in Algorithm \Box Choosing $\gamma = \delta/2$ in Equation \Box , then the event

$$E = \left\{ I | \Delta_T^2 \le 4e\eta^2 G^2 \left(T + \frac{3T^2 \ln^2(eT) \ln^2(2/\delta)}{n^2} \right) \right\}$$

satisfies $\mathbb{P}[I \in E] \geq 1 - \delta/2$. When $I \in E$, Lemma 2 implies \mathcal{A}_T' satisfies $(\epsilon, \delta/2)$ -DP when

$$\sigma = \frac{\sqrt{2\ln(2.5/\delta)}\Delta_T}{\epsilon}.$$

Furthermore, by Lemma 3, the final output $\mathbf{w}_{\text{priv}} = \Pi_{\mathcal{W}}(\mathcal{A}_T')$ also satisfies $(\epsilon, \delta/2)$ -DP. Therefore, for any $\epsilon > 0$ and any event O in the output space of \mathbf{w}_{priv} ,

$$\begin{split} & \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S) \in O\big] = \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S) \in O \cap I \in E\big] \\ & + \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S) \in O \cap I \notin E\big] \\ & = \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S) \in O | I \in E\big] \mathbb{P}\big[I \in E\big] \\ & + \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S) \in O | I \notin E\big] \mathbb{P}\big[I \notin E\big] \\ & \leq \Big(e^{\epsilon} \mathbb{P}[\mathbf{w}_{\mathrm{priv}}(S') \in O | I \in E\big] + \frac{\delta}{2}\Big) \mathbb{P}[I \in E] + \frac{\delta}{2} \\ & \leq e^{\epsilon} \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S') \in O \cap I \in E\big] + \frac{\delta}{2} + \frac{\delta}{2} \\ & \leq e^{\epsilon} \mathbb{P}\big[\mathbf{w}_{\mathrm{priv}}(S') \in O \cap I \in E\big] + \delta \end{split}$$

where the first inequality is because when $I \in E$, \mathbf{w}_{priv} satisfies $(\epsilon, \delta/2)$ -DP and the fact $\mathbb{P}[I \notin E] \leq \delta/2$, the second inequality is by the definition of conditional probability. The proof is complete.

4 Applications

In this section, we illustrate our main results in the above sections by considering two concrete examples of pairwise learning, namely AUC maximization and similarity metric learning. According to Theorems 2 and 5, the key here is to estimate the Rademacher complexity defined by (4).

AUC Maximization. AUC maximization aims to learn a ranking function $h_{\mathbf{w}}$ defined by $h_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = \mathbf{w}^{\top}(\mathbf{x} - \mathbf{x}')$. One expects $h_{\mathbf{w}}$ will rank positive examples higher than negative examples, i.e. $\mathbf{w}^{\top}(\mathbf{x} - \mathbf{x}') \geq 0$ for y = 1 and y' = -1. Using the hinge loss $\ell(\mathbf{w}, \mathbf{z}, \mathbf{z}') = (1 - h_{\mathbf{w}}(\mathbf{x}, \mathbf{x}'))_{+} \mathbb{I}_{[y=1 \wedge y'=-1]}$, AUC maximization can be formulated as

$$\min_{\mathbf{w} \in \mathcal{W}} \mathbb{E}_{\mathbf{z}, \mathbf{z}'} \left[(1 - \mathbf{w}^{\top} (\mathbf{x} - \mathbf{x}'))_{+} \mathbb{I}_{[y = 1 \wedge y' = -1]} \right].$$
 (10)

Denote $\kappa = \sup_{\mathbf{x}} ||\mathbf{x}||_2$. The Rademacher complexity defined by (4) for AUC maximization is given in the following lemma.

Lemma 4. Given the parameter space $W = \{ \mathbf{w} \in \mathbb{R}^d | \|\mathbf{w}\|_2 \leq D \}$, the Rademacher complexity of $\mathcal{H} = \{ h_{\mathbf{w}} | \mathbf{w} \in \mathcal{W} \}$ can be upper bounded by $\mathcal{R}_t(\mathcal{H}) \leq 2D\kappa/\sqrt{t}$.

