UNDERSTANDING IMPACTS OF DIFFERENTIAL PRIVACY: A UNIFIED FRAMEWORK WITH TWO-LAYER NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

Abstract

With the growing demand for data and the increasing awareness of privacy, differentially private learning has been widely applied in various deep models. Experiments have observed several side effects of differentially private learning, including bad learning features (performance), disparate impact, and worse adversarial robustness that hurt the trustworthiness of the trained models. Recent works have expected pre-training to mitigate these side effects. It is valuable to theoretically understand the impact of differential privacy on the training process. However, existing theoretical research only explained parts of the phenomena and failed to extend to non-convex and non-smooth neural networks. To fill this gap, we propose a unified framework to explain all the above phenomena by studying the feature learning process of differentially private stochastic gradient descent in two-layer ReLU convolutional neural networks. By analyzing the test loss, we find both its upper and lower bound decrease with feature-to-noise ratios (FNRs). We then show that disparate impact comes from imbalanced FNRs among different classes and subpopulation groups. Additionally, we show that the suboptimal learned features and reduced adversarial robustness are caused by the randomness of privacy-preserving noise introduced into the learned features. Moreover, we demonstrate that pre-training cannot always improve the model performance, especially with increased feature differences in the pre-training and fine-tuning datasets. Numerical results on both synthetic and real-world datasets validate our theoretical analyses.

032 033 034

035

043

044 045

046

047

048

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

1 INTRODUCTION

Modern deep learning models have achieved remarkable success in various applications, such as image classification (He et al., 2022) and natural language processing Vaswani et al. (2017). Many of these applications require training on datasets that contain sensitive private information. Differentially private learning, as an approach, seeks to train these models while ensuring rigorous privacy guarantees (Abadi et al., 2016). A standard differentially private learning algorithm is Differential Privacy Stochastic Gradient Descent (DP-SGD) (Abadi et al., 2016), which injects noise into network parameter updates during optimization.

While DP-SGD provides robust privacy guarantees, it introduces side effects, including

- *Bad learned features (Tramer & Boneh, 2020):* DP-SGD trained models might learn bad features that are worse than handcrafted ones.
- Disparate impact (Bagdasaryan et al., 2019; Sanyal et al., 2022): DP-SGD trained models achieve different accuracy on different classes and different subpopulation groups.
- *Worse adversarial robustness (Tursynbek et al., 2020):* DP-SGD trained models might be less adversarially robust than non-private models.
- 051 052

Moreover, recent work (De et al., 2022) has shown that pre-training improves the accuracy of DP-SGD trained models by 30% compared with training from scratch.

Understanding these phenomena is essential for deploying trustworthy, differentially private learning systems. Although several studies have explored side effects such as disparate impact (Esipova et al., 2022; Bagdasaryan et al., 2019), existing models are not applicable to ReLU neural networks due to their inherent non-convex and non-smooth nature. Moreover, a unified theoretical framework that explains all of the aforementioned phenomena is still lacking.

Thus, there remains an open problem:

How to theoretically explain the aforementioned phenomena in DP-SGD trained ReLU neural networks with a unified framework?

Facing the challenges in non-convex and non-smooth neural networks, we answer this problem in two steps. First, we derive test loss bounds of two-layer ReLU Convolutional Neural Network (CNN)s trained with DP-SGD in Section 3. Then, we extend the test loss analysis in Section 4 to explain the aforementioned phenomena.

069 1.1 RELATED WORK

060 061

062

063

068

099

102

103

107

In this section, we discuss the related works from two perspectives. Interested readers can refer to Appendix A for a detailed discussion.

073 Analysis on differentially private learning side effects. After observing the side effects of 074 differentially private learning, some researchers provided theoretical explanations. For example, Tran 075 et al. (2021) employed Taylor expansion to understand the disparate impact from the optimization 076 local loss landscape for twice differentiable loss functions. Sanyal et al. (2022) studied unfairness 077 in long-tailed data distribution in an asymptotic setting where the number of training data tends to infinity. Zhang & Bu (2022) studied the adversarial robustness of private linear classifiers. However, these analyses relied on some restricted assumptions that are not applicable to neural networks due to 079 their non-convex and non-smooth nature. In this paper, we provide a unified framework to explain all the mentioned side effects by characterizing the feature learning process of DP-SGD trained 081 two-layer ReLU CNNs.

Feature learning in neural networks. Feature learning in neural networks explores how neural networks learn data-related patterns during training, offering insights into phenomena such as momentum (Jelassi & Li, 2022), benign overfitting (Cao et al., 2022; Kou et al., 2023), adversarial training (Allen-Zhu & Li, 2022; Li & Li, 2023), and data augmentation Zou et al. (2023). In this work, we make an initial step towards studying the generalization performance and effects in DP-SGD trained two-layer neural networks.

Although the convergence and feature learning process are analyzed in two-layer neural networks before, their approaches are not applicable to DP-SGD trained neural networks as the key properties for feature growth may not hold due to the perturbation of added noise. In this paper, we design new techniques to study the generalization performance of DP-SGD trained models.

094 1.2 CONTRIBUTIONS AND MAIN RESULTS

We provide a unified framework to explain the side effects in DP-SGD trained two-layer ReLU CNNs.
 Specifically,

- We explore the feature learning process of DP-SGD trained neural networks. Considering the technical challenges posed by non-convex and non-smooth ReLU CNN and the random Differential Privacy (DP) noise, we propose a new proof technique to derive the standard and adversarial test loss bounds. The high-level idea is to use a piece-wise linear function to bound the non-linear loss function. We theoretically prove that the upper and lower bounds of test loss depend on data feature size and DP-SGD noise.
- We provide theoretical explanations for the side effects and the effectiveness of pre-training in DP-SGD trained neural networks.
 - The noise from DP-SGD highly perturbs the learned activation patterns during iterations, resulting in *bad learned features (performance)*.
 - 2

110 111

108

- 112
- 113 114

115

116 117

118

- Three factors, i.e., the gradient clipping, the group or class sizes, and feature sizes, contribute to *disparate impact* in DP-SGD trained models by influencing the test loss within a given group or class.
- The DP-SGD trained models exhibit *worse adversarial robustness* because they learn non-robust, class-irrelevant features from the random DP noise.
 - The fine-tuning performance on a pre-trained neural network decreases with increased feature difference. This implies that the benefits from pre-training may actually come from distribution overlaps between pre-training and fine-tuning datasets.
- 1.3 NOTATION

119 We use lowercase letters, lowercase boldface letters, and uppercase boldface letters to denote scalars, 120 vectors, and matrices, respectively. We use [m] to denote the set $\{1, \dots, m\}$. Given two sequences $\{x_n\}$ and $\{y_n\}$, we denote $x_n = \mathcal{O}(y_n)$ if $|x_n| \leq C_1 |y_n|$ for some positive constant C_1 and $x_n =$ 121 $\Omega(y_n)$ if $|x_n| \ge C_2 |y_n|$ for some positive constant C_2 . We use $x_n = \Theta(y_n)$ if both $x_n = \mathcal{O}(y_n)$ and 122 $x_n = \Omega(y_n)$ hold. We use $\mathcal{O}(\cdot), \hat{\Theta}(\cdot), \hat{\Omega}(\cdot)$ to omit logarithmic factors in these notations. Given a 123 set \mathcal{T} , we use $|\mathcal{T}|$ to denote its cardinality. For a vector $\mathbf{x} \in \mathbb{R}^d$, we denote its $\ell_p (p \ge 1)$ norm as 124 $\|\mathbf{x}\|_p = \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$. The notation $(\mathbf{x}, y) \sim \mathcal{D}$ indicates that the data sample (\mathbf{x}, y) is generated 125 126 from a distribution \mathcal{D} . 127

128

2 Model

129 130

134

140

141

142

143

144 145

146

147 148 149

We consider a one-hidden layer CNN with data structured as follows. Similar data structures have
been applied in (Allen-Zhu & Li, 2020; Jelassi & Li, 2022; Jelassi et al., 2022; Li & Li, 2023; Zou
et al., 2023; Cao et al., 2022).¹

Data distribution. We consider a 2-class classification problem over 2-patch inputs. Each labelled data is denoted as (\mathbf{x}, y) , with label $y \in \{1, 2\}$ and data vector $\mathbf{x} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \mathbb{R}^{d \times 2}$. We assume that each label class y has 2-group features, i.e., the majority (common) group features $\mathbf{u}_{y,\text{Maj}}$ and the minority (rare) group features $\mathbf{u}_{y,\text{Min}}$. A sample (\mathbf{x}, y) is generated from a data distribution \mathcal{D} as follows.

- 1. The label y is randomly sampled from $\{1, 2\}$. With probability $p_c > 0$, the label is selected as y = 1; otherwise, it is selected as y = 2.
- 2. Each input data patch $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \in \mathbb{R}^d$ contains either feature or noise.
 - Feature patch: One data patch (x⁽¹⁾ or x⁽²⁾) is randomly selected as the feature patch. With probability p_f > 0.5, this patch contains majority features u_{y,Maj} for y ∈ {1,2}. Otherwise, this patch contains minority features u_{y,Min} for y ∈ {1,2}.
 - Noisy patch: The remaining patch $\boldsymbol{\xi}$ is generated from a Gaussian distribution $\mathcal{N}(0, \sigma_p^2 \mathbf{H})$, where $\mathbf{H} = \mathbf{I} \sum_{i=1}^2 \sum_{j \in \{\text{Maj}, \text{Min}\}} \mathbf{u}_{i,j} \mathbf{u}_{i,j}^\top \cdot \|\mathbf{u}_{i,j}\|_2^{-2}$.

Without loss of generality, we assume all the feature vectors are orthogonal, i.e., $\langle \mathbf{u}_{i,j}, \mathbf{u}_{i',j'} \rangle = 0$ for all $i \in \{1,2\}, j \in \{\text{Maj}, \text{Min}\}$ when $(i, j) \neq (i', j')$. Additionally, we consider the features of majority and minority satisfy $p_f \|\mathbf{u}_{i,\text{Maj}}\|_2 > (1 - p_f) \|\mathbf{u}_{i,\text{Min}}\|_2$ for all $i \in \{1,2\}$.

Moreover, we define distributions $\mathcal{D}_{i,j}$, $i \in \{1, 2\}$, $j \in \{\text{Maj}, \text{Min}\}$, with probability density functions given by:

$$\mathbb{P}_{\mathcal{D}_{i,j}}[(\mathbf{x}, y)] = \mathbb{P}_{\mathcal{D}}[(\mathbf{x}, y)|y = i$$
, Feature patch of $\mathbf{x} = \mathbf{u}_{i,j}]$.

Learner model. We consider a two-layer CNN with ReLU activation as the learner model. The
 first layer of the CNN comprises m neurons (filters) for class 1 and m neurons for class 2. Each
 neuron processes the two data patches separately. The parameters of the CNN's second layer are

156

¹⁶¹

¹See (Allen-Zhu & Li, 2020) for detailed experimental justifications of the data distribution.

fixed as 1/m. Given an input vector $\mathbf{x} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$, the model outputs a vector $[F_1, F_2]$ whose k^{th} element is

$$F_k(\mathbf{W}, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \sum_{j=1}^2 \sigma\left(\left\langle \mathbf{w}_{k,r}, \mathbf{x}^{(j)} \right\rangle\right),\tag{1}$$

where $\sigma(\cdot) = \max\{\cdot, 0\}$ denotes the activation function ReLU(\cdot). We use W to denote the collection of all model weights and $\mathbf{w}_{k,r}$ to denote the weight vector the r^{th} neuron associated with $F_k(\mathbf{W}, \mathbf{x})$.

Training objective. Given a training dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ drawn from the distribution \mathcal{D} , we train the neural networks by minimizing the empirical risk with cross-entropy loss, i.e.,

$$\mathcal{L}_{\mathcal{S}} = \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}(\mathbf{W}, \mathbf{x}_i, y_i) = \frac{1}{n} \sum_{i=1}^{n} [-\log(\operatorname{prob}_{y_i}(\mathbf{W}, \mathbf{x}_i))],$$
(2)

where $prob(\cdot)$ represents the softmax predictions with the output of the neural network, i.e.,

$$\operatorname{prob}_{y_i}(\mathbf{W}, \mathbf{x}_i) = \frac{\exp(F_{y_i}(\mathbf{W}, \mathbf{x}_i))}{\sum_{k=1}^{2} \exp(F_k(\mathbf{W}, \mathbf{x}_i))}.$$
(3)

Differential privacy and training algorithm. DP (Dwork et al., 2014) (see the definition below) stands as the benchmark for quantifying privacy leakage, offering rigorous privacy guarantees.

Definition 2.1 ((α, δ) -Differential privacy). A randomized algorithm $\mathcal{M} : \mathcal{Z} \to \mathcal{R}$ is (α, δ) -DP if for every pair of neighboring datasets $Z, Z' \in \mathcal{Z}$ that differ in one entry and for any subset of output $\mathcal{S} \subseteq \mathcal{R}$, we have

$$\mathbb{P}\left[\mathcal{M}(Z)\in\mathcal{S}\right]\leq e^{\alpha}\mathbb{P}\left[\mathcal{M}(Z')\in\mathcal{S}\right]+\delta.$$
(4)

DP-SGD, the standard and most popular training algorithm (Abadi et al., 2016) (refer to Appendix C for details), achieves DP through an update given by

$$\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \frac{\eta}{B} \cdot \sum_{(\mathbf{x}, y) \in \mathcal{S}^{(t)}} \operatorname{clip}_{C} \left(\nabla \mathcal{L} \left(\mathbf{W}^{(t)}, \mathbf{x}, y \right) \right) + \eta \cdot \mathbf{n}^{(t)},$$
(5)

where η represents the learning rate and $\operatorname{clip}_C(\mathbf{x})$ is the gradient clipping function with clipping threshold *C* on vector \mathbf{x} , i.e.,

$$\operatorname{clip}_{C}(\mathbf{x}) = \frac{\mathbf{x}}{\max\left\{1, \left\|\mathbf{x}\right\|_{2}/C\right\}}.$$
(6)

In (5), $S^{(t)}$ represents the randomly selected mini-batch datasets with a batch size *B* from the training dataset *S* generated from the data distribution D, and $\mathbf{n}^{(t)}$ is the noise used for privacy protection following $\mathcal{N}(0, \sigma_n^2 \mathbf{I})$ at iteration *t*. We initialize network parameters with Gaussian initialization, where all entries of $\mathbf{W}^{(0)}$ are sampled from i.i.d. Gaussian distributions $\mathcal{N}(0, \sigma_0^2 \mathbf{I})$.

