

CHARACTERIZING THE TRAINING DYNAMICS OF PRIVATE FINE-TUNING WITH LANGEVIN DIFFUSION

Anonymous authors

Paper under double-blind review

ABSTRACT

We show that differentially private full fine-tuning (DP-FFT) can distort pre-trained backbone features based on both theoretical and empirical results. We identify the cause of the distortion as the misalignment between the pre-trained backbone and the randomly initialized linear head. We prove that a sequential fine-tuning strategy can mitigate the feature distortion: first-linear-probing-then-fine-tuning (DP-LP-FFT). A new approximation scheme allows us to derive approximate upper and lower bounds on the training loss of DP-LP and DP-FFT, in a simple but canonical setting of 2-layer neural networks with ReLU activation. Experiments on real-world datasets and architectures are consistent with our theoretical insights. We also derive new upper bounds for 2-layer linear networks without the approximation. Moreover, our theory suggests a trade-off of privacy budget allocation in multi-phase fine-tuning methods like DP-LP-FFT.

1 INTRODUCTION

Today, many differentially-private (DP) machine learning pipelines proceed in two phases: (1) A model is pre-trained (non-privately) on a public dataset. (2) The model is then fine-tuned on private data, using DP optimization techniques such as DP stochastic gradient descent (DP-SGD) and its variants (Hoory et al., 2021; De et al., 2022; Tang et al., 2023; Zhang et al., 2024b). Pre-training a backbone model on public data enables differentially private fine-tuning to achieve improved performance across various downstream tasks (Yu et al., 2022) and is proven to be necessary in some cases (Ganesh et al., 2023a).

Despite these advances, the effect of DP on fine-tuning training dynamics remains poorly understood. Several key questions are yet to be answered: (1) how does randomness (both of initialization and DP optimization) impact the pre-trained representations? (2) What are the convergence rates of common fine-tuning methods, such as DP full fine-tuning (DP-FFT) and DP linear probing (DP-LP, where feature representations are frozen, and only the linear head is fine-tuned)? (3) Prior work suggests that combining an early stage of DP-LP with a later stage of DP-FFT yields better privacy-utility tradeoffs (Tang et al., 2023), yet there is no theoretical understanding of this phenomenon, nor is it clear how to optimally combine these fine-tuning methods.

Answering these questions theoretically requires an analysis that can capture the fine-grained optimization dynamics of DP fine-tuning. We seek a model of DP finetuning that satisfies 2 properties.

1. **Architecture-sensitivity:** The convergence dynamics must differentiate between representation learning in the backbone and learning in the linear head. The analyses of Bassily et al. (2014), Wang et al. (2022), Fang et al. (2023), Ganesh et al. (2023b) focus only on the network’s dimension, failing to capture this distinction.
2. **Ability to model nonlinearities:** The model should account for the nonlinearities introduced by multi neural layers, unlike existing methods that simplify analysis by linearizing neural networks (Ye et al., 2023a; Wang et al., 2024).

We propose a novel approximation of DP-SGD training dynamics based on linearizing Langevin diffusion around *the noise term*. This approach offers new insights into DP fine-tuning and significantly simplifies analysis by converting stochastic differential equations into ordinary differential equations (ODEs). We validate our theoretical predictions with real experiments.

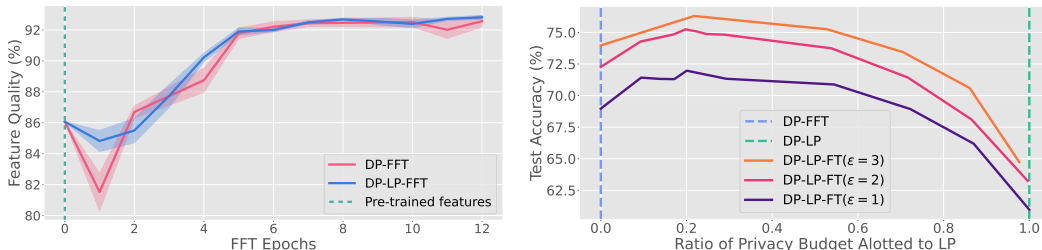


Figure 1: **Left:** Backbone feature quality evaluated by top-1 kNN accuracy on the downstream task, for ResNet-50, through public pre-training on ImageNet-1K and differentially private fine-tuning on STL-10. **Right:** Privacy budget trade-off in DP-LP-FFT, predicted in our theory, for WideResNet-16-4 on CIFAR-10 (Tang et al., 2023).

Main contributions. In summary, our key contributions are:

1. **New approximation technique:** In Section 2, we derive a first-order ODE via an asymptotic expansion of the stochastic noise in Langevin diffusion. Unlike previous methods, which linearize neural network parameters, our technique preserves the multi-layer structure of deep learning models while simplifying the analysis. This approach, commonly used in physics and control theory (Skorokhod et al., 2002), is novel in the context of private machine learning and bridges the gap between non-private neural network theory and the private regime.
2. **Understanding of feature distortion:** In Section 3, we provide a theoretical understanding of how DP fine-tuning affects feature representations. Using our approximation, we prove that, in 2-layer ReLU networks, randomly initialized linear heads distort pre-trained backbone features in the early stages of DP-FFT. Empirically Figure 1 demonstrates that feature quality evaluated on private data initially degrades during DP-FFT but later improves and surpasses pre-fine-tuning quality. Our theory also predicts that running a single epoch of DP-LP before transitioning to DP-FFT can mitigate this initial feature distortion, as shown empirically in the DP-LP-FFT curve of Figure 1 (left). This insight extends the findings of Kumar et al. (2022), who showed that LP-FFT reduces feature distortion in non-private, OOD scenarios, to in-distribution settings for both DP and non-DP cases.
3. **Theoretical convergence bounds:** In Section 4, we present new upper and lower bounds on the training loss of DP-LP and DP-FFT for 2-layer ReLU networks using our approximation technique. We also prove upper bounds for 2-layer linear networks without the approximation. To the best of our knowledge, this is the first convergence analysis of DP-SGD on non-linear neural network architectures.
4. **Mitigating feature distortion by combining fine-tuning methods:** Prior work by Tang et al. (2023) empirically showed that combining DP-LP and DP-FFT (DP-LP-FFT) can achieve better test accuracy than either method alone. In Figure 1b, we demonstrate that allocating approximately 20% of the privacy budget to DP-LP yields optimal test accuracy. In Section 5, we provide a partial theoretical explanation for this phenomenon. Specifically, our bounds suggest that DP-FFT may underperform relative to DP-LP at lower privacy budgets, while DP-LP-FFT can outperform both methods under moderate privacy budgets. These predictions are empirically verified across various architectures and benchmarks in Section 5.3.

1.1 RELATED WORK

Similar empirical phenomena have been explored in non-private, out-of-distribution (OOD) contexts by Aghajanyan et al. (2021), Kumar et al. (2022), Trivedi et al. (2023), and Chen et al. (2024). Kumar et al. (2022) demonstrated that non-DP fine-tuning distorts pre-trained features, leading to degraded OOD performance. But their theory relies on the assumption that OOD test data exists in an orthogonal subspace to the fine-tuning training data, leaving their results unable to explain why, in many transfer learning tasks, linear-probe fine-tuning (LP-FFT) still outperforms both LP and full fine-tuning (FFT) in in-distribution (ID) settings. Our work seeks to fill this research gap.

108 Wang et al. (2024) examined how pre-trained representations enhance DP fine-tuning within the
 109 neural collapse framework, though their analysis was restricted to the final layer. Meanwhile, Tang
 110 et al. (2023) empirically observed the privacy budget trade-off for WideResNet models pre-trained
 111 on synthetic data, but without accompanying theoretical insights.

112 Analyses by Wang et al. (2019), Chen et al. (2020a), Ganesh et al. (2023b), and Fang et al. (2023)
 113 rely on standard convexity/non-convexity and smoothness assumptions, which abstract away the
 114 simultaneous dynamics between the backbone and linear head. Other works (Ye et al., 2023b; Wang
 115 et al., 2024) focus on linearized models, limiting their ability to capture the nuanced interactions
 116 between these components. Our explanation of representation alignment builds on the theoretical
 117 foundation of Min et al. (2024), which we extend to a DP context using novel approximation tools.

119 2 CONTINUOUS MODELING OF DIFFERENTIALLY PRIVATE FINE-TUNING

121 **Notation.** We use ∂ to denote both the deterministic and stochastic differential operators. The
 122 dot product between vectors x, y is $x^\top y$, the Euclidean norm of vector x is $\|x\|_2$, and the infinity
 123 norm is $\|x\|_\infty$. The trace of a matrix is denoted by tr , and the ReLU activation is ϕ . For any
 124 twice differentiable function $f(x)$, its gradient is denoted $\nabla_x f$ and its Hessian as $H_x f$. \sqcup denotes
 125 the disjoint union. $[i] := \{1, \dots, i\}$. The cosine similarity between two vectors u, v is defined as
 126 $\cos(u, v) = \frac{u^\top v}{\|u\|_2 \|v\|_2}$. We denote the privacy cost estimated by Rényi divergence as r .

128 **DP-SGD Dynamics.** Differential privacy (DP) is a widely used framework for evaluating privacy
 129 leakage in a dataset accessed through queries (Dwork & Roth, 2014). In machine learning, DP
 130 ensures that an adversary cannot confidently determine whether specific training samples are part
 131 of the dataset. **Differentially Private Stochastic Gradient Descent** (DP-SGD), introduced by Abadi
 132 et al. (2016), is the standard algorithm for training deep neural networks while maintaining privacy.

133 Our fine-tuning theory is built on an analysis of DP-SGD dynamics. Although real-world algo-
 134 rithms are discrete, continuous approximations—such as stochastic differential equations (SDE)
 135 like Langevin diffusion—are often used to study these dynamics (Chourasia et al., 2021; Ye et al.,
 136 2023b). In a similar vein, Kumar et al. (2022) use gradient flow, a continuous approximation of
 137 SGD, to study fine-tuning in a non-private context.