Note in the case of (10), it is easy to check $R_t(\ell \circ \mathcal{H}) \leq 4GD\kappa/\sqrt{t}$ by Ledoux-Talagrand inequality (Ledoux and Talagrand) 2013). Combining this lemma with Theorems 2 and 5, one can derive the following excess risk and utility bound for Algorithms 1 and 2 in the context of non-smooth AUC maximization.

Corollary 1. Consider the problem of AUC maximization (10). If one runs Algorithm I with $T = n^2$ and $\eta = \mathcal{O}(n^{-3/2})$, then, with high probability we have

$$\epsilon_{risk}(\bar{\mathbf{w}}_T) = \tilde{\mathcal{O}}\left(\sqrt{\frac{\kappa}{n}}\right).$$

Corollary 2. For the problem of AUC maximization (10), if one runs Algorithm 2 with $T = n^2$, $\eta = \mathcal{O}(n^{-3/2})$ and σ given by (5), then, with high probability we have

$$\epsilon_{risk}(\mathbf{w}_{priv}) = \tilde{\mathcal{O}}\left(\frac{\sqrt{\kappa d}}{\sqrt{n}\epsilon}\right).$$

Similarity Metric Learning. We now turn to another notable example of pairwise learning called similarity metric learning. It aims to learn a (squared) Mahalanobis distance metric which is defined by $h_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^{\top} \mathbf{w}(\mathbf{x} - \mathbf{x}')$ parametrized by a positive semi-definite matrix $\mathbf{w} \in \mathbb{R}^{d \times d}$. The intuition behind similarity metric learning is that the distance between samples from the same class should be small and the distance between examples from distinct classes should be large. Using the hinge loss $\ell(\mathbf{w}, \mathbf{z}, \mathbf{z}') = (1 + \tau(y, y')h_{\mathbf{w}}(\mathbf{x}, \mathbf{x}'))_+$, it can be formulated as

$$\min_{\mathbf{w} \in \mathcal{W}} \mathbb{E}_{\mathbf{z}, \mathbf{z}'} \left[(1 + \tau(y, y')(\mathbf{x} - \mathbf{x}')^{\top} \mathbf{w}(\mathbf{x} - \mathbf{x}'))_{+} \right], \quad (11)$$

where $\tau(y, y') = 1$ if y = y' and -1 otherwise.

Lemma 5. Consider the parameter space defined via the nuclear norm $W = \{ \mathbf{w} \in \mathbb{R}^{d \times d}, \|\mathbf{w}\|_{S_1} \leq D \}$, where $\|\mathbf{w}\|_{S_1}$ denotes the nuclear norm of a matrix \mathbf{w} . The complexity of $\mathcal{H} = \{ h_{\mathbf{w}} : \mathbf{w} \in \mathcal{W} \}$ is bounded by

$$\mathcal{R}_t(\mathcal{H}) = \mathcal{O}\left(\frac{D\left\|\mathbb{E}[\|X\|_2^2 X X^\top]\right\|_{S_\infty}^{\frac{1}{2}} \sqrt{\log d}}{\sqrt{t}}\right), \quad (12)$$

where $\|\cdot\|_{S_{\infty}}$ denotes the largest singular value.

The proof of Lemma 5 is postponed to Appendix A.5.

As direct corollaries of Lemma 5, we can derive generalization bounds for metric learning from Theorems 2 and 5. For brevity, denote $\chi = \left\| \mathbb{E}[\|X\|_2^2 X X^\top] \right\|_{S_\infty}$. We derive the following results of SGD for pairwise learning in the context of non-smooth metric learning.