3 TEST LOSS ANALYSIS

In this section, we analyze the test loss of DP-SGD trained CNNs, serving as a stepping-stone for analyzing the impacts of DP-SGD in Section 4.

3.1 STANDARD TEST LOSS ANALYSIS

We first define the standard test loss of the trained model on data distribution $\mathcal{D}_{i,j}$, for any $i \in [2]$ and $j \in \{\text{maj, min}\}$ as follows.

$$\mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}) = \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[\mathcal{L}(\mathbf{W},\mathbf{x},y)\right].$$
(7)

Our results for test loss are based on the following conditions and assumption.

Condition 3.1. Suppose that there exists a positive constants $c_1 > 0, 0 < c_2 < 1$ and a sufficiently large constants $c_3, c_4 > 0$ such that

- 1. The CNNs have sufficiently large number of neurons, i.e., $m \ge c_1 \cdot d$.
- 2. The dimension d satisfies $d \ge 50$.
 - 3. The batch size satisfies $B \ge c_2 \cdot n$.

Feature size, patch noise and DP-SGD noise satisfies ||u_{i,j}||₂ = Θ(√dσ_p), σ_n ≤ c₃σ_p, for i ∈ [2], j ∈ {maj, min}.

5. The learning rate satisfies
$$\eta \leq \left(c_4(C + \sqrt{d}\sigma_n)(\max_{i,j} \|\mathbf{u}_{i,j}\|_2 + \sqrt{d}\sigma_p)\right)^{-1}$$
.

Condition 1 and 2 are common in theoretical analyses (e.g., (Allen-Zhu et al., 2019; Kou et al., 2023)) and also match the practical setting that modern neural networks are usually over-parameterize, i.e., have more parameters than the number of training examples. Condition 3 guarantees that batch size is proportional to the training dataset size so that stochastic gradients can take advantage of large training datasets. Condition 4 ensures that (1) the feature in not small so that the model can learn features (Kou et al., 2023); (2) the noise added by DP is in a reasonable range, i.e., privacy budgets are not too tight ($\epsilon \ge 0.5$). Condition 5 ensures that the model update is bounded.

Assumption 3.2 (*s*-non-perfect model). We assume the model is almost surely not perfect on a test example, i.e.,

$$\mathcal{L}\left(\mathbf{W}^{(t)}, \mathbf{x}, y\right) \ge s, \ \forall (\mathbf{x}, y) \sim \mathcal{D},$$
(8)

with some constant s > 0.

Assumption 3.2 is considered a mild assumption. Given that DP-SGD introduces randomness in the model training, the resulting trained model is randomized and is unlikely to attain a zero cross-entropy loss on a test example.

Next, we define Feature-to-Noise Ratios (FNRs) and the clipping factor to facilitate the analyses.
Definition 3.3 (Feature-to-noise ratios). Feature-to-noise ratio of the class *i* in group *j* is defined as

246 247

251

252 253

259

260 261 262

263

264 265 266

267

237 238

219 220

221

222 223

224

225 226

227

$$\mathcal{F}_{i,j} = \frac{\|\mathbf{u}_{i,j}\|_2}{\sqrt{d\sigma_n}},\tag{9}$$

where $\gamma_{i,j}$ denotes the expected proportion of class *i* group *j* data in the whole training dataset, i.e., $\gamma_{1,\text{maj}} = p_c p_f, \gamma_{1,\text{min}} = p_c (1 - p_f), \gamma_{2,\text{maj}} = (1 - p_c) p_f, \gamma_{2,\text{min}} = (1 - p_c) (1 - p_f).$ Definition 3.4 (Clipping factor). The clipping factor for the class *i* in group *i* is defined as

Definition 3.4 (Clipping factor). The clipping factor for the class i in group j is defined as

$$\Lambda_{i,j} = \frac{C}{\|\mathbf{u}_{i,j}\|_2 + \sigma_p \sqrt{d}}.$$
(10)

Here, the clipping factor $\Lambda_{i,j}$ quantifies the maximum change of gradient magnitude on data from class *i* and group *j*.

Based on the conditions, assumption, and definitions, we characterize upper bound test loss of DP-SGD trained models in Theorem 3.5.

Theorem 3.5. Under Condition 3.1 and Assumption 3.2, with a probability at least $1 - \exp(-\hat{\Omega}(d))$, for any $i \in \{1, 2\}, j \in \{maj, min\}$, the test loss of a DP-SGD trained model satisfies

$$\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(T)}\right) \leq \bar{L}_{i,j}(\mathbf{W}^{(T)}),\tag{11}$$

where

$$\bar{L}_{i,j}(\mathbf{W}^{(T)}) = \underbrace{\exp\left(-\Omega\left(\frac{\Lambda_{i,j}\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{\sqrt{m}}T\right)\right) \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right)}_{\text{Vanishing error}} + \underbrace{\mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{mn}\gamma_{i,j}\Lambda_{i,j}}\right)}_{\text{Generalization error}}$$
(12)

268
269 +
$$\mathcal{O}\left(\frac{\sqrt{m}}{\Lambda_{i,j}\gamma_{i,j}\mathcal{F}_{i,j}}\right)$$

Privacy protection error

274

275

281

284 285

286

287

288 289

290

291

292

293

308

313

319

322

We defer the proof of Theorem 3.5 to Appendix G.2. As the randomness of noise may improve the model, the high probability test loss lower bound becomes trivial. We instead provide a lower bound of the expected test loss in the following. Theorem 3.6 Under Condition 3.1 and Assumption 3.2 with the number of iterations T

Theorem 3.6. Under Condition 3.1 and Assumption 3.2, with the number of iterations $T \ge \Omega\left(-1/\log\left(1-\Omega\left(\eta\min_{i,j}\{\gamma_{i,j} \|\mathbf{u}_{i,j}\|_2^2\}/m\right)\right)\right)$ and a probability at least $1-\tilde{\mathcal{O}}(1/d)$, for any $i \in \{1,2\}, j \in \{maj, min\}$, the expected test loss a DP-SGD trained model satisfies

$$\mathbb{E}\Big[\mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(T)})\Big] \ge \exp\left(-\Omega\left(\frac{\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}T\right)\right) \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(0)}) + \underbrace{\Omega\left(\frac{\sigma_{p}^{2}}{\gamma_{i,j}\mathcal{F}_{i,j}^{2}}\right)}_{\text{Privacy protection error}} - \mathcal{O}\left(\frac{1}{\gamma_{i,j}}\sqrt{\frac{d}{n}}\right).$$
(13)

Remark 3.7. Theorems 3.5 and 3.6 show that the upper and lower bound of test loss both decrease with the feature-to-noise ratio $\mathcal{F}_{i,j}$. Moreover, Theorem 3.5 illustrates the presence of non-vanishing error terms in the test loss bound, i.e., generalization error and privacy protection error. Specifically,

- Generalization error arises from data noise patch. This error decreases with $O(1/\sqrt{n})$, aligning with the generalization error bounds derived in neural networks (Arora et al., 2019; Xu & Mannor, 2012).
- Privacy protection error stems from DP noise. With a fixed privacy budget ((α, δ) in Definition 2.1), the noise variance σ_n^2 increases with the number of iterations T as $\sigma_n^2 = \Theta(T)$ (Abadi et al., 2016). Consequently, the error due to privacy protection increases with the number of iterations as $\Theta(\sqrt{T})$. This rate aligns with the non-vanishing training error of DP-SGD on Lipschitz-smooth objectives (Bu et al., 2024).

Privacy-Utility Tradeoff Privacy composi-295 tion (Theorem 1 in (Abadi et al., 2016)) shows 296 that the noise variance σ_n^2 increases with $\Theta(T)$ 297 under (α, δ) -DP. Then, the privacy protection 298 error in (12) is in order $\mathcal{O}(\sqrt{T}\log(1/\delta)/(n\alpha))$, 299 implying that 1) larger privacy budget leads to 300 larger loss; 2) larger dataset has small error 301 thanks to privacy amplification by subsampling; 302 3) the loss increases with the number of itera-303 tions. In addition, we demonstrate a sharp phase 304 transition between benign and harmful privacy 305 protection to model utility, as shown in Figure 1. Details about the simulation settings are pro-306 vided in Appendix D.2. 307

3.2 Adversarial Test Loss Analysis



Figure 1: Illustration of the privacy-utility phase transition between benign and harmful privacy protection. The yellow region represents a benign regime where the test loss is small. The blue region represents a harmful regime where the test loss is large.

In this subsection, we analyze the impact of
 DP-SGD on adversarial robustness as a basis for
 understanding the side offset of means adversarial

understanding the side effect of worse adversarial robustness in Section 4.

Adversarial robustness refers to a machine learning model's ability to withstand carefully manipulated
 input samples, commonly known as adversarial examples.

Definition 3.8 (Adversarial example). For an data example (\mathbf{x}, y) , the corresponding adversarial example is $\tilde{\mathbf{x}} = \mathbf{x} + \boldsymbol{\zeta}$, where $\boldsymbol{\zeta} = \arg \max_{\|\boldsymbol{\zeta}\|_p \leq \bar{\boldsymbol{\zeta}}} \mathcal{L}(\mathbf{W}, \mathbf{x} + \boldsymbol{\zeta}, y)$ is the adversarial perturbation and $\bar{\boldsymbol{\zeta}}$ is the perturbation radius.

We evaluate the performance of the trained model on adversarial robustness using adversarial test loss, which is defined as

$$\mathcal{L}_{\mathcal{D}}^{\mathrm{adv}}(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\boldsymbol{\zeta}\|_{p} \leq \bar{\boldsymbol{\zeta}}} \mathcal{L}(\mathbf{W}, \mathbf{x} + \boldsymbol{\zeta}, y) \right].$$
(14)

We give our main result on adversarial test loss in the following theorem.

Theorem 3.9. Under Condition 3.1 and Assumption 3.2, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, for any $i \in \{1, 2\}, j \in \{maj, min\}$, the adversarial test loss of a DP-SGD trained model with learning rate $\eta = \Theta(1)$ satisfies

$$\mathcal{L}_{\mathcal{D}_{i,j}}^{\mathrm{adv}}\left(\mathbf{W}^{(T)}\right) \leq \bar{L}_{i,j}(\mathbf{W}^{(T)}) + \underbrace{\mathcal{O}\left(\left[\frac{T}{m}C + \frac{\sqrt{Td}}{m}\sigma_n + \sqrt{d}\sigma_0\right]\bar{\zeta}d^{1-1/p}\right)}_{\mathrm{By\ adversarial\ perturbation}}.$$
(15)

We defer the proof of Theorem 3.9 to Appendix G.4.

Remark 3.10. Theorem 3.9 shows that the error induced by adversarial perturbation increases at a rate of $\mathcal{O}(T)$ and contains a DP noise term $\sqrt{Td}/m\sigma_n$, which increases at a rate of $\mathcal{O}(\sqrt{T})$. Moreover, the test loss bound aligns with the excess risk bound of adversarial training (Xiao et al., 2022).

4 UNDERSTANDING DP-SGD IMPACTS

In Sections 3, we establish the test loss bounds (Theorems 3.5 and 3.9) for a DP-SGD trained model. In this section, we extend the results to interpret the side effects and the effectiveness of pre-training. We extend the upper bound in Theorem 3.5, and we can also get a similar lower bound by extending the lower bound in Theorem 3.6 to obtain similar insights.

4.1 INTERPRETATION OF BAD FEATURES (PERFORMANCE)

Theorem 3.5 implies that the learned features are perturbed by the noise introduced during DP-SGD training. This noise prevents CNNs from learning perfect features, ultimately resulting in a non-vanishing error. Be specific, this loss increases when the feature is seriously affected by the DP noise, i.e., a smaller FNR^{DP}_{*i,j*} (see Definition 3.3).

4.2 INTERPRETATION OF DISPARATE IMPACT

Define the distribution of class i as $\mathbb{P}_{\mathcal{D}_i}[(\mathbf{x}, y)] = \mathbb{P}_{\mathcal{D}}[(\mathbf{x}, y)|y = i]$, $i \in \{1, 2\}$. We evaluate model performance on different classes by bounding the test loss on data distribution within a class.

Corollary 4.1 (Disparate impact of different classes). Under Condition 3.1 and Assumption 3.2, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, for any $i \in \{1, 2\}$, the test loss of a DP-SGD trained model with learning rate $\eta = \Theta(1)$ satisfies

$$\mathcal{L}_{\mathcal{D}_{i}}\left(\mathbf{W}^{(T)}\right) \leq \frac{1}{\sum_{j \in \{\min, \max \}} \gamma_{i,j}} \sum_{j \in \{\min, \max \}} \gamma_{i,j} \cdot \bar{L}_{i,j}\left(\mathbf{W}^{(T)}\right).$$
(16)

Similarly, define the distribution of group j as D_j , $j \in \{\text{maj}, \text{min}\}$. We evaluate model performance across different groups (majority and minority) by the test loss on data distribution within a group.