138 **Definition 2.1** (Langevin diffusion (Ganesh et al., 2023b)). Langevin diffusion is an SDE that
 139 models the dynamics of a system influenced by both deterministic and random forces (Lemons
 140 & Gythiel, 1997). For DP-SGD, we define an p -dimensional Langevin diffusion as follows:

$$141 \quad \partial\theta = -\nabla_\theta \mathcal{L}(\theta|f) \partial t + \sqrt{2\sigma^2} \partial Q_t, \quad (1)$$

142 where $\theta \in \mathbb{R}^p$ represents the neural network parameters, f is the network architecture, $\mathcal{L}(\cdot|f) : \mathbb{R}^p \rightarrow \mathbb{R}$
 143 is the training loss, and $\sigma > 0$ is the noise multiplier (Abadi et al., 2016). $\{Q_t\}_{t \geq 0}$ is the
 144 standard Brownian motion in \mathbb{R}^m modeling the Gaussian noise mechanism.

146 By Itô’s lemma (Ito, 1951), the Langevin diffusion of the training loss is given by

$$147 \quad \partial\mathcal{L} = [-\|\nabla_\theta \mathcal{L}(\theta|f)\|_2^2 + \sigma^2 \text{tr}(H_\theta \mathcal{L})] \partial t + \sqrt{2\sigma^2} (\nabla_\theta \mathcal{L}(\theta|f))^\top \partial Q_t. \quad (2)$$

149 Ye et al. (2023b) study how random initialization affects DP-SGD performance in linearized neural
 150 networks via Langevin diffusion. To facilitate theoretical analysis, they linearize the entire neural
 151 network using 1st-order Taylor expansions at the initial parameter θ_0 .

$$152 \quad f(x) \approx f_{\text{lin}}(x) := f(x) \Big|_{\theta=\theta_0} + \frac{\partial f(x)}{\partial \theta} \Big|_{\theta=\theta_0} \cdot (\theta - \theta_0). \quad (3)$$

155 Recently, this linearization technique has gained popularity for explaining key deep learning phe-
 156 nomena (Ortiz-Jimenez et al., 2021). However, fully linearizing the model removes critical multi-
 157 layer interactions, making this approach unsuitable for our analysis.

158 To address this, we treat the optimization trajectory as a dynamical system and consider noise in
 159 gradient updates as random perturbations. Applying a zeroth-order asymptotic expansion of Equa-
 160 tion (1) with respect to the noise multiplier σ (Freidlin et al., 2012), we approximate:

$$161 \quad \partial\theta \approx \tilde{\partial}\theta = -\nabla \mathcal{L}(\tilde{\theta}|f) \partial t. \quad (4)$$

This zeroth-order expansion simplifies the analysis of complex, stochastic, and non-linear equations. By substituting the approximate parameter $\tilde{\theta}$ into Equation (2), our modeling preserves the noisy behavior characteristic of DP-SGD.

3 REPRESENTATION ALIGNMENT

In this section, we introduce the concept of representation alignment, present our theoretical findings, and validate them with experiments. Representation alignment refers to the process by which the classification head aligns itself with the pre-trained backbone features. During the DP-FFT process, this alignment creates a characteristic trend in feature quality: initially, the randomly initialized linear head distorts the pre-trained features, but as it better aligns with the backbone, the distortion diminishes, and the overall quality of the backbone features improves over time.

3.1 THEORY

Our goal is to understand (1) how does DP fine-tuning distort the pre-trained features in the backbone, and (2) under what conditions this distortion can be mitigated. We consider the simple binary classification setup from Min et al. (2024), which provides a clear and intuitive understanding of representation alignment. The results generalize to our experiments in Section 3.2. Specifically, we use a 2-layer fully-connected neural network with h hidden nodes and ReLU activation ϕ ,

$$f(x) = v^\top g(x) = v^\top \phi(W^\top x) = \sum_{j=1}^h v_j \phi(w_j^\top x). \quad (5)$$

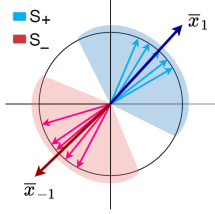


Figure 2: Visualization of Assumption 3.1.

fine-tuning on a dataset $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^n$ with n inputs $x_i \in \mathbb{R}^{d_x}$, and binary labels $y_i \in \{-1, 1\}$. The objective is to minimize the training loss $\mathcal{L}(\tilde{\theta}|f) := \sum_{i=1}^n \ell(y_i, f(x_i))$, using the exponential loss $\ell(y, \hat{y}) := \exp(-y\hat{y})$. Similar results hold for logistic loss (Min et al., 2024). For simplicity, we make the following assumption.

Assumption 3.1 (Data correlation (Min et al., 2024)). For any pair of data $(x_i, y_i), (x_j, y_j)$, the inputs are positively/negatively correlated if the labels are the same/different.

$$\inf_{i, j \in [n]} \left[(y_1 y_2) \cdot \frac{x_1^\top x_2}{\|x_1\|_2 \|x_2\|_2} \right] := \mu > 0. \quad (6)$$

We define two cones in \mathbb{R}^{d_x} that separate subspaces spanned by data points in the positive and negative classes, respectively: $S_+ = \{z \in \mathbb{R}^{d_x} : \forall i \in [n], \mathbb{I}_{x_i^\top z > 0} = \mathbb{I}_{y_i = 1}\}$, $S_- = \{z \in \mathbb{R}^{d_x} : \forall i \in [n], \mathbb{I}_{x_i^\top z > 0} = \mathbb{I}_{y_i = -1}\}$. Min et al. (2024) prove that $S_+ \cap S_- = \emptyset$, and $x_i \in S_{+/-}$ if $y_i = 1/-1$ (see Figure 2). We define the mean data directions of class $c \in \{-1, 1\}$ by $\bar{x}_c := \sum_{i \in [n]} x_i \cdot \mathbb{I}_{y_i = c}$.

We assume that a ‘‘clustering’’ behavior emerges in the pre-trained features, which allows the features to work well in transfer learning (Galanti et al., 2022). This phenomenon is well-documented in the neural collapse literature (Kothapalli, 2023), suggests that pre-trained features w_j tend to converge around the mean direction for data in class $c(j)$.

Assumption 3.2 (Collapsed neural features). For each w_j in Equation (5) where $j \in [h]$ (with h denoting the dimension of the linear head), it holds that $w_j \in S_+$ or $w_j \in S_-$. We define $c(j) = 1$ if $w_j \in S_+$, and $c(j) = -1$ if $w_j \in S_-$. Thus, there is a partition $[h] = F_+ \sqcup F_-$ over the index set $[h]$, such that for each w_j ,

$$\begin{cases} j \in F_+ & \text{if } w_j \in S_+, \\ j \in F_- & \text{if } w_j \in S_-. \end{cases} \quad (7)$$

Feature quality. Assumption 3.2 says that data with positive label (resp. negative) only activates the j -th neuron if $j \in F_+$ (resp. $j \in F_-$). As a result, any positive data pair, (x, y) and (x, y') with $y = y'$, activate the same set of neurons. From a contrastive learning viewpoint, it makes the representations of them semantically similar (Saunshi et al., 2019). Namely, when the features w_j

and data inputs x_i are normalized unit vectors, the difference between representations of a positive data pair is bounded by:

$$\|g(x) - g(x')\|_\infty \leq \max_{y_i=c(j)=y} \cos(w_j, x_i), \quad (8)$$

which represents the maximum cosine similarity between the features w_j and the data points.

However, FFT or DP-FFT with random initialization may reduce the feature quality.

Theorem 3.3 (Random initialization causes feature distortion). *If Assumption 3.1 and Assumption 3.2 hold, and the linear head is randomly initialized by $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$, then with probability $1 - 2^{-h}$, $\forall \beta > 0, \exists j \in [h], \Delta t > 0$ such that during the time interval $(0, \Delta t)$, DP-FFT distorts w_j reducing its alignment with the data cluster. The cosine similarity between w_j and the data cluster mean $\bar{x}_{c(j)}$ decreases monotonically:*

$$\left. \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right|_t < 0, \quad \forall t \in (0, \Delta t) \quad (9)$$

For a pre-trained w_j that aligns with $c(j)$ -labeled data, DP-FFT (as modeled by Equation (4)) makes it deviate from $\bar{x}_{c(j)}$, the mean direction of those data. w_j is optimal when $\cos(w_j, \bar{x}_{c(j)}) = 1$. This result holds for both DP and non-DP settings and explains the potential feature distortion observed in in-distribution and non-private settings, such as those studied by Kumar et al. (2022)). The stochastic analysis of non-smooth loss, activation, cosine similarity functions is challenging without our approximation.

Next, we show that running (DP-)LP before (DP-)FFT could mitigate feature distortion.

Theorem 3.4 (DP-LP first mitigates feature distortion). *Suppose Assumption 3.1 and Assumption 3.2 hold, and the linear head is randomly initialized by $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ for any $\beta > 0$. There exists $\Delta t > 0$ such that after running DP-LP for time Δt , switching to full fine-tuning ensures that DP-FFT does not distort the pre-trained features. Specifically, $\cos(w_j, \bar{x}_{c(j)})$ is non-decreasing for all $j \in [h]$:*

$$\left. \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right|_t \geq 0, \quad \forall t \in (\Delta t, +\infty) \quad (10)$$

See complete proofs of Theorem 3.3 and Theorem 3.4 in Appendix C.1.

3.2 EXPERIMENTS

In this section, we show empirical evidence supporting Theorems 3.3 and 3.4.

Pre-training and Model. We pre-train Vision Transformers (ViT) and ResNet-50 backbones on ImageNet-1K using Self-Supervised Learning methods, including BYOL (Grill et al., 2020) and MoCo v2 (Chen et al., 2020b), as well as distillation methods (Touvron et al., 2021). Then we fine-tune the backbone with a linear classification head on CIFAR-10 and STL-10 using DP-SGD.

Experiment protocols. We conduct public pre-training for 100 epochs with a batch size of 256. Following this, we implement DP-SGD using the pre-trained weights and a randomly initialized linear head for 30 epochs. Each DP fine-tuning process is repeated with 5 random seeds and a batch size of 1000. We evaluate the backbone features on both the pre-training and fine-tuning datasets, measuring feature quality through top-1 kNN accuracy (Chen et al., 2023).