Corollary 3. Consider the similarity metric learning problem (11). If one runs Algorithm 1 for $T = n^2$ and $\eta = \mathcal{O}(n^{-3/2})$, then, with high probability we have

$$\epsilon_{risk}(\bar{\mathbf{w}}_T) = \tilde{\mathcal{O}}\left(\sqrt{\frac{\chi \log(d)}{n}}\right).$$

Corollary 4. Consider the similarity metric learning problem (11). If one runs Algorithm 2 with $T = n^2$, $\eta = \mathcal{O}(n^{-3/2})$ and σ given by (5), then, with high probability we have

$$\epsilon_{risk}(\mathbf{w}_{priv}) = \tilde{\mathcal{O}}\left(\frac{\sqrt{\chi d \log(d)}}{\sqrt{n}\epsilon}\right).$$

Remark 1. We now show the advantage of our result as compared to the existing results. Based on the argument in Lei and Ying (2016), it can be shown

$$\mathcal{R}_t(\mathcal{H}) = \mathcal{O}\left(\frac{D \sup_{\mathbf{x}} \|\mathbf{x}\|^2 \sqrt{\log d}}{\sqrt{t}}\right). \tag{13}$$

The difference between (12) and (13) is that we replace $\sup_{\mathbf{x}} \|\mathbf{x}\|^2$ by the term $\|\mathbb{E}[\|X\|_2^2 X X^{\top}]\|_{S_{\infty}}^{\frac{1}{2}}$. Notice $\|\mathbb{E}[\|X\|^2 X X^{\top}]\|_{S_{\infty}} \geq \frac{1}{d} \mathrm{tr} \left(\mathbb{E}[X X^{\top} X X^{\top}]\right) = \frac{1}{d} \mathbb{E}\left[\|X\|_2^4\right] \gtrsim d^2$, then the upper bound of (12) satisfies the relation $\gtrsim \sqrt{d \log d}/\sqrt{t}$ and in the extreme case this lower bound can be achieved within a constant factor. As a comparison, the upper bound in (13) satisfies the relation $\gtrsim d\sqrt{(\log d)/t}$. That is, our argument outperforms the existing results by enjoying a milder dependency on the dimensionality for using nuclear-norm constraints, which is appealing in the high-dimensional setting. If we use Frobenius-norm constraint in defining $\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^{d \times d}, \|\mathbf{w}\|_F \leq D_2\}$, then one can show that $\mathcal{R}_t(\mathcal{H}) = \mathcal{O}(D_2 \sup_{\mathbf{x}} \|\mathbf{x}\|^2/\sqrt{t})$ (Lei and Ying)

2016). This matches the bound (13) within a logarithmic factor except that D there is replaced by D_2 . Since $\|\mathbf{w}\|_F \leq \|\mathbf{w}\|_{S_1}$, the argument in Lei and Ying (2016) leads to a misleading argument that Frobenius-norm constraint is always preferable to the nuclear-norm constraint. It was posed as an open question on whether one can derive a generalization bound for similarity metric learning showing the advantage of nuclear-norm constraint over Frobenius-norm constraint (Cao et al., 2016). We provide an affirmative solution to this open question in Lemma 5.

5 Conclusions

In this paper, we provide the first-ever-known stability analysis of SGD for pairwise learning with nonsmooth losses and obtain optimal excess risk bounds $\mathcal{O}(1/\sqrt{n})$. We extend our analysis to unbounded parameter space and achieve a rate of $\tilde{\mathcal{O}}(n^{-1/3})$. We apply our stability results to study differentially private SGD algorithms in pairwise learning. Our output perturbation method achieves utility bound $\mathcal{O}(\sqrt{d/(\sqrt{n}\epsilon)})$, which matches the previous results in Huai et al. (2020) for smooth losses. Finally, we provide two examples to illustrate our stability and differential privacy results. In particular, the analysis for the example of metric learning shows the advantage of nuclear norm constraint over Frobenius norm constraint which solved an open question raised in Cao et al. (2016).

Here we only considered SGD with replacement. It would be interesting to extend our analysis to SGD without replacement which is drawing increasing interests. The utility bound is suboptimal as compared with pointwise learning with non-smooth losses, which is $\tilde{\mathcal{O}}\left(\max\{1/\sqrt{n},\sqrt{d}/(n\epsilon)\}\right)$. It remains an open question to us if the same bound can be achieved in pairwise learning.