Corollary 4.2 (Disparate impact of subpopulation groups). Under Condition 3.1 and Assumption 3.2, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, for any $j \in \{\text{maj}, \min\}$, the test loss of a DP-SGD trained model with learning rate $\eta = \Theta(1)$ satisfies

$$\mathcal{L}_{\mathcal{D}_{j}}\left(\mathbf{W}^{(T)}\right) \leq \frac{1}{\sum_{i=1}^{2} \gamma_{i,j}} \sum_{i=1}^{2} \gamma_{i,j} \cdot \bar{L}_{i,j}\left(\mathbf{W}^{(T)}\right).$$
(17)

Recall the expression of *L* in (12). Corollaries 4.1 uncovers three sources of disparate impact: *gradient clipping* $\Lambda_{i,j}$, *data imbalance* $\gamma_{i,j}$, and *feature disparity* $\|\mathbf{u}_{i,j}\|_2$, which correspond to three key components of FNRs and clipping factors (see Definition 3.3 and 3.4). We discuss them in detail below. **Gradient clipping.** Theorem 3.5 shows that the test loss depends on the clipping factor $\Lambda_{i,j}$. This implies that a class or group with a larger gradient norm will experience more aggressive clipping, leading to poorer feature learning performance.

Data imbalance. As the privacy protection error terms in Theorems 3.5 and 3.6 decrease with the data proportion $\gamma_{i,j}$, a group or a class with more data enjoys better performance and vice versa. This raises a risk of worse model performance on the skewed data sources² and the long-tailed distributed applications (Feldman & Zhang, 2020).³

Feature disparity. The feature-to-noise ratio $\mathcal{F}_{i,j}$ (see Definition 3.3) depends on the feature sizes of the data $\mathbf{u}_{i,j}$. In real-world applications, data from different classes or groups may have significantly different feature sizes, resulting in divergent model performances.

Remark 4.3. Recent papers, e.g., (Esipova et al., 2022) have pointed out that the disparate impact
 mainly arises from gradient clipping. However, we show that even without gradient clipping, different
 groups still exhibit different performances due to data imbalance and the intrinsic differences in
 feature sizes.

4.3 INTERPRETATION OF ADVERSARIAL ROBUSTNESS

As shown in Theorem 3.9, DP-SGD leads to a high adversarial test loss. We interpret worse adversarial robustness from two perspectives.

Feature learning. As pointed out in Allen-Zhu & Li (2020), an adversarially robust model typically removes the non-robust class-irrelevant noises and learns robust features. However, DP-SGD injects much noise during the training process, making the neural networks inevitably learn non-robust class-irrelevant noises.

Network parameter growth. Due to the noise added by DP-SGD, the network parameters have a consistently growing norm with the number of iterations. As adversarial perturbation ζ attacks the model by changing the neurons' activated inner products, i.e.,

$$F_{k}(\mathbf{W}, \mathbf{x} + \boldsymbol{\zeta}) = \frac{1}{m} \sum_{r=1}^{m} \sum_{j=1}^{2} \left[\sigma \left(\left\langle \mathbf{w}_{k,r}, \mathbf{x}^{(j)} \right\rangle + \underbrace{\left\langle \mathbf{w}_{k,r}, \boldsymbol{\zeta}^{(j)} \right\rangle}_{\text{By adversarial perturbation}} \right) \right], \quad (18)$$

higher network parameter norms result in increased vulnerability to adversarial attacks.

4.4 PRE-TRAINING AND FINE-TUNING

416 As demonstrated in Berrada et al. (2023), pre-training on public datasets can significantly reduce the accuracy drop and mitigate the side effects caused by DP-SGD. With pre-training, the neural network 417 utilizes good pre-trained features and the loss at initialization $\mathcal{L}_{\mathcal{D}_{i,i}}(\mathbf{W}^{(0)})$ is relatively small. As 418 a result, fine-tuning only needs a small number of iterations, leading to smaller privacy protection 419 errors and thus, mild side effects. In Tramèr et al. (2022), the authors pointed out that many of 420 the existing pre-training and fine-tuning datasets satisfy the private fine-tuning data distribution is 421 essentially a subset of the public data distribution (e.g. pre-training on ImageNet and fine-tuning on 422 CIFAR-10). So, in this subsection, we explore how the distribution shifts affect the private fine-tuning 423 performance. 424

425 Consider a two-layer CNN model (defined in Section 2) that is trained on a simplified data distribution 426 \mathcal{D}_{pt} that each class only has one feature with the same magnitude, i.e., \mathbf{u}_1 for class 1, and \mathbf{u}_2 for class 2

427

386

394 395

399

404

405

406

413 414

 ²Take ImageNet (Deng et al., 2009) as an example. More than 45% of ImageNet data comes from the United
 States, corresponding to only 4% of the world's population; In contrast, China and India contribute just 3% of
 ImageNet data (Zou & Schiebinger, 2018). Thus, DP-SGD trained models on ImageNet may perform poorly on
 tasks concerning China and India.

^{431 &}lt;sup>3</sup>For example, in the SUN dataset (Xiao et al., 2010), the number of examples in each class displays a long-tailed structure (Feldman, 2020).

and $\|\mathbf{u}_1\|_2 = \|\mathbf{u}_2\|_2$, with SGD and fine-tuned on similar data distribution $\mathcal{D}_{\rm ft}$ (data is first generated from $\mathcal{D}_{\rm pt}$ and then rotated with an angle θ) (Details are discussed in Appendix B). Based on the results in Kou et al. (2023), a ReLU CNN trained with gradient descent learns the feature in a constant order. Therefore, we assume that the pre-trained trained model is $\tilde{\mathbf{w}}_{j,r} = C_1 \cdot \mathbf{u}_j + C_3 \cdot \mathcal{N}(0, \sigma_p \mathbf{H})$ for simplicity. Then we can characterize the fine-tuning test loss as follows (We define Λ_i in a similar manner that $\Lambda_i = C/\|\mathbf{u}_i\|_2 + \sigma_p \sqrt{d}$).

Proposition 4.4. Suppose, then the fine-tuned test loss satisfies

$$\mathcal{L}_{\mathcal{D}_{ft}}(\mathbf{W}^{(T)}) \le \exp\left(-\Omega\left(\frac{\Lambda_i \|\mathbf{u}_i\|_2^2}{\sqrt{m}}T\right)\right) \cdot \tilde{L} + \mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{mn}\Lambda_i}\right) + \mathcal{O}\left(\frac{\sqrt{md}\sigma_n}{\Lambda_i \|\mathbf{u}_i\|_2}\right),\tag{19}$$

where

438

452

453

454

455 456 457

458 459

460

461 462 463

464

$$\tilde{L} = -\frac{1}{2} \ln \left(\frac{\exp(C_1 \cos \theta \|\mathbf{u}_1\|_2^2)}{\exp(C_1 \cos \theta \|\mathbf{u}_1\|_2^2) + \exp(C_1 \sin \theta \|\mathbf{u}_1\|_2^2 + C_3 \sigma_p^2)} \right) -\frac{1}{2} \ln \left(\frac{\exp(C_1 \cos \theta \|\mathbf{u}_2\|_2^2)}{\exp(C_1 \cos \theta \|\mathbf{u}_2\|_2^2) + \exp(C_3 \sigma_p^2)} \right).$$
(20)

Remark 4.5. A worth noting fact is that L increases with θ , meaning that the fine-tuning test loss decreases with increased feature difference, i.e., θ . Even worse, when $\tilde{L} > \mathcal{L}_{\mathcal{D}_2}(\mathbf{W}^{(0)})$, the pre-training can lead to worse performance than training from scratch. Hence, as indicated in Tramèr et al. (2022), the remarkable good performance of pre-training on private learning may come from the distribution overlap between pre-training and fine-tuning datasets.

5 EXPERIMENTS

In this section, we use synthetic data and real-world datasets MNIST (LeCun et al., 1998) and CIFAR-10 (Krizhevsky et al., 2009) to verify our theory.

5.1 SYNTHETIC DATASETS

Synthetic data generation. We generate a synthetic dataset following the data distribution described in Section 2. Specifically, we set the sizes of both training and test datasets to 450. We set the feature vector length in each patch to 100. The feature vector sizes and dataset sizes of the majority and minority groups of classes 0 and 1 are specified as $\|\mathbf{u}_{1,\text{maj}}\|_2 = 4$, $\gamma_{1,\text{maj}} = 44\%$, $\|\mathbf{u}_{1,\text{min}}\|_2 =$ 2, $\gamma_{1,\text{min}} = 22\%$, $\|\mathbf{u}_{2,\text{maj}}\|_2 = 1.5$, $\gamma_{2,\text{maj}} = 22\%$ and $\|\mathbf{u}_{2,\text{min}}\|_2 = 0.5$, $\gamma_{2,\text{min}} = 11\%$. The difference in feature vector sizes and dataset sizes characterize feature disparity and data imbalance. In addition, we set the standard deviation of the noise patch to $\sigma_p = 0.2$.

We train a two-layer CNN with ReLU activation function and cross-entropy loss (see Section 2). We set the number of neurons as 64, i.e., m = 32. We utilize the default initialization in PyTorch. We train the CNN with DP-SGD with batch size B = 128 for 20 epochs.

Test loss. We experiment on the model's test loss of groups with different features in Figure 2 (a).
We observe that the test loss generally increases with the DP noise standard deviation, aligning with the findings in Theorem 3.5 where the upper bound of test loss depends on the corresponding FNR.
Furthermore, the class and the group with a larger feature size incur smaller test losses. It is also worth noting that the gaps among the groups become more significant with increasing noise standard deviation.

Adversarial robustness. In Figure 2(b), we assess the adversarial robustness by attacking the model with the projected gradient descent method (generate the adversarial examples by maximizing the loss with projected gradient descent) with $\bar{\zeta} = 0.02$. We observe that the DP-SGD trained model experiences a degradation in adversarial robustness on certain groups/classes. This aligns with the results from Theorem 3.5 that the upper bound of the model's adversarial test loss increases with DP noise standard deviation.



(a) Test loss versus DP noise standard devia-(b) Adversarial test loss attacked with the tion σ_n . Projected Gradient Descent Method.

Figure 2: Model standard test loss and adversarial test loss.

5.2 REAL-WORLD DATASETS

Setup. For MNIST and CIFAR-10, we train LeNet and a CNN following the architecture in (Tramer & Boneh, 2020) with DP-SGD. We fix the privacy budget as $\alpha = 3, \delta = 10^{-5}$. We fix the gradient clipping threshold as C = 0.1 (Tramer & Boneh, 2020). We fix the batch size as 256 and try various learning rates. We use the DP-SGD implementation in Opacus (Yousefpour et al., 2021). We obtain the adversarial example through the projected gradient descent method with $\bar{\zeta} = 4/255$.

Impact of feature size. In real images, it is hard to distinguish which part of the images represents the features. We simply see the object as features and regard the background as noise. To emulate the image backgrounds, we apply zero-padding to the periphery of the input images and the resize the image back to the original size (as shown in Figure 3 in the appendix). The table below shows both the test accuracy and adversarial test accuracy with different padding ratios. As indicated in the theory, the model accuracy decreases with the padding ratios (feature sizes).

Padding ratio	0%	13%	26%	38%	50%	62%	75%
MNIST	97%	97%	97%	96%	95%	94%	86%
CIFAR-10	58%	56%	54%	52%	51%	48%	46%
MNIST (adversari	al) 95%	93%	89%	77%	50%	20%	1%
CIFAR-10 (advers	arial) 3%	3%	2%	2%	1%	0%	0%

Pre-training vs. fine-tuning. We maintain the training dataset unrotated. We rotate the test images and split them into fine-tuning training dataset and test dataset. We additionally use ResNet-18 (He et al., 2016) for CIFAR-10. The table below shows the fine-tuning test performance under different rotation angles. As indicated in the theory, the model accuracy decreases with the rotation angle.

Rotation angle	0°	22.5°	45°	67.5°	90°
MNIST	99%	97%	94%	94%	94%
CIFAR-10 (CNN Tramer & Boneh (2020))	70%	59%	51%	49%	51%
CIFAR-10 (ResNet-18)	91%	66%	43%	37%	44%

In this paper, we investigate the side effects of DP-SGD in two-layer ReLU CNNs. We show that the side effects of DP-SGD trained CNNs depend on data's feature, data noise, and privacy-preserving noise. Our results uncover three sources of disparate impact: *gradient clipping, data imbalance*, and *feature disparity*. In addition, we show that the privacy-preserving noise introduces randomness into the learned features, leading to bad learned features and worse adversarial robustness. Moreover, we demonstrate that the fine-tuning performance with pre-trained models decreases with increased feature differences in the pre-training and fine-tuning datasets. Numerical results on both synthetic and real-world datasets validate our theoretical analyses.

CONCLUSIONS AND FUTURE WORKS

5407REPRODUCIBILITY STATEMENT541

We put the source code in the Supplementary Material. All the experimental results can be reproduced with the source code.

References

542

543

544

546

550

586

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and
 self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In 2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS), pp. 977–988. IEEE, 2022.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over parameterization. In *International conference on machine learning*, 2019.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, 2019.
- Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. Differential privacy has disparate
 impact on model accuracy. *Advances in neural information processing systems*, 32:15479–15488,
 2019.
- Leonard Berrada, Soham De, Judy Hanwen Shen, Jamie Hayes, Robert Stanforth, David Stutz,
 Pushmeet Kohli, Samuel L Smith, and Borja Balle. Unlocking accuracy and fairness in differentially
 private image classification. *arXiv preprint arXiv:2308.10888*, 2023.
- Zhiqi Bu, Hua Wang, Zongyu Dai, and Qi Long. On the convergence and calibration of deep learning
 with differential privacy. *Transactions on machine learning research*, 2023, 2023.
- Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. Automatic clipping: Differentially private deep learning made easier and stronger. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yuan Cao, Zixiang Chen, Misha Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. *Advances in neural information processing systems*, 35:25237–25250, 2022.
- Rachel Cummings, Varun Gupta, Dhamma Kimpara, and Jamie Morgenstern. On the compatibility of privacy and fairness. In *Adjunct publication of the 27th conference on user modeling, adaptation and personalization*, pp. 309–315, 2019.
- Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlocking high accuracy differentially private image classification through scale. *arXiv preprint arXiv:2204.13650*, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- 593 Maria S Esipova, Atiyeh Ashari Ghomi, Yaqiao Luo, and Jesse C Cresswell. Disparate impact in differential privacy from gradient misalignment. *arXiv preprint arXiv:2206.07737*, 2022.

- Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings* of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, pp. 954–959, 2020.
- 597 Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long
 598 tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891,
 599 2020.
- Georgi Ganev, Bristena Oprisanu, and Emiliano De Cristofaro. Robin hood and matthew effects:
 Differential privacy has disparate impact on synthetic data. In *International Conference on Machine Learning*, pp. 6944–6959. PMLR, 2022.
- Antonious Girgis, Deepesh Data, Suhas Diggavi, Peter Kairouz, and Ananda Theertha Suresh.
 Shuffled model of differential privacy in federated learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 2521–2529. PMLR, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 pp. 770–778, 2016.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in
 deep learning. In *International Conference on Machine Learning*, pp. 9965–10040. PMLR, 2022.
- Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure.
 Advances in Neural Information Processing Systems, 35:37822–37836, 2022.
- Yiwen Kou, Zixiang Chen, Yuanzhou Chen, and Quanquan Gu. Benign overfitting in two-layer relu
 convolutional neural networks. In *International Conference on Machine Learning*, pp. 17615–17659. PMLR, 2023.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Binghui Li and Yuanzhi Li. Why clean generalization and robust overfitting both happen in adversarial training. *arXiv preprint arXiv:2306.01271*, 2023.
- Paul Mangold, Michaël Perrot, Aurélien Bellet, and Marc Tommasi. Differential privacy has bounded
 impact on fairness in classification. In *International Conference on Machine Learning*, 2023.
- Lucas Rosenblatt, Julia Stoyanovich, and Christopher Musco. A simple and practical method for reducing the disparate impact of differential privacy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- Amartya Sanyal, Yaxi Hu, and Fanny Yang. How unfair is private learning? In *Uncertainty in Artificial Intelligence*, pp. 1738–1748. PMLR, 2022.
- Florian Tramer and Dan Boneh. Differentially private learning needs better features (or much more data). *arXiv preprint arXiv:2011.11660*, 2020.
- Florian Tramèr, Gautam Kamath, and Nicholas Carlini. Position: Considerations for differentially
 private learning with large-scale public pretraining. In *Forty-first International Conference on Machine Learning*, 2022.
- Cuong Tran, My Dinh, and Ferdinando Fioretto. Differentially private empirical risk minimization under the fairness lens. *Advances in Neural Information Processing Systems*, 34:27555–27565, 2021.
- 647 Nurislam Tursynbek, Aleksandr Petiushko, and Ivan Oseledets. Robustness threats of differential privacy. *arXiv preprint arXiv:2012.07828*, 2020.

648 649 650	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30:6000–6010, 2017.
651 652 653 654	Jiancong Xiao, Yanbo Fan, Ruoyu Sun, Jue Wang, and Zhi-Quan Luo. Stability analysis and generalization bounds of adversarial training. <i>Advances in Neural Information Processing Systems</i> , 35:15446–15459, 2022.
655 656 657	Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In <i>2010 IEEE computer society conference on</i> <i>computer vision and pattern recognition</i> , pp. 3485–3492. IEEE, 2010.
658 659	Huan Xu and Shie Mannor. Robustness and generalization. Machine learning, 86:391–423, 2012.
660 661 662	Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad, Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, et al. Opacus: User-friendly differential privacy library in pytorch. <i>arXiv preprint arXiv:2109.12298</i> , 2021.
663 664 665	Hanlin Zhang, Xuechen Li, Prithviraj Sen, Salim Roukos, and Tatsunori Hashimoto. A closer look at the calibration of differentially private learners. <i>arXiv preprint arXiv:2210.08248</i> , 2022.
666 667	Yuan Zhang and Zhiqi Bu. Differentially private optimizers can learn adversarially robust models. <i>arXiv preprint arXiv:2211.08942</i> , 2022.
668 669 670	Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. The benefits of mixup for feature learning. <i>arXiv preprint arXiv:2303.08433</i> , 2023.
671 672	James Zou and Londa Schiebinger. Design ai so that it's fair. Nature, 559(7714):324–326, 2018.
673	
674	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
6097	
600	
700	
701	

702 ADDITIONAL RELATED WORK А

703

711

704 **Side effects in differentially private learning.** Side effects of DP have been widely studied in deep 705 learning literature. Disparate impact was initially observed in classification tasks Bagdasaryan et al. 706 (2019) and generative tasks (Ganev et al., 2022). Then, researchers have proposed several methods, 707 such as a regularization approach (Tran et al., 2021), re-weighting, and stratification methods (Esipova 708 et al., 2022; Rosenblatt et al., 2024) to mitigate the disparate impacts. A recent paper (Berrada et al., 709 2023) showed that the DP models pre-trained with large datasets and fine-tuned with large batch 710 size can have marginal disparate effects. We discuss the implication of this work to pre-training in Appendix ??.

712 The interplay between DP and *fairness* is also a widely studied topic. Cummings et al. (2019) showed 713 that exact fairness is not compatible with DP under the PAC learning setting. Sanyal et al. (2022) 714 showed that it is not possible to build accurate learning algorithms that are both private and fair when 715 data follows a specific kind of long-tailed distribution. Mangold et al. (2023) bounded the difference in fairness levels between private and non-private models under the assumption that the confidence 716 margin is Lipschitz-continuous. 717

718 Some other side effects have also been studied. Tramer & Boneh (2020) showed that Differentially 719 Private Learning (DPL) may perform worse after bad feature learning compared with learning 720 handcraft features. Tursynbek et al. (2020) studied adversarial robustness in DPL and showed 721 that models trained by DPL may be more vulnerable compared with non-private models. Zhang & Bu (2022) studied the adversarial robustness of private linear classifiers and showed differentially 722 privately fine-tuned pre-trained models may be robust under certain parameter settings. In addition, 723 Zhang et al. (2022) studied calibration of DPL and observed miscalibration across a wide range 724 of vision and language tasks. Bu et al. (2023) studied DPL with neural tangent kernel (NTK) and 725 demonstrated that a large clipping threshold may benefit the calibration of DPL. 726

In these works, researchers have studied the convergence of DPL and explored explanations for the 727 aforementioned side effects. However, these analyses relied on some restricted assumptions that are 728 not applicable to neural networks because (1) differentially private neural networks training is not in 729 the NTK regime as the noise keeps the network parameter far away from the initialization during 730 training; (2) training loss of ReLU neural networks is non-convex and non-smooth, contradicting the 731 assumptions in most analyses. In this work, we aim to overcome these challenges and explain the 732 side effects of DP-SGD on a two-layer ReLU CNN. 733

734 735

736

740

744

745

746

747 748

749

750

751 752

В DETAILS ABOUT PRE-TRAINING AND FINE-TUNING DATA DISTRIBUTIONS

737 To illustrate the impact of pre-training, we simplify the data distribution as follows. We consider the 738 following pre-training and fine-tuning data distributions. We control their difference by a parameter θ. 739

Pre-training data distribution. We consider a 2-class classification problem over 2-patch inputs. 741 Each labelled data is denoted as (\mathbf{x}, y) , with label $y \in \{1, 2\}$ and data vector $\mathbf{x} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{1, 2\}$ 742 $\mathbb{R}^{d \times 2}$. A sample (\mathbf{x}, y) is generated from a data distribution \mathcal{D} as follows. 743

- 1. The label y is randomly sampled from $\{1, 2\}$. With probability 1/2, the label is selected as y = 1; otherwise, it is selected as y = 2.
- 2. Each input data patch $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \in \mathbb{R}^d$ contains either feature or noise.
 - Feature patch: One data patch ($\mathbf{x}^{(1)}$ or $\mathbf{x}^{(2)}$) is randomly selected as the feature patch. This patch contains a feature \mathbf{u}_y for $y \in \{1, 2\}$.
 - Noisy patch: The remaining patch $\boldsymbol{\xi}$ is generated from a Gaussian distribution $\mathcal{N}(0, \sigma_p^2 \mathbf{H})$, where $\mathbf{H} = \mathbf{I} \sum_{i=1}^{2} \mathbf{u}_i \mathbf{u}_i^{\top} \cdot \|\mathbf{u}_i\|_2^{-2}$.

753 **Fine-tuning data distribution.** We consider a 2-class classification problem over 2-patch inputs. 754 Each labelled data is denoted as (\mathbf{x}, y) , with label $y \in \{1, 2\}$ and data vector $\mathbf{x} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \in \{1, 2\}$ 755 $\mathbb{R}^{d \times 2}$. A sample (\mathbf{x}, y) is generated from a data distribution \mathcal{D} as follows.



810 E SOME DATASETS AND NEURAL NETWORK ARCHITECTURES OBSERVING 811 SIDE EFFECTS

Bagdasaryan et al. (2019) has identified disparate effects in datasets such as MNIST, CelebA, and Twitter posts using ResNet and LSTM networks. Unfairness has been observed in CelebA and CIFAR-10 datasets using ResNet networks (Sanyal et al., 2022). Zhang et al. (2022) observed miscalibration in QNLI, QQP, SST-2 with fine-tuned RoBERTa-base.

F OVERVIEW OF CHALLENGES AND PROOF SKETCH

In this section, we outline the main challenges in studying feature learning of DP-SGD on CNNs and the key proof techniques employed to overcome the challenges.

F.1 CHALLENGE 1: NON-SMOOTHNESS OF THE RELU ACTIVATION FUNCTION

The first challenge arises from the non-smoothness of the ReLU activation function. Some existing papers (e.g., (Girgis et al., 2021)) analyzed DP-SGD with Lipschitz-smoothness-based approaches. However, this kind of approach is not applicable to ReLU neural networks.

In a two-layer CNNs, we can track the neurons' feature learning process through their gradients. For any $i \in \{1, 2\}, r \in [m]$, the gradient on $\mathbf{w}_{i,r}^{(t)}$ can be decomposed as

833

835

813

814

815

816

817 818

819 820

821

822 823

824 825

826

827

828

829

$$\nabla_{\mathbf{w}_{i,r}^{(t)}} \mathcal{L}_{\mathcal{S}}(\mathbf{W}^{(t)}) = \underbrace{\sum_{j \in \{\text{maj,min}\}} (\mu_{i,j} \mathbf{u}_{i,j} - \mu_{3-i,j} \mathbf{u}_{3-i,j})}_{\text{Data features}} + \underbrace{\sum_{k=1}^{n} \boldsymbol{\xi}_{k}(\bar{\rho}_{i,j,k} \mathbb{I}(y_{k}=i) - \underline{\rho}_{i,j,k} \mathbb{I}(y_{k}\neq i))}_{\text{Data noise}}$$

where $\mu_{i,j}, \mu_{3-i,j}, \bar{\rho}_{i,j,k}, \underline{\rho}_{i,j,k} \ge 0$ are constants. The neurons tend to learn both class-relevant data features and data noise of the targeted class while unlearning others.

Some existing approaches (e.g., (Cao et al., 2022)) bound the feature learning process by characterizing the leading neurons that learn the most features. However, this approach only works for ReLU^q (q > 2) activation functions, where the leading neurons dominate other neurons during training. As the ReLU function is piece-wise linear, these approaches fail in ReLU neural networks.

To overcome this challenge, we study the feature learning process by analyzing the dynamics of the model outputs $F_i(\mathbf{W}, \mathbf{x}), i \in [2]$ defined in Section 2

845 846

847

857

F.2 CHALLENGE 2: RANDOMNESS FROM DP-SGD

The second challenge stems from the randomness introduced by DP-SGD. The learning process is significantly perturbed due to the random noise in DP-SGD. Kou et al. (2023) attempted to bound the feature learning process based on the monotonicity of the weights of feature vectors. However, due to the randomness from DP-SGD, the weights are not consistently increasing.

To address this challenge, we track the increments of the model outputs for any data point $(\mathbf{x}, y) \sim \mathcal{D}_{i,j}$ instead. The key proposition is presented as follows.

Proposition F.1. For any $(\mathbf{x}, y) \sim \mathcal{D}_{i,j}, i \in [2], j \in \{maj, min\}$, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, exp $(-\tilde{\Omega}(d))$,

• The increment of model output for the targeted class y satisfies

$$\Delta_{y}^{(t)}(\mathbf{x}) = F_{y}\left(\mathbf{W}^{(t+1)}, \mathbf{x}\right) - F_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)$$

$$\geq \Omega\left(\frac{\eta}{\sqrt{m}} \cdot \gamma_{i,j} \cdot \Lambda_{i,j} \cdot \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)\right]\right) \cdot \|\mathbf{u}_{i,j}\|_{2}^{2} \quad (21)$$

$$- \mathcal{O}\left(\frac{\eta}{m} \|\mathbf{u}_{i,j}\|_{2}^{2} + \eta\sigma_{n}\sqrt{d} \|\mathbf{u}_{i,j}\|_{2} + \frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \eta d\sigma_{n}\sigma_{p}\right).$$

• The increment of model output for the other class 3 - y satisfies

$$\Delta_{3-y}^{(t)}(\mathbf{x}) = F_{3-y}\left(\mathbf{W}^{(t+1)}, \mathbf{x}\right) - F_{3-y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)$$
$$\leq \mathcal{O}\left(\eta\sigma_n\sqrt{d} \|\mathbf{u}_{i,j}\|_2 + \frac{\eta}{m\sqrt{n}}d\sigma_p^2 + \eta d\sigma_n\sigma_p\right).$$
(22)

Proposition F.1 shows that the model output on the data (\mathbf{x}, y) with respect to the targeted class, $F_y(\mathbf{W}, \mathbf{x})$, tends to increase over iterations (see term 1 in the RHS of (21)). However, due to the noise perturbation introduced by DP-SGD, the model increment $\Delta_{u}^{(t)}$ becomes smaller and cannot always be positive (because of the second negative term in RHS of (21)). In addition, the model output on the data (\mathbf{x}, y) with respect to the other class may increase due to the randomness from batches and DP-SGD. The results of Proposition F.3 allow us to track the test loss increments, as shown in the following subsection.