Private fine-tuning initially distorts features. Figure 3 qualitatively visualizes the effect of DP-FFT on feature quality with respect to the private test data. We pre-train (BYOL) a ResNet-50 backbone on ImageNet-1K and DP fine-tune (DP-SGD, $\epsilon = 1$) it on STL-10. We qualitatively assess the features of the private test data within the ResNet-50 backbone by visualizing the backbone mappings (outputs from the penultimate layer) of data points using UMAP (McInnes et al., 2020). For simplicity, we only plot 3 classes in CIFAR-10.

Figure 3 indicates that during the initial phases of DP-FFT, the randomly initialized linear head interferes with the pre-trained features in the backbone network, leading to a degradation in feature quality on both the pre-training and fine-tuning datasets. This observation validates Theorem 3.3. Concurrently, the linear head begins adapting to these pre-trained features, a process we refer to as

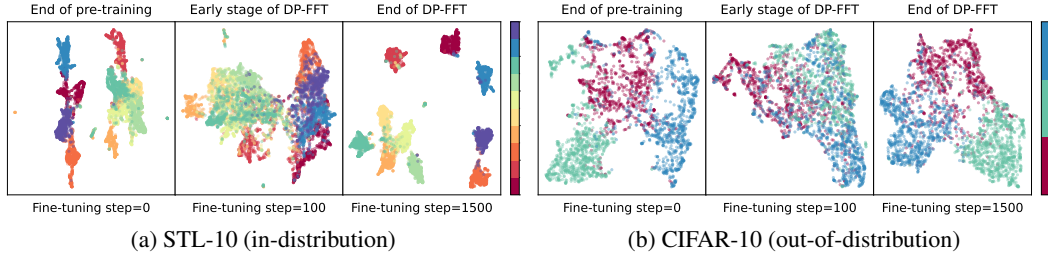


Figure 3: We pre-train (BYOL) a ResNet-50 backbone on ImageNet-1K and DP fine-tune (DP-SGD, $\epsilon = 1$) it on STL-10. We qualitatively evaluate the features in the ResNet-50 backbone by visualizing the backbone mappings (penultimate layer outputs) of data points via UMAP (McInnes et al., 2020). These results suggest that DP-FFT distorts feature quality before improving it, as predicted by Theorem 3.3.

“**representation alignment.**” As this alignment progresses, the backbone starts to regain a portion of its original feature quality, which had been degraded by DP noise and shifts in data distribution.

Linear probing mitigates feature distortion. To illustrate the benefits of linear probing, we first run DP-LP for 1 epoch before transitioning to DP-FFT for the remaining epochs. In the initial steps of DP-FFT, the feature distortion is significantly weaker (Figure 1a) if we first run DP-LP. This supports the claim of Theorem 3.4.

We also evaluate features on the pre-training domain (see Figure 5).

4 DP FINE-TUNING CONVERGENCE RATES

Section 3 showed that DP-LP-FFT can mitigate feature distortion. A natural question is, for a fixed privacy budget, how do DP-LP and DP-FFT affect the convergence of fine-tuning loss function? We study this question under two models: (1) our zeroth-order approximation of Langevin diffusion (Section 4.1), and (2) a two-layer neural network without our zeroth-order approximation (Section 4.1.1). The second result will be used to study the budget allocation of DP-LP-FFT in Section 5. To our knowledge, these are the first convergence guarantees (approximate or not) for DP fine-tuning on explicit nonlinear neural network architectures.

Privacy guarantees We begin by establishing the privacy guarantees of Langevin diffusion by bounding the Rényi divergence of its trajectory distributions on neighboring datasets (Mironov, 2017). Both Ganesh et al. (2023b) and Ye et al. (2023b) show that the Rényi divergence increases linearly over time. We use this guarantee for all fine-tuning variants.

Theorem 4.1 (Rényi privacy guarantee (Ganesh et al., 2023b)). *Suppose we initialize a pair of neural network parameters θ, θ' by some i.i.d. distributions Θ_0, Θ'_0 . We fine-tune θ, θ' respectively on neighboring datasets $\mathcal{D}, \mathcal{D}'$ via Langevin diffusion. Denote the distribution of the trajectory of θ by $\Theta_{[0,T]}$ over $[0, T]$. Similarly, denote the trajectory distribution of θ' by $\Theta'_{[0,T]}$. Then for any $\alpha \geq 1$, the Rényi divergence R_α is bounded linearly in time,*

$$r := R_\alpha(\Theta_{[0,T]} \parallel \Theta'_{[0,T]}) = O\left(\frac{\alpha \Delta_g T}{\sigma^2}\right) \quad (11)$$

where σ is the noise multiplier, and $\Delta_g \geq \|\nabla \mathcal{L}(\theta; \mathcal{D}) - \nabla \mathcal{L}(\theta; \mathcal{D}')\|$ is the upper bound of gradient difference between neighboring datasets. Thus, for any $\delta \in (0, 1)$, the Langevin diffusion satisfies

$$\left(\frac{\alpha \Delta_g T}{4\sigma^2} + \frac{\log(1/\delta)}{\alpha - 1}, \delta\right) - \text{differential privacy.} \quad (12)$$

4.1 CONVERGENCE RATES UNDER THE ZERO-ORDER APPROXIMATION

We follow the approximation scheme outlined in Equation (4) to derive convergence results for two-layer ReLU neural networks.

Theorem 4.2 (Approximate DP-LP loss convergence). *If Assumption 3.1 and Assumption 3.2 hold at $t = 0$, we can bound the loss after running DP-LP for $t = T$:*

$$\frac{1}{\frac{1}{\mathcal{L}_c(0)}e^{-B_1T} + \frac{A_1}{B_1}(1 - e^{-B_1T})} \leq \mathcal{L}_c(T) \leq \frac{1}{\frac{1}{\mathcal{L}_c(0)}e^{-B_2T} + \frac{A_2}{B_2}(1 - e^{-B_2T})} \quad (13)$$

where $\mathcal{L}_c(t)$ denotes the training loss of data points labeled $c \in \{-1, 1\}$, $\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2$, and

$$\begin{cases} A_1 = \sum_{w_j \in S_c} [\max_{y_i=c} w_j^\top x_i]^2 \\ B_1 = \frac{1}{2}\sigma^2 \left\{ \sum_{y_i=c} \|\text{relu}(W^\top x_i)\|_2^{-2} \right\}^{-1} \\ A_2 = \sum_{w_j \in S_c} [\min_{y_i=c} w_j^\top x_i]^2 \\ B_2 = \frac{1}{2}\sigma^2 \left\{ \sum_{y_i=c} \|\text{relu}(W^\top x_i)\|_2^4 \right\}^{1/2} \end{cases} \quad (14)$$

are constants for DP-LP.

When we set $n = h = 2$, $y_1 = -y_2$, $w_1 = x_1 = -w_2 = -x_2$, the upper and lower bounds are equal and we achieve a tight bound on the DP-LP loss.

Theorem 4.3 (Approximate DP-FFT loss convergence). *For simplicity, we assume that $\|x_i\|_2 = R$ for all $i \in [n]$. If Assumption 3.1 and Assumption 3.2 hold, and we consider a balanced initialization $\|W\|_F^2 = \|v_0\|_2^2$ (Min et al., 2023a) at $t = 0$, then*

(i) we lower bound the loss after running DP-FFT for $T > 0$:

$$\mathcal{L}_c(T) \geq \frac{1}{\frac{1}{\mathcal{L}_c(0)}e^{(1-\exp(\lambda_c T))A_l C_l / \lambda_c} + \frac{B_l}{C_l} [1 - e^{(1-\exp(\lambda_c T))A_l C_l / \lambda_c}]} \quad (15)$$

where we define $A_l = \|W_0\|_F^2$, $B_l = 2R^2$, $C_l = \frac{R^2\sigma^2(1+\mu^2)}{2}$ and $\lambda_c = 2R\mathcal{L}_c(0)$.

(ii) we upper bound the loss after running DP-FFT for $T > 0$:

$$\mathcal{L}_c(T) \leq \frac{1}{\frac{B_u}{C_u}(1 - e^{-A_c C_u T}) + \frac{1}{\mathcal{L}_c(0)}e^{-A_c C_u T}} \quad (16)$$

where we define $A_c = \sum_{w_j \in S_c} [v_{j,t=0}^2 + \|w_j\|_2^2]$, $B_u = R^2\mu^2$ and $C_u = \frac{1}{2}R^2\sigma^2$.

4.1.1 THEORY WITHOUT THE ZERO-ORDER APPROXIMATION (2-LAYER LINEAR NETWORK)

We complement the results in Section 4.1 by removing the zeroth-order approximation in a simpler setup: 2-layer linear networks for a regression task. We define a linear network by replacing the ReLU activation ϕ with an identity function in Equation (5). We collect the data inputs in a matrix $X \in \mathbb{R}^{n \times d_x}$ and put the labels in a vector $Y \in \mathbb{R}^n$. For simplicity, we assume that $n \geq d$ and $X^T X = I_{d_x \times d_x}$. We consider the MSE training loss $\mathcal{L}(v, W) := \frac{1}{2} \sum_{i \in [n]} (v^\top W^\top x_i - y_i)^2 = \frac{1}{2} \|XWv - Y\|_2^2$.

Note that the loss function is nonconvex in the parameters being fine-tuned, so the gradient descent training becomes a nonlinear dynamical system. This significantly complicates theoretical analysis. Prior works have dealt with the challenging analysis by using heavy approximations (Bu et al., 2023; Ye et al., 2023b). We overcome these theoretical difficulties by using conservation laws and geometric properties of Langevin dynamics (see Appendix for more detail).