Acknowledgement

The authors would like to thank the anonymous reviewers for their constructive comments and suggestions. This work is supported by NSF grants IIS-1816227 and IIS-2008532. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agency.

References

Agarwal, S., Dugar, D., and Sengupta, S. (2010). Ranking chemical structures for drug discovery: a new machine learning approach. *Journal of chemical information and modeling*, 50(5):716–731.

- Agarwal, S. and Niyogi, P. (2009). Generalization bounds for ranking algorithms via algorithmic stability. *Journal of Machine Learning Research*, 10(Feb):441–474.
- Bach, F. and Moulines, E. (2013). Non-strongly-convex smooth stochastic approximation with convergence rate o (1/n). In Advances in neural information processing systems, pages 773–781.
- Bartlett, P. L. and Mendelson, S. (2002). Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482.
- Bassily, R., Feldman, V., Guzmán, C., and Talwar, K. (2020). Stability of stochastic gradient descent on nonsmooth convex losses.
- Bassily, R., Feldman, V., Talwar, K., and Thakurta, A. G. (2019). Private stochastic convex optimization with optimal rates. In Advances in Neural Information Processing Systems, pages 11279–11288.
- Bottou, L. and Cun, Y. L. (2004). Large scale online learning. In Advances in neural information processing systems.
- Bu, Z., Dong, J., Long, Q., and Su, W. J. (2020). Deep learning with gaussian differential privacy. *Harvard data science review*, 2020(23).
- Cao, Q., Guo, Z.-C., and Ying, Y. (2016). Generalization bounds for metric and similarity learning. *Machine Learning*, 102(1):115–132.
- Cesa-Bianchi, N., Conconi, A., and Gentile, C. (2004). On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057.
- Charles, Z. and Papailiopoulos, D. (2018). Stability and generalization of learning algorithms that converge to global optima. In *International Conference on Machine Learning*, pages 745–754. PMLR.
- Clémençon, S., Lugosi, G., and Vayatis, N. (2008). Ranking and empirical minimization of u-statistics. The Annals of Statistics, pages 844–874.
- De la Pena, V. and Giné, E. (2012). Decoupling: from dependence to independence. Springer Science & Business Media.
- Dong, J., Roth, A., and Su, W. J. (2019). Gaussian differential privacy. arXiv preprint arXiv:1905.02383.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer.
- Dwork, C., Roth, A., et al. (2014). The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science, 9(3–4):211–407.

- Feldman, V., Koren, T., and Talwar, K. (2020). Private stochastic convex optimization: Optimal rates in linear time. arXiv preprint arXiv:2005.04763.
- Gao, W., Jin, R., Zhu, S., and Zhou, Z.-H. (2013). One-pass auc optimization. In *International Conference on Machine Learning*, pages 906–914.
- Hardt, M., Recht, B., and Singer, Y. (2016). Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Ma*chine Learning, pages 1225–1234.
- Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30.
- Hu, T., Fan, J., Wu, Q., and Zhou, D.-X. (2013). Learning theory approach to minimum error entropy criterion. *Journal of Machine Learning Research*, 14(Feb):377–397.
- Huai, M., Wang, D., Miao, C., Xu, J., and Zhang, A. (2020). Pairwise learning with differential privacy guarantees. In AAAI.
- Jin, R., Wang, S., and Zhou, Y. (2009). Regularized distance metric learning: Theory and algorithm. In Advances in neural information processing systems, pages 862–870.
- Kar, P., Sriperumbudur, B., Jain, P., and Karnick, H. (2013). On the generalization ability of online learning algorithms for pairwise loss functions. In *International Conference on Machine Learning*, pages 441–449.
- Kuzborskij, I. and Lampert, C. (2018). Datadependent stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 2815–2824. PMLR.
- Lacoste-Julien, S., Schmidt, M., and Bach, F. (2012). A simpler approach to obtaining an o (1/t) convergence rate for the projected stochastic subgradient method. arXiv preprint arXiv:1212.2002.
- Ledoux, M. and Talagrand, M. (2013). Probability in Banach Spaces: isoperimetry and processes. Springer Science & Business Media.
- Lei, Y., Ledent, A., and Kloft, M. (2020). Sharper generalization bounds for pairwise learning. Advances in Neural Information Processing Systems, 33.
- Lei, Y. and Ying, Y. (2016). Generalization analysis of multi-modal metric learning. *Analysis and Applications*, 14(04):503–521.
- Lei, Y. and Ying, Y. (2020). Fine-grained analysis of stability and generalization for stochastic gradient descent. In *ICML*.
- Lin, J., Camoriano, R., and Rosasco, L. (2016). Generalization properties and implicit regularization for