F.3 CHALLENGE 3: NON-LINEARITY OF CROSS-ENTROPY AND SOFTMAX FUNCTIONS

Due to the non-linearity of cross-entropy and softmax functions, the model output increment bounds in Proposition F.1 cannot be directly applied to bounding test loss.

To tackle this challenge, we bound the non-linear functions with a piece-wise linear function, as stated in Lemma F.2.

Lemma F.2. Under Assumption 3.2, we have

$$\mathcal{L}\left(\mathbf{W}^{(t+1)}, \mathbf{x}, y\right) - \mathcal{L}\left(\mathbf{W}^{(t)}, \mathbf{x}, y\right) \leq c_1 \cdot \sigma\left(\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_y^{(t)}(\mathbf{x})\right) - c_2 \cdot \sigma\left(\Delta_y^{(t)}(\mathbf{x}) - \Delta_{3-y}^{(t)}(\mathbf{x})\right),$$

for some constants $c_1, c_2 > 0$.

Lemma F.2 allows us to apply the model output increment bounds in Proposition F.1 to bound the test loss, as shown in Proposition F.3.

Proposition F.3. Under Condition 3.1 and Assumption 3.2, with probability at least $1 - \exp(-\Omega(d))$, for any $i \in [2], j \in \{maj, min\}$, we have

$$\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t+1)}\right) - \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t)}\right) \leq -\Omega\left(\frac{\eta}{\sqrt{m}} \cdot \gamma_{i,j} \cdot \Lambda_{i,j} \cdot \|\mathbf{u}_{i,j}\|_{2}^{2}\right) \cdot \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)\right] \\ + \underbrace{\mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2} + \eta\sigma_{n}\sqrt{d} \|\mathbf{u}_{i,j}\|_{2} + \frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \eta d\sigma_{n}\sigma_{p}\right)}_{:=\phi}.$$

With the fact that under Assumption 3.2,

$$1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right) = \Theta\left(1\right) \cdot \mathcal{L}\left(\left(\mathbf{W}^{(t)}\right), \mathbf{x}, y\right)$$
(23)

(25)

holds, Proposition F.3 can be applied to establish the following test loss bound

$$\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t+1)}\right) \leq \left(1 - \Omega\left(\frac{\eta}{\sqrt{m}} \cdot \gamma_{i,j} \cdot \Lambda_{i,j} \cdot \|\mathbf{u}_{i,j}\|_{2}^{2}\right)\right) \cdot \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t)}\right) + \phi.$$
(24)

Recursively applying (24) over T iterations yields Theorem 3.5.

PROOF G

G.1 PRELIMINARIES

Lemma G.1 (Gaussian distribution tail bound). A variable x following $\mathcal{N}(0, \sigma_0^2)$ satisfies

 $\mathbb{P}[x \ge t\sigma_0], \mathbb{P}[x \le -t\sigma_0] \le \exp\left(-\frac{t^2}{2}\right), \forall t \ge 0.$

P18 Lemma G.2 (Chi-squared distribution tail bound). A variable $\mathbf{x} \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_d)$ satisfies With probability at least $1 - \exp(-td/10)$, we have

$$\mathbb{P}\left[\left\|\mathbf{x}\right\|_{2} \ge \sigma_{0}\sqrt{2(t+1)d}\right] \le \exp(-(t+1)d/10), \forall t \ge 0.$$
(26)

Lemma G.3 (Half-normal distribution concentration bound). Suppose $x_1, x_2, \dots, x_n \sim \mathcal{N}(0, \sigma_0^2)$. *Then,*

$$\mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}|x_{i}|-\sqrt{\frac{2}{\pi}}\sigma_{0}\geq t\sigma_{0}\right], \mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}|x_{i}|-\sqrt{\frac{2}{\pi}}\sigma_{0}\leq -t\sigma_{0}\right]\leq \exp\left(-\frac{t^{2}}{2}\right), \forall t\geq 0.$$
(27)

Proof. First, half-normal variables $|x_i|, \forall i \in [n]$ are sub-Gaussian as a half-normal variable has a negative tail bounded by $-\sqrt{2/\pi}$ and a Gaussian delay positive tail. Then, by Hoeffding's inequality, we have

$$\mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}|x_{i}|-\sqrt{\frac{2}{\pi}}\sigma_{0}\geq t\sigma_{0}\right], \mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}|x_{i}|-\sqrt{\frac{2}{\pi}}\sigma_{0}\leq -t\sigma_{0}\right]\leq \exp\left(-\frac{t^{2}}{2}\right), \forall t\geq 0.$$
(28)

Lemma G.4. Let x_1, \dots, x_m be m independent zero-mean Gaussian variables. Denote z_i as indicators for signs of x_i , i.e., for all $i \in [m]$,

$$z_i = \begin{cases} 1, & x_i > 0, \\ 0, & x_i \le 0. \end{cases}$$
(29)

Then, we have

$$\mathbb{P}\left[\sum_{i=1}^{m} z_i \ge \frac{m}{4}\right] \ge 1 - \exp\left(-2m\right).$$
(30)

Proof. Because $z_i, i \in [m]$ are bounded in $[0,1], z_i, i \in [m]$ are sub-Gaussian variables. By Hoeffding's inequality, we have

$$\mathbb{P}\left[m \cdot \left(\frac{1}{m}\sum_{i=1}^{m} z_i\right) \le m \cdot \left(\frac{1}{2} - \epsilon\right)\right] \le \exp\left(\frac{2m^2\epsilon^2}{m(1/16)}\right).$$
(31)

Let $\epsilon = 1/4$, we have

$$\mathbb{P}\left[\sum_{i=1}^{m} z_i \le \frac{m}{4}\right] \le \exp\left(-2m\right).$$
(32)

Therefore, we have

$$\mathbb{P}\left[\sum_{i=1}^{m} z_i \ge \frac{m}{4}\right] \ge 1 - \exp\left(-2m\right).$$
(33)

This completes the proof.

Lemma G.5. Let x_1 be a Gaussian variable following $\mathcal{N}(0, \sigma_1)$ and x_2 be a Gaussian variable following $\mathcal{N}(0, \sigma_2)$. Then, with probability at least $1 - \exp(t_1^2/2) - \exp(t_2^2/2)$, we have

$$\langle x_1, x_2 \rangle \le t_1 t_2 \sigma_1 \sigma_2. \tag{34}$$

Lemma G.5 can be proved by using Lemma G.1.

Lemma G.6. For N Independent and Identically Distributed (IID) random variables $x_1, \dots, x_N \in [0, 1]$ with expectation μ , we have

$$\mathbb{P}\left[\frac{1}{N}\sum_{i=1}^{N}x_{i}-\mu\leq\sqrt{\frac{t}{N}}\right]\geq1-\exp\left(-2t\right)$$
(35)

with t > 0.

~

972 Lemma G.6 can be proved by Hoeffding's inequality. 973 **Lemma G.7.** For any constant $t \in (0, 1]$ and $x \in [-a, b]$, a, b > 0, we have 974 $\log(1 + t \cdot (\exp(x) - 1)) < \Gamma(x)x,$ (36)975 where $\Gamma(x) = \mathbb{I}(x > 0) + [\log(1 + t \cdot (\exp(-a) - 1))/(-a)] \cdot \mathbb{I}(x < 0).$ 976 977 *Proof.* First, considering $x \ge 0$, we have 978 $\frac{\partial \log(1+t\cdot(\exp(x)-1))}{\partial t} = \frac{\exp(x)-1}{1+t\cdot(\exp(x)-1)} \ge 0.$ 979 (37)980 Thus, $\log (1 + t \cdot (\exp(x) - 1)) \le x, \forall x > 0$. Second, considering x < 0, we have 981 $\frac{\partial^2 \log(1 + tv(\exp(x) - 1))}{\partial x^2} = \frac{(1 - t)t\exp(x)}{[1 + t(\exp(x) - 1)]^2} \ge 0.$ So $\log(1 + t(\exp(x) - 1))$ is a convex function of x. We can conclude that 982 (38)983 984 985 $\log(1 + t(\exp(x) - 1)) \le \frac{\log(1 + t(\exp(-a) - 1))}{-a}x, \forall x < 0.$ (39)986 987 This completes the proof. 988 **Lemma G.8.** For $x \in [x_0, 1]$ and $x_0 > 0$, we have 989 $1 - x \ge \frac{1 - x_0}{-\log(x_0)} \cdot (-\log(x)).$ 990 (40)991 992 Lemma G.8 can be proved by applying the convexity of $-\log(x)$. 993 **Lemma G.9.** For a geometric sequence defined as $z_{t+1} = \beta z_t$ for a constant $\beta < 1$, we have 994 $\sum_{t=1}^{T} z_t = \frac{1-\beta^T}{1-\beta} \cdot z_1.$ 995 (41)996 997 Lemma G.9 is obtained from the property of Geometric sequences. 998 **Lemma G.10.** For $x \leq x_0$, we have 999 $\frac{1}{\exp(x)+1} \ge \frac{1}{\exp(\bar{x}_0)+1} - \frac{\exp(\bar{x}_0)}{(\exp(\bar{x}_0)+1)^2} \cdot (x-\bar{x}_0),$ 1000 (42)1001 1002 where $\bar{x}_0 = |x_0|$. 1003 Lemma G.10 can be prove by the monotonicity and convexity of function $f(x) = 1/\exp(x) + 1$ 1004 with x > 0. 1005 **Lemma G.11.** In each iteration t, with probability at least $1 - \exp(-\Omega(d))$, for any $(\mathbf{x}, y) \in \mathcal{D}_{i,j}$, we have $\left\|\nabla_{\mathbf{W}^{(t)}}\mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y)\right\|_{2} \leq \mathcal{O}\left(\frac{1}{\sqrt{m}} \cdot \left(\left\|\mathbf{u}_{i,j}\right\|_{2} + \sigma_{p}\sqrt{d}\right)\right).$ 1008 (43)1009 1010 Lemma G.11 follows from Lemma G.2. 1011 **Lemma G.12.** For a variable $x \in [a, b](a < 0, b > 0)$, the function $f(x) = \log(1 + x)$ satisfies 1012 $f(x) \geq \frac{\log(1+b)}{{}^{\mathbf{h}}} x \cdot \mathbb{I}(x \geq 0) + \frac{\log(1+a)}{-a} x \cdot \mathbb{I}(x < 0).$ 1013 (44)1014 1015 Lemma G.12 can be proved by the monotonicity and concavity of the $log(\cdot)$ function. 1016 **Lemma G.13.** For any $(\mathbf{x}, y) \sim \mathcal{D}$, With probability at least 1 - 1/d, we have 1017 $\frac{\sigma_p^2 d}{2} \le \|\boldsymbol{\xi}\|_2 \le \frac{3\sigma_p^2 d}{2}$ (45)1018 1019 1020 *Proof.* By Bernstein's inequality, with probability 1 - 1/d, we have 1021 $|\|\boldsymbol{\xi}\|_{2} - \sigma_{p}^{2}(d-2)| = \mathcal{O}(\sigma_{p}^{2}\sqrt{d\log(2d)}).$ (46)1022 As d > 50, we have 1023 $\frac{\sigma_p^2 d}{2} \le \|\boldsymbol{\xi}\|_2 \le \frac{3\sigma_p^2 d}{2},$ 1024 (47)1025

with probability 1 - 1/d.

19

1026 G.2 PROOF OF THEOREM 3.5

In this subsection, we will prove Theorem 3.5. For convenience, we first define the clipping multiplier of data (\mathbf{x}, y) as

$$h(C, \mathbf{x}, y) = \frac{1}{\max\{1, \|\nabla \mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y)\|_2 / C\}}.$$
(48)

1034 Then, we compute the gradient of the neural networks and prove a bound for it.

1036 G.2.1 NETWORK GRADIENT

• (---(t))

1037 The stochastic gradient on $\mathbf{w}_{q,r}, q \in \{1, 2\}$ at iteration t is

$$\nabla_{\mathbf{w}_{q,r}^{(t)}} \mathcal{L}_{\mathcal{S}}(\mathbf{W}^{(t)}) = -\frac{1}{mB} \cdot \sum_{(\mathbf{x},y)\in\mathcal{S}^{(t)}} \left[\mathbb{I}\left(y=q\right) \cdot \left(1 - \operatorname{prob}_{q}(\mathbf{W}^{(t)},\mathbf{x})\right) \cdot \sum_{j=1}^{2} \sigma'\left(\left\langle \mathbf{w}_{q,r}^{(t)},\mathbf{x}^{(j)}\right\rangle\right) \cdot \mathbf{x}^{(j)} \right] + \frac{1}{mB} \cdot \sum_{(\mathbf{x},y)\in\mathcal{S}^{(t)}} \left[\mathbb{I}\left(y\neq q\right) \cdot \operatorname{prob}_{q}(\mathbf{W}^{(t)},\mathbf{x}) \cdot \sum_{j=1}^{2} \sigma'\left(\left\langle \mathbf{w}_{q,r}^{(t)},\mathbf{x}^{(j)}\right\rangle\right) \cdot \mathbf{x}^{(j)} \right].$$

$$(49)$$

G.2.2 BOUND OF THE CLIPPING MULTIPLIER $h(C, \mathbf{x}, y)$

1050 By definition (48), we know that

$$h(C, \mathbf{x}, y) \le 1. \tag{50}$$

In addition, from Lemma G.11, we know that with probability at least $1 - \exp(\Omega(d))$,

$$h(C, \mathbf{x}, y) \ge \Omega\left(\frac{C\sqrt{m}}{\|\mathbf{u}_{i,j}\|_2 + \sigma_p\sqrt{d}}\right).$$
(51)

1058 G.2.3 LOSS INCREMENT

For any data $(\mathbf{x}, y) \sim \mathcal{D}_{i,j}, i \in \{1, 2\}, j \in \{\text{Maj}, \text{Min}\}$, with some rearrangement, we can express the increment of the loss as

$$\mathcal{L}(\mathbf{W}^{(t+1)}, \mathbf{x}, y) - \mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y) = -\log\left(\operatorname{prob}_{y}\left(\mathbf{W}^{(t+1)}, \mathbf{x}\right)\right) + \log\left(\operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)\right)$$
$$= \log\left(1 + \left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}\right)\right)\left(\exp\left(\Delta_{3-y}^{(t)}\left(\mathbf{x}\right) - \Delta_{y}^{(t)}\left(\mathbf{x}\right)\right) - 1\right)\right),$$
(52)

where $\Delta_y^{(t)}(\mathbf{x}) = F_y^{(t+1)}(\mathbf{x}) - F_y^{(t)}(\mathbf{x})$, $\Delta_{3-y}^{(t)}(\mathbf{x}) = F_{3-y}^{(t+1)}(\mathbf{x}) - F_{3-y}^{(t)}(\mathbf{x})$ represent the model output increments at iteration t. As we can see in (52), one key factor that control the loss increment is $\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_y^{(t)}(\mathbf{x})$. We then bound the term it as follows. We first decompose $\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_y^{(t)}(\mathbf{x})$ as following.

where $\boldsymbol{\xi}$ is the noise patch of a data sample (\mathbf{x}, y) generated from $\mathcal{D}_{i,j}$. We then upper bound A, B and lower bound C, D to find the upper bound of $\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_y^{(t)}(\mathbf{x})$.