Pretrained features. We evaluate a backbone W by the least square error:

$$\gamma(W) := \inf_{u \in \mathbb{R}^h} \mathcal{L}(u, W) = Y^T (I_{n \times n} - XW(XW)^\dagger)Y. \quad (17)$$

where $(\cdot)^\dagger$ denotes the pseudo inverse of a matrix. This metric measures the optimal loss for LP when fixing the current features. $\gamma = \gamma(W_0)$ denotes the initial least square error. We suppose W_0 has orthonormal columns, following prior works (Tripuraneni et al., 2020; Kumar et al., 2022).

Theorem 4.4 (DP-LP loss convergence). *If we randomly initialize the linear head $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ and we run linear probing for time T , then*

$$\mathbb{E}[\mathcal{L}(T)] \leq \frac{1}{2}(h\beta + \|Y\|^2)e^{-T} + (\gamma + h\sigma^2)(1 - e^{-T}) \quad (18)$$

In this theorem, the first term describes that the loss tends to exponentially decrease, while the second term describes the limiting behavior induced by linear probing and the added noise.

Theorem 4.5 (DP-FFT loss convergence). *If $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ and Assumption E.7 holds, and we run fine-tuning (Equation (96)) for time T , then the loss converges:*

$$\mathbb{E}[\mathcal{L}(T)] \leq \frac{1}{2}(h\beta + \|Y\|_2^2)e^{-AT} + L^\square(1 - e^{-AT}) \quad (19)$$

$$\text{where } \begin{cases} A = h\beta - 1 - \sqrt{2}\sigma^2(1 + d_x) > 0 \\ L^\square = \sigma^2 \frac{(1+d_x)\|X^T Y\|_2 + d_x}{A} \end{cases} .$$

This upper bound has a similar form to Equation (18) while the factor A of the exponential terms depends on the initialization and the noise. When we take limit $\sigma \rightarrow 0$ in Theorem 4.4 and 4.5, the Langevin diffusion degenerates to a gradient flow and the loss converges exponentially to zero as $T \rightarrow \infty$. This recovers known results from the non-private optimization literature (Min et al., 2023a).

The bounds in Section 4.1 and Section 4.1.1 exhibit different dependencies on the hidden dimension h and the data dimension d_x due to the differing curvature properties of the loss functions in each setup. The underlying reason is that the noise term introduced by Itô’s formula (Equation (2)) is influenced by the curvature of the loss function. While the square function has constant curvature, the exponential function does not, leading to varying noise impacts.

5 BUDGET ALLOCATION BETWEEN DP-LP AND DP-FFT

Finally, we consider the DP-LP-FFT fine-tuning strategy, which first applies DP-LP for some portion r of the privacy budget (i.e. for some number of training iterations), then uses the remaining privacy budget for DP-FFT. In this section, we ask: given a fixed privacy budget, how should we allocate it across DP-LP and DP-FFT? Our results, both theoretical and empirical, suggest that at low total privacy budget, one should allocate more of the total privacy budget to DP-LP.

5.1 RESULTS UNDER ZERO-ORDER APPROXIMATION

We first show how to allocate privacy budget to avoid the feature distortion analyzed in Section 3, using the zeroth-order approximation.

Theorem 5.1 (Estimated privacy budget allocated to DP-LP). *If Assumption 3.1 and Assumption 3.2 hold at $t = 0$, then for any $\rho \in (0, 1)$, with probability $(1 - \rho)^h$, we can avoid feature distortion by spending*

$$r \propto \sigma^4 \sqrt{\ln(2/\rho)} \quad (20)$$

amount r of privacy budget on DP-LP, where σ is the noise multiplier. That is, we ensure that $\forall j \in [h]$, and any $t > 0$ after DP-LP,

$$\left. \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right|_t \geq 0 \quad (21)$$

According to Theorem 5.1, a greater proportion of the privacy budget should be allocated to DP-LP when the total privacy budget is smaller.

5.2 RESULTS WITHOUT APPROXIMATION (2-LAYER LINEAR NETWORK)

Complementing the result of Section 5.1, we use the 2-layer linear model of Section 4.1.1 to show that DP-LP-FFT may work better in some settings than linear probing or full fine-tuning alone.

Linear probing first can accelerate fine-tuning by aligning the linear head. The following result provides a convergence bound for DP-LP-FFT when we linear-probe for time t_{lp} , and then fully fine-tune for time t .

Proposition 5.2 (Convergence of DP-LP-FFT). *Suppose we randomly initialize the linear head $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ and Assumption E.7 hold. We run linear probing for time t_{lp} and then fine-tuning (Equation equation 96) for time t , then the loss is upper bounded by:*

$$\mathbb{E}[\mathcal{L}(t)] \leq \mathbb{E}[\mathcal{L}_{lp}]e^{-At} + L^\square(1 - e^{-At}) \quad (22)$$

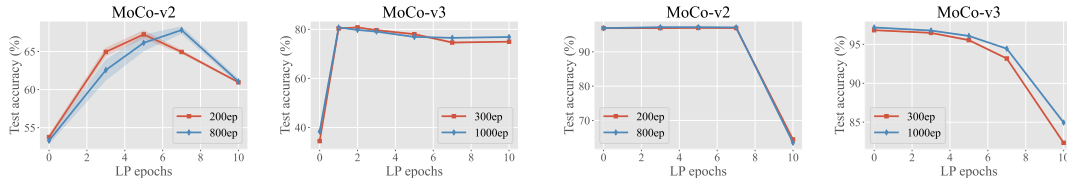
where \mathcal{L}_{lp} is the expected loss after linear probing, $A = h\beta - 1 - \sqrt{2}\sigma^2(1 + d_x)$, and $L^\square = \sigma^2 \frac{(1+d_x)\|X^T Y\|_2 + d_x}{A}$. The coefficient $A = \mathbb{E}[\lambda_{\max}(D)] > 0$ increases as t_{lp} increases when we run linear probing in a finite time interval $t_{lp} < \ln \left[3 + \frac{h(\sigma^2 - \beta)}{\|W_0^\top X^T Y\|_2^2} \right]$.

Corollary 5.3. *Suppose we randomly initialize the linear head $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ and Assumption E.7 hold. Then the two-phase method, first-linear-probing-then-finetuning (LP-FFT), could achieve a tighter loss upper bound than linear probing or fine-tuning in expectation if we first run linear probing for $t_{lp} < \ln \left[3 + \frac{h(\sigma^2 - \beta)}{\|W_0^\top X^T Y\|_2^2} \right]$.*

Corollary 5.3 suggests that when we fix other hyperparameters (e.g. the total training time T), the performance of LP-FFT depends on the noise scale σ . If σ is large enough such that $T < \ln \left[3 + \frac{k(\sigma^2 - \beta)}{\|B_0 X^T Y\|_2^2} \right]$, then LP may be the best; if σ is small enough such that $\ln \left[3 + \frac{k(\sigma^2 - \beta)}{\|B_0 X^T Y\|_2^2} \right] \leq 0$, then FT may be the best; LP-FT could achieve the best performance when the noise scale is in a proper interval $\sigma^2 \in \left(\beta - 2 \frac{\|B_0 X^T Y\|_2^2}{k}, \beta + (e^T - 3) \frac{\|B_0 X^T Y\|_2^2}{k} \right)$.

In our theory without approximation, these predictions are based only on upper bounds, so we cannot conclusively say that any fine-tuning approach outperforms another. Nonetheless, our theoretical results in two approaches suggest that the smaller the total budget, the more privacy budget should be allotted to DP-LP.

5.3 EXPERIMENTS



(a) Private utility curves ($\sigma = 0.3$)

(b) Non-private utility curves

Figure 4: Utility curves for pretraining on ImageNet-1K and fine-tuning on CIFAR-10 over ResNet-50, with pretrained features from MoCo-v2 and MoCo-v3 (Chen et al., 2020b; Chen* et al., 2021). We compare the performance from pre-trained weights of different pre-training epochs (200/800 epochs for MoCo-v2, 300/1k epochs for MoCo-v3). The x-axis sweeps the number of LP epochs from 0 to 10; the remaining epochs (out of 10) use FFT.

To illustrate the privacy budget trade-off, we empirically evaluate the benefits of DP-LP-FFT on real data and architectures.

DP-LP-FFT outperforms other fine-tuning methods: Pre-training on synthetic data. We follow the setup in Tang et al. (2023) and generate utility curves for $\epsilon = 1, 2, 3$ (Figure 1b). We pre-train WideResNet with synthetic images generated from StyleGAN-oriented (Baradad et al., 2021), and fine-tune it with DP-SGD on CIFAR-10. The x-axis sweeps the fraction of privacy budget allocated to DP-LP, and the remaining budget is used for DP-FFT. We find that at various privacy levels, DP-LP-FFT gives a clear advantage over either DP-FFT or DP-LP alone.

Figure 1b presents a different trend from our theoretical prediction, where we expect the optimal budget ratio for DP-LP to increase as the privacy noise grows. A possible intuitive explanation is

that, in the Figure 1b experiments, the pre-training data is synthetic, making it 'distant' from the CIFAR-10 fine-tuning data distribution. This divergence may violate our assumption that the pre-trained weights w_j are well-aligned with the fine-tuning data x_i .

DP-LP-FFT outperforms other fine-tuning methods: Pre-training on ImageNet-1K. Figure 4 illustrates the utility curves on ResNet-50 for $\sigma = 0, 0.3$.¹ To demonstrate utility curves for DP-LP-FFT, we vary the number of epochs of linear probing from $e_{LP} = 0$ to $e_{LP} = 10$; all remaining epochs (out of 10 total) are allocated to full fine-tuning, i.e., $e_{FFT} = 10 - e_{LP}$. Note that full fine-tuning corresponds to $e_{LP} = 0$ (the leftmost point of our subplots), and linear probing corresponds to $e_{LP} = 10$. We observe that for non-private optimization (Figure 4b), full fine-tuning achieves the highest test accuracy. However, for DP-SGD (Figure 4a), linear probing outperforms full fine-tuning, and DP-LP-FFT outperforms both DP-LP and DP-FFT.