- multiple passes sgm. In *International Conference on Machine Learning*, pages 2340–2348.
- Lin, J., Lei, Y., Zhang, B., and Zhou, D.-X. (2017). Online pairwise learning algorithms with convex loss functions. *Information Sciences*, 406:57–70.
- Liu, T., Lugosi, G., Neu, G., and Tao, D. (2017). Algorithmic stability and hypothesis complexity. In *ICML*.
- Lust-Piquard, F. and Pisier, G. (1991). Non commutative khintchine and paley inequalities. *Arkiv för matematik*, 29(1-2):241–260.
- Mukherjee, S. and Wu, Q. (2006). Estimation of gradients and coordinate covariation in classification. *Journal of Machine Learning Research*, 7(Nov):2481–2514.
- Mukherjee, S. and Zhou, D.-X. (2006). Learning coordinate covariances via gradients. *Journal of Machine Learning Research*, 7(Mar):519–549.
- Qiu, R. and Wicks, M. (2014). Cognitive networked sensing and big data. Springer.
- Rakhlin, A., Shamir, O., and Sridharan, K. (2012). Making gradient descent optimal for strongly convex stochastic optimization. In Proceedings of the 29th International Conference on Machine Learning, pages 449–456.
- Rejchel, W. (2012). On ranking and generalization bounds. *Journal of Machine Learning Research*, 13:1373–1392.
- Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. (2009). Stochastic convex optimization. In COLT.
- Shen, W., Yang, Z., Ying, Y., and Yuan, X. (2020). Stability and optimization error of stochastic gradient descent for pairwise learning. *Analysis and Applications*, 18(05):887–927.
- Tropp, J. A. (2015). An introduction to matrix concentration inequalities. arXiv preprint arXiv:1501.01571.
- Verma, N. and Branson, K. (2015). Sample complexity of learning mahalanobis distance metrics. In Advances in neural information processing systems, pages 2584–2592.
- Wainwright, M. J. (2019). *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Wang, Y., Khardon, R., Pechyony, D., and Jones, R. (2012). Generalization bounds for online learning algorithms with pairwise loss functions. In *Conference* on *Learning Theory*, pages 13–1.

- Weinberger, K. Q. and Saul, L. K. (2009). Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Re*search, 10:207–244.
- 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, pages 1307–1322.
- Xing, E. P., Jordan, M. I., Russell, S. J., and Ng, A. Y. (2003). Distance metric learning with application to clustering with side-information. In Advances in neural information processing systems.
- Ying, Y. and Li, P. (2012). Distance metric learning with eigenvalue optimization. *The Journal of Machine Learning Research*, 13:1–26.
- Ying, Y., Wen, L., and Lyu, S. (2016). Stochastic online auc maximization. In Advances in Neural Information Processing Systems.
- Ying, Y. and Zhou, D.-X. (2006). Online regularized classification algorithms. *IEEE Transactions on In*formation Theory, 52(11):4775–4788.
- Ying, Y. and Zhou, D.-X. (2016). Online pairwise learning algorithms. *Neural computation*, 28(4):743–777.
- Zhao, P., Jin, R., Yang, T., and Hoi, S. C. (2011). Online auc maximization. In *Proceedings of the 28th international conference on machine learning (ICML-11)*.