Here, we prove that $\Delta_{3-y}^{(t)}\left(\mathbf{x}\right) - \Delta_{y}^{(t)}\left(\mathbf{x}\right)$ is bounded.

Lemma G.14. For any $(\mathbf{x}, y) \sim \mathcal{D}$, with probability at least $1 - \exp{-\tilde{\Omega}(d)}$, we have

$$\left|\Delta_{3-y}^{(t)}\left(\mathbf{x}\right) - \Delta_{y}^{(t)}\left(\mathbf{x}\right)\right| \le \mathcal{O}\left(\eta(C + \sqrt{d}\sigma_{n})\left(\max_{i,j} \left\|\mathbf{u}_{i,j}\right\|_{2} + \sqrt{d}\sigma_{p}\right)\right).$$
(54)

Proof. With Lemma G.2, we have

$$\begin{aligned} |\Delta_{3-y}^{(t)}\left(\mathbf{x}\right) - \Delta_{y}^{(t)}\left(\mathbf{x}\right)| &\leq 2\eta \left(C + \left\|\mathbf{n}^{(t)}\right\|_{2}\right) \left(\max_{i,j} \left\|\mathbf{u}_{i,j}\right\|_{2} + \left\|\boldsymbol{\xi}\right\|_{2}\right) \\ &\leq \mathcal{O}\left(\eta (C + \sqrt{d}\sigma_{n})(\max_{i,j} \left\|\mathbf{u}_{i,j}\right\|_{2} + \sqrt{d}\sigma_{p})\right), \end{aligned} \tag{55}$$

with probability at least $1 - \exp(-\tilde{\Omega}(d))$. By the learning rate condition in Condition 3.1, we can conclude that $|\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_{y}^{(t)}(\mathbf{x})|$ is upper bounded by a constant.

For the term A, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we have the following inequality,

$$A = \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)} - \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{i,j}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \cdot \mathbf{u}_{i,j} \right) \right] \\ + \eta \cdot \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right) \right] - \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right] - \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right] \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right] - \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right| \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)}, \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \\ = \frac{1}{m} \sum_{r=1}^{m} \left[\left| \left\langle \eta \mathbf{n}^{(t)} \right\rangle \right] \right] \\ = \frac{1}{m} \sum_{r=1}^{m$$

where (a) is obtained by the monotonicity of ReLU activation function; (b) is because ReLU function is 1-Lipschitz continuous; (c) is due to Lemma G.3.

For the term B, we have that with probability at least $1 - \exp(-\hat{\Omega}(d))$, $B = \frac{1}{m} \sum_{r=1}^{m} \left| \sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)} - \frac{\eta}{mB} \cdot \sum_{(x_{1}-y_{2}) \in \mathcal{S}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \cdot \boldsymbol{\xi}_{k} \right) \right| \right|$ $+\frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{k}, y_{k}) \in \mathcal{S}_{c}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y, r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \cdot h(C, \mathbf{x}, y) \cdot \left(1 - \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \cdot \boldsymbol{\xi}_{k} + \eta \mathbf{n}_{3-y, r}^{(t)}, \boldsymbol{\xi}_{k} \right) \right)$ $-\frac{1}{m}\sum_{m}^{m}\left[\sigma\left(\left\langle \mathbf{w}_{3-y,r}^{(t)},\boldsymbol{\xi}\right\rangle\right)\right]$ $\overset{(a)}{\leq} \frac{1}{m} \sum_{r=1}^{m} \left| \left| \left\langle \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}, y_{k}) \in S^{(t)}} \boldsymbol{\xi}_{k}, \boldsymbol{\xi} \right\rangle \right| + \left| \left\langle \eta \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right|$ $\stackrel{(b)}{\leq} \mathcal{O}\left(\frac{\eta}{m\sqrt{B}}\sqrt{d}\sigma_p + \eta\sqrt{d}\sigma_n\right) \cdot \mathcal{O}\left(\sqrt{d}\sigma_p\right)$ $=\mathcal{O}\left(\frac{\eta}{m\sqrt{B}}d\sigma_p^2 + \eta d\sigma_n \sigma_p\right),$ (57) where (a) is because $\sigma'(\cdot) \ge 0$, prob_{*u*}, prob_{3-*y*} $\in [0, 1]$ and ReLU function is 1-Lipschitz continuous; (b) is because of Cauchy-Schwarz inequality, the property of 1-norm and 2-norm, and Lemma G.3. For the term C, we have $C = \frac{1}{m} \sum_{r=1}^{m} \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)} + \frac{\eta}{mB} \cdot \sum_{(u,v) \in \mathbf{S}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \cdot h(C, \mathbf{x}_k, y_k) \cdot \left(1 - \operatorname{prob}_y \left(\mathbf{W}^{(t)}, \mathbf{x}_k \right) \right) \cdot \mathbf{u}_{i,j} \right) \right)$ $+\eta \cdot \mathbf{n}_{y,r}^{(t)}, \mathbf{u}_{i,j} \rangle - \frac{1}{m} \sum_{i=1}^{m} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right)$ (58)Based on Lemma G.4, we can conclude that with probability at least $1 - \exp(-2m)$, the number of activated neurons at iteration t are at least m/4. Then, with probability at least $1 - \exp(-\Omega(d))$, we have $C \ge \frac{1}{m} \sum_{r=1}^{m} \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)} + \frac{\eta}{mB} \sum_{(\mathbf{x}_{i}, y_{i}) \in \mathbf{S}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \right) \right)$ $\mathbf{u}_{i,j}, \mathbf{u}_{i,j} \rangle) - \frac{1}{m} \sum_{j=1}^{m} \left| \left\langle \eta \cdot \mathbf{n}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right| - \frac{1}{m} \sum_{j=1}^{m} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right)$

$$\geq \Omega\left(\frac{\eta C}{B\sqrt{m}(\|\mathbf{u}_{i,j}\|_{2} + \sigma_{p}\sqrt{d})}\right) \sum_{(\mathbf{x}_{k}, y_{k})\in\mathcal{S}_{i,j}^{(t)}} \left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}_{k}\right)\right) \|\mathbf{u}_{i,j}\|_{2}^{2} - \mathcal{O}\left(\eta\sigma_{n}\sqrt{d}\|\mathbf{u}_{i,j}\|_{2}\right)$$

$$(59)$$

1180 The second inequality is by using the bound of the clipping multiplier.

Therefore, by Lemmas G.3 and G.6, with probability at least $1 - \exp(-\Omega(d)) - \exp(-\Omega(n))$, we have

$$C \geq \Omega \left(\frac{\eta \gamma_{i,j} C \|\mathbf{u}_{i,j}\|_{2}^{2}}{\sqrt{m} (\|\mathbf{u}_{i,j}\|_{2} + \sigma_{p} \sqrt{d})} \right) \mathbb{E}_{(\mathbf{x}_{k}, y_{k}) \sim \mathcal{D}_{i,j}} \left[1 - \operatorname{prob}_{y_{k}} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right] - \mathcal{O} \left(\frac{\eta}{m} \|\mathbf{u}_{i,j}\|_{2}^{2} \right) - \mathcal{O} \left(\frac{\eta}{m} \|\mathbf{u}_{i,j}\|_{2}^{2} \right) - \mathcal{O} \left(\eta \sigma_{n} \sqrt{d} \|\mathbf{u}_{i,j}\|_{2} \right).$$
(60)

In the following, we prove the bound of D. Similar to the proof of bound of B, we have that with probability at least $1 - \exp\left(-\tilde{\Omega}(d)\right) - \exp\left(-\Omega(n)\right)$, we have

ιPΡ $x_{3-y}(\mathbf{x})$ $-\Delta y$ (**x**) probability at least $1 - \exp\left(-\tilde{\Omega}(d)\right)$,

$$\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_{y}^{(t)}(\mathbf{x})$$

$$\leq - \underbrace{\Omega\left(\frac{\eta\gamma_{i,j}C \|\mathbf{u}_{i,j}\|_{2}^{2}}{\sqrt{m}(\|\mathbf{u}_{i,j}\|_{2} + \sigma_{p}\sqrt{d})}\right) \cdot \left[\mathbb{E}_{(\mathbf{x}_{k},y_{k})\sim\mathcal{D}_{i,j}}\left[1 - \operatorname{prob}_{y_{k}}\left(\mathbf{W}^{(t)},\mathbf{x}_{k}\right)\right]\right]}_{\Phi_{1}} \qquad (62)$$

$$+ \underbrace{\mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}}\|\mathbf{u}_{i,j}\|_{2}^{2}\right) + \mathcal{O}\left(\eta\sigma_{n}\sqrt{d}\|\mathbf{u}_{i,j}\|_{2}\right) + \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \eta d\sigma_{n}\sigma_{p}\right)}_{\Phi_{2}}.$$

Armed with the loss increment bound (62), we prove the test loss bound of each data $(\mathbf{x}, y) \sim \mathcal{D}$ in the next sub section.

G.2.4 TEST LOSS BOUND

Under Assumption 3.2, for any $(\mathbf{x}, y) \sim \mathcal{D}$, we have

$$-\operatorname{prob}_{y}\left(\mathbf{W}^{(t)},\mathbf{x}\right) \ge 1 - \exp(-s).$$
 (63)

By (52) and Lemma G.7, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we can upper bound the loss increment by a piece-wise linear function,

 $\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t+1)}\right) - \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t)}\right)$ $= \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}_{i,j}} \left[\log \left(1 + \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x} \right) \right) \cdot \left(\exp \left(\Delta_{3-y}^{(t)} \left(\mathbf{x} \right) - \Delta_{y}^{(t)} \left(\mathbf{x} \right) \right) - 1 \right) \right) \right]$ $\stackrel{(a)}{\leq} -\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[\Gamma\left(\Delta_{3-y}^{(t)}\left(\mathbf{x}\right)-\Delta_{y}^{(t)}\left(\mathbf{x}\right)\right)\cdot\Phi_{1}\right]+\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[\Gamma\left(\Delta_{3-y}^{(t)}\left(\mathbf{x}\right)-\Delta_{y}^{(t)}\left(\mathbf{x}\right)\right)\cdot\Phi_{2}\right]$ (64)

where (a) is by Lemma G.7 and Lemma (62). Then, substituting (62) to the above inequality yields

(65)where (a) is obtain by Lemma G.8. $\underline{\Gamma} = -\log(1 + (1 - \exp(-s)) \cdot (\exp(-a) - 1))/a, \gamma =$ $-\exp(-s)/\ln(1-\exp(-s))$ and a is the lower bound of $\Delta_{3-y}^{(t)} - \Delta_y^{(t)}$ (By Lemma G.14, $\Delta_{3-y}^{(t)} - \Delta_y^{(t)}$) is lower bounded by a constant). Therefore, we have

$$\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t+1)}\right) \leq \left(1 - \Omega\left(\frac{\eta\gamma_{i,j}\Lambda_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{\sqrt{m}}\right)\right) \cdot \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t)}\right) + \mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right) \\ + \mathcal{O}\left(\eta\sigma_{n}\sqrt{d} \|\mathbf{u}_{i,j}\|_{2}\right) + \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \eta d\sigma_{n}\sigma_{p}\right).$$

(66)

Then, combining all T iterations and using Lemma G.9, we have

$$\begin{aligned}
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(T)}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(T)}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) &\leq \left(1 - \Omega\left(\frac{\eta\gamma_{i,j}\Lambda_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{\sqrt{m}}\right)\right)^{T} \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \left(\mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right) + \mathcal{O}\left(\eta\sigma_{n}\sqrt{d} \|\mathbf{u}_{i,j}\|_{2}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \left(\mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right)\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\frac{\sqrt{\frac{d}{mn}}}{\sqrt{\frac{d}{nn}}}\frac{1}{\gamma_{i,j}\Lambda_{i,j}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}\Lambda_{i,j}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}\Lambda_{i,j}}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}\Lambda_{i,j}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}\Lambda_{i,j}}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) + \mathcal{O}\left(\sqrt{\frac{d}{mn}}\frac{1}{\gamma_{i,j}}}\right) \\
\mathbf{\mathcal{L}}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(0)}\right) \\
\mathbf{\mathcal{L}}_{\mathcal$$

Setting the parameters with Condition 3.1 yields the conclusion. This completes the proof.