Model	ResNet ₁₈			MobileNet _{v3}			Transformer _{DeiT}		
	∞	1.29	0.57	∞	1.29	0.57	∞	1.29	0.26
LP	68.54 _{0.02}	67.90 _{0.12}	66.60 _{0.04}	71.12 _{0.31}	69.54 _{0.08}	67.32 _{0.03}	95.74 _{0.04}	93.61 _{0.08}	94.21 _{0.08}
LP-FFT	72.66 _{0.12}	68.65 _{0.08}	59.79 _{1.03}	71.30 _{0.11}	71.18 _{0.06}	66.94 _{0.08}	96.82 _{0.08}	93.66 _{0.15}	93.62 _{0.05}
FFT	73.69 _{0.03}	59.79 _{1.03}	53.82 _{0.37}	77.02 _{0.31}	63.06 _{0.05}	45.12 _{0.07}	96.17 _{0.08}	90.31 _{0.53}	84.19 _{0.82}

Table 1: Test accuracies of DP-LP, DP-LP-FFT, and DP-FFT on various architectures.

Comparing DP fine-tuning methods. As suggested by Theorem 5.1 and Corollary 5.3, as the noise scale σ increases, the best fine-tuning strategy changes from DP-FFT (small σ , low privacy regime) to DP-LP-FFT, to DP-LP (large σ , high privacy regime). To qualitatively test this prediction, we sweep over different noise scales σ and fix other hyperparameters in each benchmark and model architecture. We sort the rows by the number of parameters of each model and the noise scale in an ascending order. For each experiment setting, we report average test accuracies with standard errors. As expected, among the three fine-tuning methods (Table 1), DP-FFT almost always does the best under small noise scales (including the non-private setting where $\sigma = 0$), DP-LP-FFT does the best under moderate noise scales, and DP-LP does the best under large noise scales. The close non-DP (ϵ) performance of FFT and LP-FFT on transformer architectures is consistent with previous observations in Kumar et al. (2022, Table 1). We also provide results with LoRA (see Table 2).

6 CONCLUSION AND DISCUSSION

We characterize the training dynamics of DP fine-tuning under a simplified theoretic setup (2-layer neural networks, separable datasets with -1/1 labels) using a Langevin diffusion-based approximation of DP-SGD, with an asymptotic expansion of random perturbations in dynamical systems as an approximation for Langevin diffusion. Our theory identifies and explains the phenomenon of representation distortion and alignment during DP fine-tuning, which we confirm empirically. Our work takes a step towards understanding how different private fine-tuning strategies can be mixed to improve performance, which could be useful for designing or mixing other strategies, such as memory-efficient zeroth-order optimization with differential privacy (Zhang et al., 2024a).

Limitations and open questions There are several open questions we cannot cover in this work, such as generalizing our results to multi-layer neural networks with our approximation technique, the effect of other loss functions on the fine-tuning dynamics, and loss lower bounds for DP-LP/FFT without the zeroth-order approximation. Moreover, it is unclear how to apply our theory to other fine-tuning methods like LoRA (Hu et al., 2022b), as well as generative models for which neural collapse does not happen. Understanding whether the zeroth-order approximation can facilitate analysis in these settings is an interesting and important question for future work.

Reproducibility Statement. We have included full proofs for all theoretical results and sufficient experimental details in appendices to reproduce our results. We will also release our code under a permissive open-source license upon acceptance.

¹The model performance is compromised because we replace the BatchNorm (Ioffe & Szegedy, 2015) in the pre-trained weights with GroupNorm (Wu & He, 2018). BatchNorm relies on batch statistics, which conflicts with the principles of differential privacy.

REFERENCES

- 540
541
542 Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar,
543 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM*
544 *SIGSAC Conference on Computer and Communications Security, CCS '16*, pp. 308–318, New
545 York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450341394. doi:
546 10.1145/2976749.2978318. URL <https://doi.org/10.1145/2976749.2978318>.
- 547 Armen Aghajanyan, Akshat Shrivastava, Ancht Gupta, Naman Goyal, Luke Zettlemoyer, and
548 Sonal Gupta. Better fine-tuning by reducing representational collapse. In *International Confer-*
549 *ence on Learning Representations, 2021*. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=OQ08SN70M1V)
550 [OQ08SN70M1V](https://openreview.net/forum?id=OQ08SN70M1V).
- 551 Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit
552 acceleration by overparameterization. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of*
553 *the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine*
554 *Learning Research*, pp. 244–253. PMLR, 10–15 Jul 2018. URL [https://proceedings.](https://proceedings.mlr.press/v80/arora18a.html)
555 [mlr.press/v80/arora18a.html](https://proceedings.mlr.press/v80/arora18a.html).
- 556 Manel Baradad, Jonas Wulff, Tongzhou Wang, Phillip Isola, and Antonio Torralba. Learning to see
557 by looking at noise. In *Advances in Neural Information Processing Systems, 2021*.
- 558 Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient
559 algorithms and tight error bounds. In *Proceedings of the 2014 IEEE 55th Annual Symposium on*
560 *Foundations of Computer Science, FOCS '14*, pp. 464–473, USA, 2014. IEEE Computer Society.
561 ISBN 9781479965175. doi: 10.1109/FOCS.2014.56. URL [https://doi.org/10.1109/](https://doi.org/10.1109/FOCS.2014.56)
562 [FOCS.2014.56](https://doi.org/10.1109/FOCS.2014.56).
- 563 Louis Béthune, Thomas Massena, Thibaut Boissin, Aurélien Bellet, Franck Mamalet, Yannick Pru-
564 dent, Corentin Friedrich, Mathieu Serrurier, and David Vigouroux. DP-SGD without clipping:
565 The lipschitz neural network way. In *The Twelfth International Conference on Learning Repre-*
566 *sentations, 2024*. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=BEyEziZ4R6)
567 [BEyEziZ4R6](https://openreview.net/forum?id=BEyEziZ4R6).
- 568 Zhiqi Bu, Hua Wang, Zongyu Dai, and Qi Long. On the convergence and calibration of deep learning
569 with differential privacy. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
570 URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=K0CAGgjYS1)
571 [K0CAGgjYS1](https://openreview.net/forum?id=K0CAGgjYS1).
- 572 Annie S Chen, Yoonho Lee, Amrith Setlur, Sergey Levine, and Chelsea Finn. Project and probe:
573 Sample-efficient adaptation by interpolating orthogonal features. In *The Twelfth International*
574 *Conference on Learning Representations, 2024*. URL [https://openreview.net/forum?](https://openreview.net/forum?id=f6CBQYxXvr)
575 [id=f6CBQYxXvr](https://openreview.net/forum?id=f6CBQYxXvr).
- 576 Xiangyi Chen, Zhiwei Steven Wu, and Mingyi Hong. Understanding gradient clipping in private
577 sgd: a geometric perspective. In *Proceedings of the 34th International Conference on Neural*
578 *Information Processing Systems, NIPS '20*, Red Hook, NY, USA, 2020a. Curran Associates Inc.
579 ISBN 9781713829546.
- 580 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum
581 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- 582 Xinlei Chen*, Saining Xie*, and Kaiming He. An empirical study of training self-supervised vision
583 transformers. *arXiv preprint arXiv:2104.02057*, 2021.
- 584 Yubei Chen, Zeyu Yun, Yi Ma, Bruno Olshausen, and Yann LeCun. Minimalistic unsupervised rep-
585 resentation learning with the sparse manifold transform. In *The Eleventh International Confer-*
586 *ence on Learning Representations, 2023*. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=nN_nBVKAhhd)
587 [nN_nBVKAhhd](https://openreview.net/forum?id=nN_nBVKAhhd).
- 588 Rishav Chourasia, Jiayuan Ye, and Reza Shokri. Differential privacy dynamics of
589 langevin diffusion and noisy gradient descent. In M. Ranzato, A. Beygelzimer,
590 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural In-*
591 *formation Processing Systems*, volume 34, pp. 14771–14781. Curran Associates, Inc.,
592 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/7c6c1a7bdfdel75bed616b39247ccacel-Paper.pdf)
593 [file/7c6c1a7bdfdel75bed616b39247ccacel-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/7c6c1a7bdfdel75bed616b39247ccacel-Paper.pdf).

- 594 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised
595 feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), *Proceedings*
596 *of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of
597 *Proceedings of Machine Learning Research*, pp. 215–223, Fort Lauderdale, FL, USA, 11–13 Apr
598 2011. PMLR. URL <https://proceedings.mlr.press/v15/coates11a.html>.
- 599 Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlock-
600 ing High-Accuracy Differentially Private Image Classification through Scale. *arXiv preprint*
601 *arXiv:2204.13650*, 2022.
- 602 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learn-
603 ing deep homogeneous models: Layers are automatically balanced. In S. Bengio,
604 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Ad-*
605 *vances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
606 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/](https://proceedings.neurips.cc/paper_files/paper/2018/file/fe131d7f5a6b38b23cc967316c13dae2-Paper.pdf)
607 [file/fe131d7f5a6b38b23cc967316c13dae2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/fe131d7f5a6b38b23cc967316c13dae2-Paper.pdf).
- 608 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Found. Trends*
609 *Theor. Comput. Sci.*, 9(3–4):211–407, aug 2014. ISSN 1551-305X. doi: 10.1561/0400000042.
610 URL <https://doi.org/10.1561/0400000042>.
- 611 Huang Fang, Xiaoyun Li, Chenglin Fan, and Ping Li. Improved convergence of differential private
612 SGD with gradient clipping. In *The Eleventh International Conference on Learning Representa-*
613 *tions*, 2023. URL <https://openreview.net/forum?id=FRLswckPXQ5>.
- 614 M.I. Freidlin, J. Szücs, and A.D. Wentzell. *Random Perturbations of Dynamical Systems*.
615 Grundlehren der mathematischen Wissenschaften. Springer, 2012. ISBN 9783642258473. URL
616 <http://books.google.de/books?id=p8LFMILAiMEC>.
- 617 Geoff French, Michal Mackiewicz, and Mark Fisher. Self-ensembling for visual domain adap-
618 tation. In *International Conference on Learning Representations*, 2018. URL [https://](https://openreview.net/forum?id=rkpoTaxA-)
619 openreview.net/forum?id=rkpoTaxA-.
- 620 Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in trans-
621 fer learning. In *International Conference on Learning Representations*, 2022. URL [https://](https://openreview.net/forum?id=SwIp410B6aQ)
622 openreview.net/forum?id=SwIp410B6aQ.
- 623 Arun Ganesh, Mahdi Haghifam, Milad Nasr, Sewoong Oh, Thomas Steinke, Om Thakkar,
624 Abhradeep Guha Thakurta, and Lun Wang. Why is public pretraining necessary for private model
625 training? In *International Conference on Machine Learning*, pp. 10611–10627. PMLR, 2023a.
- 626 Arun Ganesh, Abhradeep Thakurta, and Jalaj Upadhyay. Universality of langevin diffusion for
627 private optimization, with applications to sampling from rashomon sets. In Gergely Neu and
628 Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume
629 195 of *Proceedings of Machine Learning Research*, pp. 1730–1773. PMLR, 12–15 Jul 2023b.
630 URL <https://proceedings.mlr.press/v195/ganesh23a.html>.
- 631 Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena
632 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Ghesh-
633 laghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own
634 latent: A new approach to self-supervised learning, 2020.
- 635 Shlomo Hoory, Amir Feder, Avichai Tessler, Alon Cohen, Sofia Erell, Itay Laish, Hootan Nakhost,
636 Uri Stemmer, Ayelet Benjamini, Avinatan Hassidim, and Yossi Matias. Learning and evaluating
637 a differentially private pre-trained language model. In Oluwaseyi Feyisetan, Sepideh Ghana-
638 vati, Shervin Malmasi, and Patricia Thaine (eds.), *Proceedings of the Third Workshop on Privacy*
639 *in Natural Language Processing*, pp. 21–29, Online, June 2021. Association for Computational
640 Linguistics. doi: 10.18653/v1/2021.privatenlp-1.3. URL [https://aclanthology.org/](https://aclanthology.org/2021.privatenlp-1.3)
641 [2021.privatenlp-1.3](https://aclanthology.org/2021.privatenlp-1.3).
- 642 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
643 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*
644 *ference on Learning Representations*, 2022a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=nZeVKeeFYf9)
645 [id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).

- 648 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
649 et al. Lora: Low-rank adaptation of large language models. In *International Conference on*
650 *Learning Representations*, 2022b.
- 651
- 652 Sergey Ioffe and Christian Szegedy. Batch normalization: accelerating deep network training by
653 reducing internal covariate shift. In *Proceedings of the 32nd International Conference on In-*
654 *ternational Conference on Machine Learning - Volume 37, ICML'15*, pp. 448–456. JMLR.org,
655 2015.
- 656 Kiyosi Ito. On stochastic differential equations. *Mem. Amer. Math. Soc.*, 1951(4):51, 1951. ISSN
657 0065-9266.
- 658
- 659 Vignesh Kothapalli. Neural collapse: A review on modelling principles and generalization. *Transac-*
660 *tions on Machine Learning Research*, 2023. ISSN 2835-8856. URL [https://openreview.](https://openreview.net/forum?id=QTXocpAP9p)
661 [net/forum?id=QTXocpAP9p](https://openreview.net/forum?id=QTXocpAP9p).
- 662 Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Canadian
663 Institute for Advanced Research, 2009.
- 664
- 665 Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-
666 tuning can distort pretrained features and underperform out-of-distribution. In *International Con-*
667 *ference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=UYneFzXSJWh)
668 [id=UYneFzXSJWh](https://openreview.net/forum?id=UYneFzXSJWh).
- 669 Don S. Lemons and Anthony Gythiel. Paul Langevin’s 1908 paper “On the Theory of Brownian
670 Motion” [“Sur la théorie du mouvement brownien,” *C. R. Acad. Sci. (Paris)* 146, 530–533 (1908)].
671 *American Journal of Physics*, 65(11):1079–1081, 11 1997. ISSN 0002-9505. doi: 10.1119/1.
672 18725. URL <https://doi.org/10.1119/1.18725>.
- 673
- 674 Jan R. Magnus and Heinz Neudecker. *Matrix Differential Calculus with Applications in Statis-*
675 *tics and Econometrics*. John Wiley, second edition, 1999. ISBN 0471986321 9780471986324
676 047198633X 9780471986331.
- 677 Xuerong. Mao. *Stochastic differential equations and their applications / Xuerong Mao*. Horwood
678 series in mathematics & applications. Horwood Pub., Chichester, 1997. ISBN 1898563268.
- 679
- 680 Sibylle Marcotte, Rémi Gribonval, and Gabriel Peyré. Abide by the law and follow the flow: conser-
681 vation laws for gradient flows. In *Thirty-seventh Conference on Neural Information Processing*
682 *Systems*, 2023. URL <https://openreview.net/forum?id=kMueEV8Eyy>.
- 683 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and
684 projection for dimension reduction, 2020. URL <https://arxiv.org/abs/1802.03426>.
- 685
- 686 Hancheng Min, Salma Tarmoun, Rene Vidal, and Enrique Mallada. On the explicit role of initializa-
687 tion on the convergence and implicit bias of overparametrized linear networks. In Marina Meila
688 and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*,
689 volume 139 of *Proceedings of Machine Learning Research*, pp. 7760–7768. PMLR, 18–24 Jul
690 2021. URL <https://proceedings.mlr.press/v139/min21c.html>.
- 691 Hancheng Min, Rene Vidal, and Enrique Mallada. On the convergence of gradient flow on multi-
692 layer linear models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
693 Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on*
694 *Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 24850–24887.
695 PMLR, 23–29 Jul 2023a. URL [https://proceedings.mlr.press/v202/min23d.](https://proceedings.mlr.press/v202/min23d.html)
696 [html](https://proceedings.mlr.press/v202/min23d.html).
- 697 Hancheng Min, Rene Vidal, and Enrique Mallada. On the convergence of gradient flow on multi-
698 layer linear models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
699 Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on*
700 *Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 24850–24887.
701 PMLR, 23–29 Jul 2023b. URL [https://proceedings.mlr.press/v202/min23d.](https://proceedings.mlr.press/v202/min23d.html)
[html](https://proceedings.mlr.press/v202/min23d.html).

- 702 Hancheng Min, Enrique Mallada, and Rene Vidal. Early neuron alignment in two-layer reLU net-
703 works with small initialization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=QibPzdVrRu>.
704
705
- 706 Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Sym-*
707 *posium (CSF)*, pp. 263–275, 2017. doi: 10.1109/CSF.2017.11.
708
- 709 Bernt Øksendal. *Stochastic Differential Equations: An Introduction with Applications (Universi-*
710 *text)*. Springer, 6th edition, January 2014. ISBN 3540047581. URL [http://www.amazon.](http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/3540047581)
711 [com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/3540047581](http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/3540047581).
712
- 713 Guillermo Ortiz-Jimenez, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. What
714 can linearized neural networks actually say about generalization? In M. Ranzato,
715 A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neu-*
716 *ral Information Processing Systems*, volume 34, pp. 8998–9010. Curran Associates, Inc.,
717 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/4b5deb9a14d66ab0acc3b8a2360cde7c-Paper.pdf)
718 [file/4b5deb9a14d66ab0acc3b8a2360cde7c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/4b5deb9a14d66ab0acc3b8a2360cde7c-Paper.pdf).
719
- 720 Natalia Ponomareva, Sergei Vassilvitskii, Zheng Xu, Brendan McMahan, Alexey Kurakin, and
721 Chiyaun Zhang. How to dp-fy ml: A practical tutorial to machine learning with differen-
722 tial privacy. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery*
723 *and Data Mining, KDD '23*, pp. 5823–5824, New York, NY, USA, 2023. Association
724 for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599561. URL
725 <https://doi.org/10.1145/3580305.3599561>.
726
- 727 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
728 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
729 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
730 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
731
- 732 Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar.
733 A theoretical analysis of contrastive unsupervised representation learning. In Kamalika Chaudhuri
734 and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine*
735 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5628–5637. PMLR,
736 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/saunshi19a.html>.
737
- 738 Anatoli V. Skorokhod, Frank C. Hoppensteadt, and Habib Salehi. *Random Perturbation Methods*
739 *with Applications in Science and Engineering*. Applied mathematical sciences (Springer-Verlag
740 New York Inc.) ; v. 150. Springer, 2002. ISBN 0387954279. doi: 10.1115/1.1579453.
741
- 742 Xinyu Tang, Ashwinee Panda, Vikash Sehwal, and Prateek Mittal. Differentially private image
743 classification by learning priors from random processes. *CoRR*, abs/2306.06076, 2023. URL
744 <https://doi.org/10.48550/arXiv.2306.06076>.
745
- 746 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
747 Herve Jegou. Training data-efficient image transformers & distillation through attention. In
748 Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on*
749 *Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–
750 10357. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/](https://proceedings.mlr.press/v139/touvron21a.html)
751 [touvron21a.html](https://proceedings.mlr.press/v139/touvron21a.html).
752
- 753 Nilesch Tripuraneni, Michael Jordan, and Chi Jin. On the theory of transfer learning: The im-
754 portance of task diversity. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
755 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 7852–7862. Cur-
756 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/59587bffe1c7846f3e34230141556ae-Paper.pdf)
757 [paper/2020/file/59587bffe1c7846f3e34230141556ae-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/59587bffe1c7846f3e34230141556ae-Paper.pdf).
758
- 759 Puja Trivedi, Danai Koutra, and Jayaraman J. Thiagarajan. A closer look at model adaptation using
760 feature distortion and simplicity bias. In *The Eleventh International Conference on Learning*
761 *Representations*, 2023. URL https://openreview.net/forum?id=wkg_b4-IwTZ.