G.3 PROOF OF THEOREM 3.6

Proof. Recall that in (52), we have

$$\mathcal{L}(\mathbf{W}^{(t+1)}, \mathbf{x}, y) - \mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y)$$

$$= \log \left(1 + \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x} \right) \right) \left(\exp \left(\Delta_{3-y}^{(t)} \left(\mathbf{x} \right) - \Delta_{y}^{(t)} \left(\mathbf{x} \right) \right) - 1 \right) \right)$$

$$\overset{(a)}{\geq} c_{0}^{(t)} \cdot \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x} \right) \right) \left(\exp \left(\Delta_{3-y}^{(t)} \left(\mathbf{x} \right) - \Delta_{y}^{(t)} \left(\mathbf{x} \right) \right) - 1 \right)$$

$$\overset{(b)}{=} \Omega \left(\exp \left(\Delta_{3-y}^{(t)} \left(\mathbf{x} \right) - \Delta_{y}^{(t)} \left(\mathbf{x} \right) \right) - 1 \right),$$
(68)

where $c_0^{(t)} > 0$ for any $t \in \{0\} \cup [T-1]$ are constants. Here (a) is obtained from Lemma G.12 Then, we bound $\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_y^{(t)}(\mathbf{x})$; (b) is by $1 - \exp(-s) \le 1 - \operatorname{prob}_y(\mathbf{W}, \mathbf{x}) \le 1$ with Assumption

3.2. Next, we will prove a lower bound of $\Delta_{3-y}^{(t)}(\mathbf{x}) - \Delta_{y}^{(t)}(\mathbf{x})$. Recall that in (53), we have $\Delta_{3-u}^{(t)}(\mathbf{x}) - \Delta_{u}^{(t)}(\mathbf{x})$ $= \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{3-y,r}^{(t+1)}, \mathbf{u}_{i,j} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right]}_{A} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{3-y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{B} - \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \mathbf{u}_{i,j} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \right]}_{C} - \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\xi} \right\rangle \right) - \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \right]}_{D} + \underbrace{\frac{1}{m} \sum_{r=1}^{m} \left[\sigma\left($ (69)We then find the lower bounds of A, B and upper bounds of C, D to obtain the lower bound of $\Delta_{3-y}^{\left(t\right)}\left(\mathbf{x}\right) - \Delta_{y}^{\left(t\right)}\left(\mathbf{x}\right).$ For the term A, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we have $A = \frac{1}{m} \sum_{r=1}^{m} \left| \sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)} - \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{1}, y_{1}) \in \mathbf{S}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \cdot \mathbf{u}_{i,j} \right\rangle \right) \right|$ + $\eta \cdot \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \rangle \bigg] - \frac{1}{m} \sum_{i=1}^{m} \bigg[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle \right) \bigg]$ $\stackrel{(a)}{\geq} - \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_k, y_k) \in \mathcal{S}_{i}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y, r}^{(t)}, \mathbf{u}_{i, j} \right\rangle \right) \cdot h(C, \mathbf{x}_k, y_k) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_k \right) \cdot \left\| \mathbf{u}_{i, j} \right\|_2^2$ $+\frac{1}{m}\sum_{i=1}^{m}\langle \eta\cdot\mathbf{n}_{3-y,r}^{(t)},\mathbf{u}_{i,j}\rangle,$ $\overset{(b)}{\geq} - \frac{\eta}{mB} \sum_{(\mathbf{x}_{i}, w_{i}) \in \mathbf{S}^{(t)}} \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \cdot \|\mathbf{u}_{i,j}\|_{2}^{2} + \langle \eta \cdot \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \rangle$ $\overset{(c)}{\geq} - \frac{\eta \gamma_{i,j}}{m} \mathbb{E}_{(\mathbf{x}_k, y_k) \sim \mathcal{D}_{i,j}} \left(1 - \operatorname{prob}_y \left(\mathbf{W}^{(t)}, \mathbf{x}_k \right) \right) \cdot \|\mathbf{u}_{i,j}\|_2^2 - \frac{\eta}{m} \sqrt{\frac{d}{B}} \|\mathbf{u}_{i,j}\|_2^2 + \frac{\eta}{m} \sum_{i=1}^m \langle \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \rangle,$ (70)

where (a) is due to the Condition 3.1. As $\|\mathbf{u}_{i,j}\|_2 = \Theta(\sqrt{d}\sigma_p)$ and $\sigma_p = \Omega(\sigma_n)$, we have that $\|\mathbf{u}_{i,j}\|_2 = \Omega(\sqrt{d}\sigma_n)$. (b) is by the fact that $\sigma'(\cdot)$, $h(C, \mathbf{x}_k, y_k) \le 1$; (c) is by Lemma G.6.

For the term B, with probability at least $1 - \exp\left(-\tilde{\Omega}(d)\right)$, we have

$$B = \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)} - \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{y}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \cdot \boldsymbol{\xi}_{k} \right) \right]$$

$$+ \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{3-y}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \cdot h(C, \mathbf{x}, y) \cdot \left(1 - \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \cdot \boldsymbol{\xi}_{k} + \eta \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right) \right]$$

$$- \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \right]$$

$$\geq - \frac{\eta}{mB} \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{y}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{3-y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \operatorname{prob}_{3-y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \cdot \left| \langle \boldsymbol{\xi}_{k}, \boldsymbol{\xi} \rangle \right| + \langle \eta \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi} \rangle$$

$$\geq - \mathcal{O} \left(\frac{\eta}{m\sqrt{B}} d\sigma_{p}^{2} \right) + \frac{1}{m} \sum_{r=1}^{m} \eta \langle \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi} \rangle,$$
For the term *C*, with probability at least $1 - \exp\left(-\tilde{\Omega}(d)\right)$, we have

, with probability $\operatorname{cp}\left(-\mathfrak{U}(a)\right),$

$$C = \frac{1}{m} \sum_{r=1}^{m} \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)} + \frac{\eta}{mB} \cdot \sum_{(\mathbf{x}_{k}, y_{k}) \in \mathcal{S}_{i,j}^{(t)}} \sigma'\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle\right) \cdot h(C, \mathbf{x}_{k}, y_{k}) \cdot \left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}_{k}\right)\right) \cdot \mathbf{u}_{i,j} \right) + \eta \cdot \mathbf{n}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle - \frac{1}{m} \sum_{r=1}^{m} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle\right) = \frac{1}{m} \sum_{r=1}^{m} \sigma\left(\left\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{u}_{i,j} \right\rangle\right) = \frac{1}{1882} \leq 0,$$

$$(72)$$

where the inequality is by Condition 3.1, which implies $\|\mathbf{u}_{i,j}\|_2 = \Omega(\sqrt{d}\sigma_n)$. For the term D, with probability at least $1 - \exp\left(-\tilde{\Omega}(d)\right)$, we have

$$D = \frac{1}{m} \sum_{r=1}^{m} \left[\sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)} - \frac{\eta}{mB} \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{3-y}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) h(C, \mathbf{x}_{k}, y_{k}) \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \boldsymbol{\xi}_{k} \right. \\ \left. + \frac{\eta}{mB} \sum_{(\mathbf{x}_{k},y_{k}) \in \mathcal{S}_{y}^{(t)}} \sigma' \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi}_{k} \right\rangle \right) h(C, \mathbf{x}_{k}, y_{k}) \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \boldsymbol{\xi}_{k} + \mathbf{n}_{y,r}^{(t)}, \boldsymbol{\xi} \right) \right] \\ \left. - \frac{1}{m} \sum_{r=1}^{m} \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \\ \left. - \frac{1}{m} \sum_{r=1}^{m} \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle \right) \\ \left. \leq \mathcal{O} \left(\frac{\eta}{m\sqrt{n}} d\sigma_{p}^{2} \right) + \frac{1}{m} \sum_{r=1}^{m} \eta \left\langle \mathbf{n}_{y,r}^{(t)}, \boldsymbol{\xi} \right\rangle,$$

$$(73)$$

where the inequality is by Condition 3.1 that $\sigma_n = \mathcal{O}(\sigma_p), h(C, \mathbf{x}_k, y_k) < 1.$

Combining the bounds together, we have $\Lambda^{(t)}$ (m) $\Lambda^{(t)}$ (m) $\eta \gamma_{i,j}$ (1

$$\Delta_{3-y}^{(t)}\left(\mathbf{x}\right) - \Delta_{y}^{(t)}\left(\mathbf{x}\right) \geq -\frac{\eta\gamma_{i,j}}{m} \mathbb{E}_{\left(\mathbf{x}_{k},y_{k}\right)\sim\mathcal{D}_{i,j}}\left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)},\mathbf{x}_{k}\right)\right) \cdot \|\mathbf{u}_{i,j}\|_{2}^{2} + \frac{1}{m} \sum_{r=1}^{m} \langle \eta \cdot \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \rangle - \mathcal{O}\left(\frac{\eta}{m\sqrt{B}}d\sigma_{p}^{2}\right) + \frac{1}{m} \sum_{r=1}^{m} \langle \eta \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi} \rangle - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2}\right) - \frac{1}{m} \sum_{r=1}^{m} \eta \langle \mathbf{n}_{y,t}^{(t)}, \boldsymbol{\xi} \rangle - \mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right),$$

$$(74)$$

with probability at least
$$1 - \exp(-\tilde{\Omega}(d))$$
. Substituting (74) into (68), we have
 $\mathbb{E}[\mathcal{L}(\mathbf{W}^{(t+1)}, \mathbf{x}, y)] - \mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y)$

$$\geq \mathbb{E} \left[\Omega \left(\exp \left(\Delta_{3-y}^{(t)} \left(\mathbf{x} \right) - \Delta_{y}^{(t)} \left(\mathbf{x} \right) \right) - 1 \right) \right]$$

$$\geq \Omega \left(\mathbb{E} \left[\exp \left(-\frac{\eta \gamma_{i,j}}{m} \mathbb{E}_{(\mathbf{x}_{k},y_{k}) \sim \mathcal{D}_{i,j}} \left(1 - \operatorname{prob}_{y} \left(\mathbf{W}^{(t)}, \mathbf{x}_{k} \right) \right) \cdot \| \mathbf{u}_{i,j} \|_{2}^{2} + \frac{1}{m} \sum_{r=1}^{m} \langle \eta \mathbf{n}_{3-y,r}^{(t)}, \boldsymbol{\xi} \rangle \right.$$

$$\left. + \frac{1}{m} \sum_{r=1}^{m} \langle \eta \cdot \mathbf{n}_{3-y,r}^{(t)}, \mathbf{u}_{i,j} \rangle - \frac{1}{m} \sum_{r=1}^{m} \eta \langle \mathbf{n}_{y,t}^{(t)}, \boldsymbol{\xi} \rangle - \mathcal{O} \left(\frac{\eta}{m\sqrt{n}} d\sigma_{p}^{2} + \frac{\eta}{m} \sqrt{\frac{d}{n}} \| \mathbf{u}_{i,j} \|_{2}^{2} \right) \right) - 1 \right] \right). \tag{75}$$

Here with a probability at least 1 - 1/d, we have

$$\mathbb{E}\left[\exp\left(\frac{1}{m}\sum_{r=1}^{m}\langle\eta\mathbf{n}_{3-y,r}^{(t)},\mathbf{u}_{i,j}\rangle + \frac{1}{m}\sum_{r=1}^{m}\langle\eta\mathbf{n}_{3-y,r}^{(t)},\boldsymbol{\xi}\rangle - \frac{1}{m}\sum_{r=1}^{m}\eta\langle\mathbf{n}_{y,t}^{(t)},\boldsymbol{\xi}\rangle\right)\right] \\
= \exp\left(\eta^{2}\frac{\|\mathbf{u}_{i,j}\|_{2}^{2}\sigma_{n}^{2} + 2\|\boldsymbol{\xi}\|_{2}^{2}\sigma_{n}^{2}}{2m}\right) \\
= \exp\left(\Theta\left(\eta^{2}\frac{\|\mathbf{u}_{i,j}\|_{2}^{2}\sigma_{n}^{2} + \sigma_{p}^{2}d\sigma_{n}^{2}}{2m}\right)\right),$$
(76)

where the least equality is by Lemma G.13. With a probability at least $1 - \tilde{O}(1/d)$, we have

$$\mathbb{E}[\mathcal{L}(\mathbf{W}^{(t+1)}, \mathbf{x}, y)] - \mathcal{L}(\mathbf{W}^{(t)}, \mathbf{x}, y) \\
\geq \Omega\left(-\frac{\eta\gamma_{i,j}}{m}\mathbb{E}_{(\mathbf{x}_{k}, y_{k})\sim\mathcal{D}_{i,j}}\left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t)}, \mathbf{x}_{k}\right) \cdot \|\mathbf{u}_{i,j}\|_{2}^{2} \\
+ \Omega\left(\frac{\eta^{2}\sigma_{n}^{2}\|\mathbf{u}_{i,j}\|_{2}^{2}}{2m}\right) + \Omega\left(\frac{\eta^{2}d\sigma_{n}^{2}\sigma_{p}^{2}}{2m}\right) - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2}\right) - \mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}}\|\mathbf{u}_{i,j}\|_{2}^{2}\right)\right), \\
\geq - \mathcal{O}\left(\frac{\eta\gamma_{i,j}}{m}\|\mathbf{u}_{i,j}\|_{2}^{2}\right)\mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t)}\right) - \mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}}\|\mathbf{u}_{i,j}\|_{2}^{2}\right) \\
+ \Omega\left(\frac{\eta^{2}\sigma_{n}^{2}\|\mathbf{u}_{i,j}\|_{2}^{2}}{2m}\right) + \Omega\left(\frac{\eta^{2}d\sigma_{n}^{2}\sigma_{p}^{2}}{2m}\right) - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2}\right).$$
(77)

where the second equality is by Lemma G.8 (the $(1 - \text{prob}_u(\mathbf{W}^{(t)}, \mathbf{x}))$ is almost surely lower bounded). Then, with a probability at least 1 - O(1/d), we have

$$\mathbb{E}_{\mathbf{n}^{(t)}} [\mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(t+1)})] \geq \left(1 - \mathcal{O}\left(\frac{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right)\right) \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(t)}) - \mathcal{O}\left(\frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right)$$
(78)

1456
1457
$$+ \Omega\left(\frac{\eta^2 \sigma_n^2 \|\mathbf{u}_{i,j}\|_2^2}{2m}\right) + \Omega\left(\frac{\eta^2 d\sigma_n^2 \sigma_p^2}{2m}\right) - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}} d\sigma_p^2\right).$$