- 756 A J Veretennikov. On strong solutions and explicit formulas for solutions of stochastic in-
757 tegral equations. *Mathematics of the USSR-Sbornik*, 39(3):387, apr 1981. doi: 10.1070/
758 SM1981v039n03ABEH001522.
- 759
760 Chendi Wang, Yuqing Zhu, Weijie J Su, and Yu-Xiang Wang. Neural collapse meets differential
761 privacy: Curious behaviors of NoisyGD with near-perfect representation learning. In Ruslan
762 Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and
763 Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*,
764 volume 235 of *Proceedings of Machine Learning Research*, pp. 52334–52360. PMLR, 21–27 Jul
765 2024. URL <https://proceedings.mlr.press/v235/wang24cu.html>.
- 766
767 Di Wang, Changyou Chen, and Jinhui Xu. Differentially private empirical risk minimization with
768 non-convex loss functions. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings*
769 *of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine*
770 *Learning Research*, pp. 6526–6535. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/wang19c.html>.
- 771
772 Puyu Wang, Yunwen Lei, Yiming Ying, and Hai Zhang. Differentially private sgd with non-smooth
773 losses. *Applied and Computational Harmonic Analysis*, 56:306–336, 2022. ISSN 1063-5203. doi:
774 <https://doi.org/10.1016/j.acha.2021.09.001>. URL <https://www.sciencedirect.com/science/article/pii/S1063520321000841>.
- 775
776 Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on*
777 *Computer Vision (ECCV)*, September 2018.
- 778
779 Jiayuan Ye, Zhenyu Zhu, Fanghui Liu, Reza Shokri, and Volkan Cevher. Initialization matters:
780 Privacy-utility analysis of overparameterized neural networks. In *Thirty-seventh Conference on*
781 *Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=IKvxmnHjkL>.
- 782
783 Jiayuan Ye, Zhenyu Zhu, Fanghui Liu, Reza Shokri, and Volkan Cevher. Initialization mat-
784 ters: Privacy-utility analysis of overparameterized neural networks. In A. Oh, T. Nau-
785 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neu-
786 ral Information Processing Systems*, volume 36, pp. 5419–5446. Curran Associates, Inc.,
787 2023b. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/1165af8b913fb836c6280b42d6e0084f-Paper-Conference.pdf.
- 788
789 Tian Ye and Simon Shaolei Du. Global convergence of gradient descent for asymmetric low-rank
790 matrix factorization. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.),
791 *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=sMIMAXqiqj3>.
- 792
793 Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan
794 Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and Huishuai Zhang.
795 Differentially private fine-tuning of language models. In *International Conference on Learning*
796 *Representations*, 2022. URL <https://openreview.net/forum?id=Q42f0dfjECO>.
- 797
798 Sergey Zagoruyko and Nikos Komodakis. Wide Residual Networks. In *British Machine Vision*
799 *Conference 2016*, York, France, January 2016. British Machine Vision Association. doi: 10.
800 5244/C.30.87. URL <https://enpc.hal.science/hal-01832503>.
- 801
802 Jinnian Zhang, Houwen Peng, Kan Wu, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan.
803 Minivit: Compressing vision transformers with weight multiplexing. In *2022 IEEE/CVF Con-*
804 *ference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12135–12144, 2022. doi:
805 10.1109/CVPR52688.2022.01183.
- 806
807 Liang Zhang, Bingcong Li, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. DPZero:
808 private fine-tuning of language models without backpropagation. In *Forty-first International Con-*
809 *ference on Machine Learning*, 2024a.
- 807
808 Xinwei Zhang, Zhiqi Bu, Steven Wu, and Mingyi Hong. Differentially private SGD without clip-
809 ping bias: An error-feedback approach. In *The Twelfth International Conference on Learning*
Representations, 2024b. URL <https://openreview.net/forum?id=uFbWWhyTlPn>.

A ADDITIONAL EXPERIMENT RESULTS

In this section, we provide more experiment results and detailed configurations.

Evaluations back in the pre-training distribution (Figure 5). We also evaluate the feature quality on ImageNet1-K, the pre-training dataset. The representation alignment for the pre-training domain is different: once a proper alignment is achieved, the backbone gradually recovers a portion of its original feature quality, which had been compromised due to DP noise and distribution-shift.

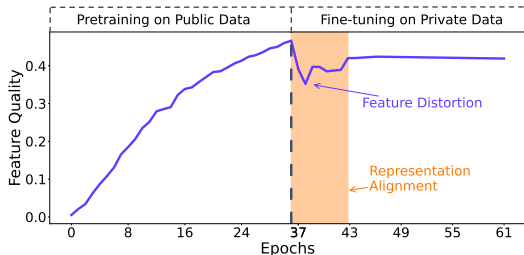


Figure 5: Backbone feature quality evaluated by average top-1 kNN accuracy on the pre-training dataset, for ResNet-50, through public pre-training on ImageNet-1K and differentially private fine-tuning on STL-10.

More experiments on parameter-efficient fine-tuning (PEFT) methods. We conduct experiments with another fine-tuning trick: differentially private LoRA (Hu et al., 2022a). We run experiments on the Mini-DeiT-Ti architecture, where we use LoRA instead of full fine-tuning. In these experiments (Table 2), our batch size is 1000, and our LoRA rank is set to 8. We observe the same trend as what we saw for full fine-tuning; namely, as we increase the noise scale (i.e., as we reduce epsilon, giving a stronger privacy guarantee), it becomes more beneficial to use LP-LoRA or even just LP.

Transformer _{DeiT}					
ϵ	∞	12.28	1.29	0.57	0.26
LP	95.81 _{0.05}	95.55 _{0.05}	94.80 _{0.06}	94.21 _{0.08}	92.48 _{0.27}
LP-LoRA	96.2 _{0.05}	95.90 _{0.03}	94.81 _{0.08}	94.18 _{0.05}	91.99 _{0.19}
LoRA	96.26 _{0.05}	95.50 _{0.06}	94.76 _{0.08}	93.05 _{0.09}	91.28 _{0.43}

Table 2: Test accuracies of LP, LP-LoRA, LoRA on Transformer_{DeiT}.

Experiment setup in Table 1. We use batch size 1000 and sweep over a range of learning rates {9, 5, 1, 0.5, 0.2, 0.15, 0.1, 0.05, 0.025}.

Summary of experiment configurations. We run experiments on five deep learning models and four transfer learning benchmarks to verify if our theoretical prediction, the existence of concave utility curves, generalizes to deep neural networks and real datasets. Each experimental setting comprises: (1) a model architecture, (2) a (larger) dataset for public pretraining, and (3) a (smaller) dataset as the private data for fine-tuning. The benchmarks we use are:

- ImageNet-1K→CIFAR-10. ImageNet-1K is a large-scale dataset. We initialize pretrained features of ResNet-50 from MoCo-v2 Chen et al. (2020b) and MoCo-v3 Chen* et al. (2021), trained on ImageNet-1K Russakovsky et al. (2015) without privacy. We then privately fine-tune the ResNet-50 on CIFAR-10.
- ImageNet-1K→STL-10. We pretrain a DeiT model on ImageNet then pretrain a Mini-DeiT-Ti model with weight distillation from the DeiT model Touvron et al. (2021); Zhang et al. (2022). After that, we privately fine-tune the Mini-DeiT-Ti model on STL-10 Coates et al. (2011) for 20 epochs.
- CIFAR-10→STL-10. We pretrain the feature extractor on CIFAR-10 Krizhevsky (2009) using stochastic gradient descent without privacy mechanisms. Then we finetune the pre-trained features and a randomly initialized linear head on STL-10. This benchmark has been studied in the context of domain adaptation French et al. (2018); Kumar et al. (2022).

864 The training subset of STL-10 only contains 500 images. To align with the small scale fine-
 865 tuning data, we run the experiments with smaller and data-efficient models: MobileNet-v3
 866 and ResNet-18.

- 867 • RandP→CIFAR-10. To reproduce the results of Tang et al. (2023) and verify the general
 868 existence of concave utility curves, we also consider a slightly non-standard pretraining
 869 protocol. We pretrain a wide residual network (WRN) Zagoruyko & Komodakis (2016) on
 870 synthetic images generated by random diffusion processes. We follow the settings in Tang
 871 et al. (2023).

872 We employ early stopping, and select the optimal learning rate based on the accuracy of the in-
 873 distribution validation.

874 B TECHNICAL RESULTS

875 **Lemma B.1** (Holder’s inequality for sums). *For a sequence $x = [x_i]_{i=1}^n$ of positive real numbers
 876 and $p > 0$, define $\|x\|_p := (\sum_{i=1}^n x_i^p)^{1/p}$. Then for any pair of positive real numbers $p > 0, q > 0$
 877 with $\frac{1}{p} + \frac{1}{q} = 1$, and any pair of sequence of positive real numbers x and y ,*

$$878 \quad \|xy\|_1 \leq \|x\|_p \|y\|_q$$

883 **Lemma B.2** (Reverse Holder’s inequality for sums). *For a sequence $x = [x_i]_{i=1}^n$ of positive real
 884 numbers and $p > 0$, define $\|x\|_p := (\sum_{i=1}^n x_i^p)^{1/p}$. Then for any pair of positive real numbers
 885 $p > 0, q > 0$ with $\frac{1}{p} - \frac{1}{q} = 1$, and any pair of sequence of positive real numbers x and y ,*

$$886 \quad \|xy\|_1 \geq \|x\|_p \|y\|_{-q}$$

888 **Lemma B.3** (Reverse QM-AM inequality for sums). *For a sequence $x = [x_i]_{i=1}^n$ of positive real
 889 numbers,*

$$890 \quad \left(\sum_{i=1}^n x_i \right)^2 \geq \sum_{i=1}^n x_i^2$$

893 **Lemma B.4** (μ -coherent data conic hull (Min et al., 2024, Lemma 5)). *Define a conic hull $K :=$
 894 $\mathcal{CH}(\{y_i x_i : i \in [n]\}) = \{\sum_{i=1}^n a_i y_i x_i : \forall a_i \geq 0, i \in [n]\}$. If Assumption 3.1 holds, i.e. the dataset
 895 is separable, then K is μ -coherent:*

$$896 \quad \forall z_1, z_2 \in K \setminus \{0\}, \quad \cos(z_1, z_2) \geq \mu$$

897 **Corollary B.5** (Orthogonally separable \implies linearly separable (Min et al., 2024)). *If Assumption 3.1
 898 holds, then $\exists \gamma > 0$ and $z \in \mathbb{S}^{D-1}$ such that*

$$899 \quad \forall i \in [n], \quad y_i \langle z, x_i \rangle \geq \gamma$$

901 *Proof of Corollary B.5.* We prove the existence statement by picking a valid pair of z, γ . Take $z :=$
 902 $\frac{y_1 x_1}{\|x_1\|_2}$. Then $\forall i \in [n]$,

$$903 \quad y_i \langle z, x_i \rangle = \|x_i\|_2 \cos(y_1 x_1, y_i x_i)$$

$$904 \quad \quad \quad // \text{by Lemma B.4}$$

$$905 \quad \geq \|x_i\|_2 \mu$$

$$906 \quad \geq \mu \cdot \min_{i \in [n]} \|x_i\|_2$$

907 Therefore $\gamma = \mu \cdot \min_{i \in [n]} \|x_i\|_2$. □

911 C APPENDIX: REPRESENTATION ALIGNMENT

912 C.1 THEORY

913 The Langevin diffusion of w_j on a n -sized data cluster ($j \in [h]$) is

$$914 \quad \dot{w}_j = \sum_{i=1}^n y_i \exp(-y_i f(x_i; W, v)) v_j \text{relu}'(w_j^\top x_i) x_i + \sigma \delta Q_t, \quad (23)$$

where Q_t is a vector containing D independent 1-dimensional Brownian motion.