Combining all the iterations, with a probability at least $1 - \tilde{\mathcal{O}}(1/d)$ we have

$$\mathbb{E}_{\mathbf{n}^{(0)},\cdots,\mathbf{n}^{(T-1)}}[\mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(T)})] \geq \left(1 - \mathcal{O}\left(\frac{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right)\right)^{T} \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(0)}) + \frac{1 - \left(1 - \mathcal{O}\left(\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}/m\right)\right)^{T}}{\mathcal{O}\left(\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}/m\right)} \qquad (79)$$

$$\left[\Omega\left(\frac{\eta^{2}\sigma_{n}^{2} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right) + \Omega\left(\frac{\eta^{2}d\sigma_{n}^{2}\sigma_{p}^{2}}{m}\right) - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right)\right].$$

With the number of iterations $T \ge \Omega\left(-1/\log\left(1 - \Omega\left(\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_2^2/m\right)\right)\right)$ and a probability at least $1 - \tilde{\mathcal{O}}(1/d)$, we have

$$\begin{split} \mathbb{E}_{\mathbf{n}^{(0)},\cdots,\mathbf{n}^{(T-1)}} [\mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(T)})] \\ &\geq \left(1 - \mathcal{O}\left(\frac{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right)\right)^{T} \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(0)}) + \Omega\left(\frac{m}{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}\right) \\ &\left[\Omega\left(\frac{\eta^{2}\sigma_{n}^{2} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right) + \Omega\left(\frac{\eta^{2}d\sigma_{n}^{2}\sigma_{p}^{2}}{m}\right) - \mathcal{O}\left(\frac{\eta}{m\sqrt{n}}d\sigma_{p}^{2} + \frac{\eta}{m}\sqrt{\frac{d}{n}} \|\mathbf{u}_{i,j}\|_{2}^{2}\right)\right] \\ &\geq \left(1 - \mathcal{O}\left(\frac{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right)\right)^{T} \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(0)}) \\ &+ \left[\Omega\left(\frac{\eta\sigma_{n}^{2}}{\gamma_{i,j}}\right) + \Omega\left(\frac{\eta d\sigma_{n}^{2}\sigma_{p}^{2}}{\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}\right) - \mathcal{O}\left(\frac{d\sigma_{p}^{2}}{\sqrt{n}\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}} + \frac{1}{\gamma_{i,j}}\sqrt{\frac{d}{n}}\right)\right] \\ &\geq \left(1 - \mathcal{O}\left(\frac{\eta\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}{m}\right)\right)^{T} \mathcal{L}_{\mathcal{D}_{i,j}}(\mathbf{W}^{(0)}) + \Omega\left(\frac{\eta d\sigma_{n}^{2}\sigma_{p}^{2}}{\gamma_{i,j} \|\mathbf{u}_{i,j}\|_{2}^{2}}\right) - \mathcal{O}\left(\frac{1}{\gamma_{i,j}}\sqrt{\frac{d}{n}}\right). \end{split}$$
This completes the proof.

This completes the proof.

G.4 PROOF OF THEOREM 3.9

Proof. Based on (52), for any $(\mathbf{x}, y) \sim \mathcal{D}_{i,j}$, we have

$$\mathcal{L}\left(\mathbf{W}^{(t+1)}, \mathbf{x} + \boldsymbol{\zeta}^{(t+1)}\left(\mathbf{x}\right), y\right) - \mathcal{L}\left(\mathbf{W}^{(t+1)}, \mathbf{x}, y\right)$$

= log $\left(1 + \left(1 - \operatorname{prob}_{y}\left(\mathbf{W}^{(t+1)}, \mathbf{x}\right)\right) \cdot \left(\exp\left(\tilde{\Delta}_{3-y}^{(t+1)}\left(\mathbf{x}\right) - \tilde{\Delta}_{y}^{(t+1)}\left(\mathbf{x}\right)\right) - 1\right)\right),$ (81)

where

$$\boldsymbol{\zeta}^{(t+1)}\left(\mathbf{x}\right) = \arg \max_{\|\boldsymbol{\zeta}\|_{p} \le \bar{\boldsymbol{\zeta}}} \mathcal{L}\left(\mathbf{W}^{(t+1)}, \mathbf{x} + \boldsymbol{\zeta}, y\right),$$
(82)

and

$$\tilde{\Delta}_{3-y}^{(t+1)}(\mathbf{x}) - \tilde{\Delta}_{y}^{(t+1)}(\mathbf{x}) = \underbrace{\frac{1}{m} \sum_{r=1}^{m} \sum_{j=1}^{2} \left[\sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t+1)}, \mathbf{x}^{(j)} + \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\rangle \right) - \sigma \left(\left\langle \mathbf{w}_{3-y,r}^{(t+1)}, \mathbf{x}^{(j)} \right\rangle \right) \right]}_{E}$$

$$- \frac{1}{m} \sum_{r=1}^{m} \sum_{j=1}^{2} \left[\sigma \left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \mathbf{x}^{(j)} + \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\rangle \right) - \sigma \left(\left\langle \mathbf{w}_{y,r}^{(t+1)}, \mathbf{x}^{(j)} \right\rangle \right) \right].$$
(83)

 \widetilde{F}

Then, we bound E and F.

For the term E, with probability at least $1 - \exp(-\Omega(d))$, we have $E \stackrel{(a)}{\leq} \frac{1}{m} \sum_{n=1}^{m} \sum_{i=1}^{2} \left| \left\langle \mathbf{w}_{3-y,r}^{(t+1)}, \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\rangle \right|$ $\stackrel{(b)}{\leq} \frac{1}{m} \sum_{n=1}^{m} \sum_{i=1}^{2} \left\| \mathbf{w}_{3-y,r}^{(t+1)} \right\|_{2} \left\| \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\|_{2}$ $\stackrel{(c)}{\leq} \frac{1}{m} \sum_{r=1}^{m} \sum_{i=1}^{2} \left\| \mathbf{w}_{3-y,r}^{(t+1)} \right\|_{2} \left\| \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\|_{p} d^{1-1/p}$ (84) $\leq \frac{2}{m} \sum_{r=1}^{m} \left(\sum_{t'=0}^{t} \left\| \mathbf{w}_{3-y,r}^{(t'+1)} - \mathbf{w}_{3-y,r}^{(t')} \right\|_{2} + \left\| \mathbf{w}_{3-y,r}^{(0)} \right\|_{2} \right) \bar{\zeta} d^{1-1/p}$ $\stackrel{(d)}{\leq} \mathcal{O}\left(\left[t\frac{\eta}{m}C + \frac{\eta}{m}\sqrt{td}\sigma_n + \sqrt{d}\sigma_0\right]\bar{\zeta}d^{1-1/p}\right),$

where (a) is because ReLU(\cdot) is 1-Lipschitz continuous; (b) is due to the Cauchy-Schwarz inequality; (c) is because of the Hölder's inequality; (d) is due to Lemma G.2.

Similarly, for the term F, with probability at least $1 - \exp(-\Omega(d))$, we have we have

$$F \geq -\frac{1}{m} \sum_{r=1}^{m} \sum_{j=1}^{2} \left\| \left\langle \mathbf{w}_{y,r}^{(t+1)}, \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\rangle \right\|$$

$$\geq -\frac{1}{m} \sum_{r=1}^{m} \sum_{j=1}^{2} \left\| \mathbf{w}_{y,r}^{(t+1)} \right\|_{2} \left\| \boldsymbol{\zeta}^{(t+1)} \left(\mathbf{x} \right)^{(j)} \right\|_{2}$$

$$\geq -\frac{2}{m} \sum_{r=1}^{m} \left(\sum_{t'=0}^{t} \left\| \mathbf{w}_{y,r}^{(t'+1)} - \mathbf{w}_{y,r}^{(t')} \right\|_{2} + \left\| \mathbf{w}_{y,r}^{(0)} \right\|_{2} \right) \bar{\zeta} d^{1-1/p}$$

$$\geq -\Omega \left(\left[t \frac{\eta}{m} C + \frac{\eta}{m} \sqrt{td} \sigma_{n} + \sqrt{d} \sigma_{0} \right] \bar{\zeta} d^{1-1/p} \right).$$

(85)

Combing with the bounds of E, F, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we have

$$\tilde{\Delta}_{3-y}^{(t+1)}\left(\mathbf{x}\right) - \tilde{\Delta}_{y}^{(t+1)}\left(\mathbf{x}\right) \le \mathcal{O}\left(\left[t\frac{\eta}{m}C + \frac{\eta}{m}\sqrt{td}\sigma_{n} + \sqrt{d}\sigma_{0}\right]\bar{\zeta}d^{1-1/p}\right).$$
(86)

Then, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we have

$$\mathcal{L}_{\mathcal{D}_{i,j}}^{\mathrm{adv}}\left(\mathbf{W}^{(t+1)}\right) - \mathcal{L}_{\mathcal{D}_{i,j}}\left(\mathbf{W}^{(t+1)}\right)$$

$$\leq \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}_{i,j}}\left[\Gamma\left(\tilde{\Delta}_{3-y}^{(t+1)}\left(\mathbf{x}\right) - \tilde{\Delta}_{y}^{(t+1)}\left(\mathbf{x}\right)\right)\left(\tilde{\Delta}_{3-y}^{(t+1)}\left(\mathbf{x}\right) - \tilde{\Delta}_{y}^{(t+1)}\left(\mathbf{x}\right)\right)\right]$$

$$\leq \mathcal{O}\left(\left[t\frac{\eta}{m}C + \frac{\eta}{m}\sqrt{td}\sigma_{n} + \sqrt{d}\sigma_{0}\right]\bar{\zeta}d^{1-1/p}\right).$$
(87)

Combining with Theorem 3.5 and setting parameters with Condition 3.1, with probability at least $1 - \exp(-\tilde{\Omega}(d))$, we have

$$\mathcal{L}_{\mathcal{D}_{i,j}}^{\mathrm{adv}}\left(\mathbf{W}^{(T)}\right) \leq \bar{L}_{i,j} + \mathcal{O}\left(\left[\frac{T}{m}C + \frac{\sqrt{Td}}{m}\sigma_n + \sqrt{d}\sigma_0\right]\bar{\zeta}d^{1-1/p}\right).$$
(88)

This completes the proof.

PROOF OF PROPOSITION 4.4 Η

Proof. We have

$$\sigma(\langle \tilde{\mathbf{w}}_{1,r}, \mathbf{u}_1' \rangle) = C_1 \cos \theta \|\mathbf{u}_1'\|_2 \|\mathbf{u}_1\|_2,$$
(89)

$$0 \le \sigma(\langle \tilde{\mathbf{w}}_{1,r}, \boldsymbol{\xi} \rangle) \le C_3 \sigma_p^2, \tag{90}$$

$$\sigma(\langle \tilde{\mathbf{w}}_{2,r}, \mathbf{u}_1' \rangle) = C_1 \sin \theta \|\mathbf{u}_1'\|_2 \|\mathbf{u}_2\|_2, \qquad (91)$$

$$0 \le \sigma(\langle \tilde{\mathbf{w}}_{2,r}, \boldsymbol{\xi} \rangle) \le C_3 \sigma_p^2, \tag{92}$$

$$\sigma(\langle \tilde{\mathbf{w}}_{1,r}, \mathbf{u}_2' \rangle) = 0,$$

$$\sigma(\langle \tilde{\mathbf{w}}_{2,r}, \mathbf{u}_2' \rangle) = C_1 \cos \theta \, \|\mathbf{u}_2\|_2 \, \|\mathbf{u}_2'\|_2 \,. \tag{94}$$

(93)

Using the above inequalities (equalities), we have

$$\mathcal{L}_{\mathcal{D}_{2}}(\tilde{\mathbf{W}}) \leq -\frac{1}{2} \ln \left(\frac{\exp(C_{1} \cos \theta \| \mathbf{u}_{1}' \|_{2} \| \mathbf{u}_{1} \|_{2})}{\exp(C_{1} \cos \theta \| \mathbf{u}_{1}' \|_{2} \| \mathbf{u}_{1} \|_{2}) + \exp(C_{1} \sin \theta \| \mathbf{u}_{1}' \|_{2} \| \mathbf{u}_{2} \|_{2} + C_{3} \sigma_{p}^{2})} \right) - \frac{1}{2} \ln \left(\frac{\exp(C_{1} \cos \theta \| \mathbf{u}_{2}' \|_{2} \| \mathbf{u}_{2} \|_{2})}{\exp(C_{1} \cos \theta \| \mathbf{u}_{2}' \|_{2} \| \mathbf{u}_{2} \|_{2}) + \exp(C_{3} \sigma_{p}^{2})} \right).$$
(95)

Based on Theorem 3.5 and $\|\mathbf{u}_1\|_2 = \|\mathbf{u}_1\|_2 = \|\mathbf{u}_1'\|_2 = \|\mathbf{u}_2'\|_2$, we have

$$\mathcal{L}_{\mathcal{D}_{\mathrm{ft}}}(\mathbf{W}^{(T)}) \leq \exp\left(-\Omega\left(\frac{\Lambda_{i} \|\mathbf{u}_{i}\|_{2}^{2}}{\sqrt{m}}T\right)\right) \cdot \tilde{L} + \mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{mn}\Lambda_{i}}\right) + \mathcal{O}\left(\frac{\sqrt{md}\sigma_{n}}{\Lambda_{i} \|\mathbf{u}_{i}\|_{2}}\right),$$
(96)

This completes the proof.

EXPERIMENTS COMPUTE RESOURCES Ι

The simulations are conducted on a commodity machine with Intel® Core i7-9700 CPU with a NVidia® 3090Ti GPU.

J **BROADER IMPACTS**

Our work studies side effects in DP-SGD trained models. One important side effect is unfairness. Our work identifies data imbalance as one source of unfairness, indicating that collecting balanced data is significant for maintaining fairness. Our work also shows the potential of design algorithms to adjust group weights to improve model fairness.