The Langevin diffusion of v on a n -sized data cluster is

$$\dot{v} = \sum_{i=1}^n y_i \exp(-y_i f(x_i; W, v)) \text{relu}(W^\top x_i) + \sigma \partial Q_t,$$

where Q_t is a vector containing h independent 1-dimensional Brownian motion.

We rewrite the Langevin diffusion by asymptotic expansion (Freidlin et al., 2012, Equation 2.1, Chapter 2.2),

$$\begin{cases} v_j \approx v_j^{(0)} + \sigma v_j^{(1)} + \dots \\ w_j \approx w_j^{(0)} + \sigma w_j^{(1)} + \dots, \end{cases} \quad (24)$$

i.e. we expand the Langevin diffusion as a linear combination of the original gradient flow and a linear stochastic diffusion.

$$\begin{cases} \dot{v}_j^{(0)} = \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \\ \dot{w}_j^{(0)} = \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \text{relu}'((w_j^{(0)})^\top x_i) x_i. \end{cases} \quad (25)$$

Lemma C.1 (Zeroth order invariance of locally linearized LD). *If we rewrite the Langevin diffusion by asymptotic expansion (Freidlin et al., 2012, Equation 2.1, Chapter 2.2),*

$$\begin{cases} v_j \approx v_j^{(0)} + \sigma v_j^{(1)} \\ w_j \approx w_j^{(0)} + \sigma w_j^{(1)}. \end{cases}$$

then the layer invariance still holds for zeroth order approximation

$$\frac{d}{dt} [(v_j^{(0)})^2 - \|w_j^{(0)}\|_2^2] = 0. \quad (26)$$

This result is similar to the imbalance matrix in gradient flow (Arora et al., 2018; Du et al., 2018; Min et al., 2023a).

We are ready to prove Theorem 3.3.

Proof of Theorem 3.3. The explicit expression of the cosine value is

$$\cos(w_j, \bar{x}_{c(j)}) = \frac{w_j^\top \bar{x}_{c(j)}}{\|w_j\|_2 \|\bar{x}_{c(j)}\|_2} \quad (27)$$

Without loss of generality, let $\|\bar{x}_{c(j)}\|_2 = 1$. To show that the cosine value decreases with high probability, we only need to prove that the derivative of $\frac{(w_j^\top \bar{x}_{c(j)})^2}{\|w_j\|_2^2}$ is negative at $t = 0$ with high probability. The explicit derivative expression is

$$\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) = \frac{2(w_j^\top \bar{x}_{c(j)})}{\|w_j\|_2^2} \left[\|w_j\|_2^2 \bar{x}_{c(j)}^\top \frac{\partial w_j}{\partial t} - \bar{x}_{c(j)}^\top w_j w_j^\top \frac{\partial w_j}{\partial t} \right] \quad (28)$$

$$= \frac{2(w_j^\top \bar{x}_{c(j)})}{\|w_j\|_2^2} \left[\|w_j\|_2^2 \bar{x}_{c(j)} - (\bar{x}_{c(j)}^\top w_j) w_j \right]^\top \frac{\partial w_j}{\partial t} \quad (29)$$

$$// \text{by Assumption 3.2} \quad (30)$$

$$\text{sign} \left(\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right) = \text{sign} \left(\left[\|w_j\|_2^2 \bar{x}_{c(j)} - (\bar{x}_{c(j)}^\top w_j) w_j \right]^\top \frac{\partial w_j}{\partial t} \right) \quad (31)$$

$$= \text{sign} \left(v_j (\|w_j\|_2^2 - (\bar{x}_{c(j)}^\top w_j)^2) \right) \quad (32)$$

$$= \text{sign}(v_j) \quad (33)$$

Since we initialize $v \sim \mathcal{N}(0, \beta I_{h \times h})$, with probability $1 - 2^{-h}$, there exists j such that $v_j < 0$ at $t = 0 \implies \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) < 0$ at $t = 0$. By the continuity of the approximated Langevin diffusion, there exists $\Delta t > 0$ such that for any $t \in (0, \Delta t)$,

$$\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) < 0. \quad (34)$$

□

972 *Proof of Theorem 3.4.* In the proof of Theorem 3.3, we show that for $w_j \in S_c$, $c \in \{-1, 1\}$,

$$973 \text{sign} \left(\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right) = \text{sign}(v_j) \cdot \text{sign}(c) \quad (35)$$

974 To mitigate the feature distortion after some time index Δt , we only need $c \cdot v_j > 0$. For DP-LP,
 975 every $\frac{\partial}{\partial t} v_j$ increases/decreases if $c = 1/-1$. Therefore, for any initialization, there exists Δt such
 976 that $\text{sign}(v_j) = \text{sign}(c)$ after time index Δt . If we switch to DP-FFT after Δt , $\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) > 0$
 977 for any $j \in [h]$. Thus $\cos(w_j, \bar{x}_{c(j)})$ is non-decreasing in DP-FFT. \square

982 D APPROXIMATE CONVERGENCE OF DP-LP AND DP-FFT

983 D.1 APPROXIMATE DP-LP CONVERGENCE

984 We add some extra notations for the following proofs:

- 985 • Positive data subset $\mathcal{I}_+ := \{i \in [n] : y_i > 0\}$
- 986 • Negative data subset $\mathcal{I}_- := \{i \in [n] : y_i < 0\}$
- 987 • Positive head cluster $\mathcal{V}_+(t) := \{j \in [h] : \text{sign}(v_j(t)) > 0\}$
- 988 • Negative head cluster $\mathcal{V}_-(t) := \{j \in [h] : \text{sign}(v_j(t)) < 0\}$
- 989 • Index function $\mathcal{S} : \mathbb{R}^D \rightarrow \{\mathcal{I}_+, \mathcal{I}_-\}$ maps feature vector to its cluster

$$990 \mathcal{S}(w) = \begin{cases} \mathcal{I}_+ & w \in S_+ \\ \mathcal{I}_- & w \in S_- \\ \emptyset & \text{otherwise} \end{cases}$$

991 We first derive the upper bound for approximate DP-LP.

992 *Upper bound proof of Theorem 4.2.* We construct a lower bound of the drift terms in the zeroth
 993 order approximation

$$1000 \|\nabla_v \mathcal{L}^{(0)}\|_2^2 = \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \quad (36)$$

$$1001 = \sum_{j=1}^h \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \quad (37)$$

$$1002 \geq \sum_{j=1}^h \left[\min_{i \in \mathcal{S}(w_j^{(0)})} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2 \quad (38)$$

$$1003 = \sum_{j=1}^h \left[\min_{i \in \mathcal{S}(w_j^{(0)})} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2 \quad (39)$$

$$1004 = \sum_{j \in \mathcal{V}_+} \left[\min_{i \in \mathcal{I}_+} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 (\mathcal{L}_+^{(0)})^2 + \sum_{j \in \mathcal{V}_-} \left[\min_{i \in \mathcal{I}_-} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 (\mathcal{L}_+^{(0)})^2 \quad (40)$$

$$1005 \geq \min \left\{ \sum_{j \in \mathcal{V}_+} \left[\min_{i \in \mathcal{I}_+} \text{relu}((w_j^{(0)})^\top x_i) \right]^2, \sum_{j \in \mathcal{V}_-} \left[\min_{i \in \mathcal{I}_-} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right\} [(\mathcal{L}_+^{(0)})^2 + (\mathcal{L}_-^{(0)})^2] \quad (41)$$

$$\geq \frac{1}{2} \min \left\{ \sum_{j \in \mathcal{V}_+} \left[\min_{i \in \mathcal{I}_+} \text{relu}((w_j^{(0)})^\top x_i) \right]^2, \sum_{j \in \mathcal{V}_-} \left[\min_{i \in \mathcal{I}_-} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right\} \left[\mathcal{L}_+^{(0)} + \mathcal{L}_-^{(0)} \right]^2 \quad (42)$$

$$= \frac{1}{2} \min \left\{ \sum_{j \in \mathcal{V}_+} \left[\min_{i \in \mathcal{I}_+} \text{relu}((w_j^{(0)})^\top x_i) \right]^2, \sum_{j \in \mathcal{V}_-} \left[\min_{i \in \mathcal{I}_-} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right\} (\mathcal{L}^{(0)})^2 \quad (43)$$

We construct an upper bound of the diffusion terms in the zeroth order approximation

$$\begin{aligned} & \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \left\{ \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\} \cdot \left\{ \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right\} \\ & \quad // \text{by Lemma B.1} \\ & \leq \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \ell^2(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\}^{1/2} \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^4 \right\}^{1/2} \\ & \quad // \text{by Lemma B.3} \\ & \leq \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\} \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^4 \right\}^{1/2} \\ & = \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^4 \right\}^{1/2} \end{aligned}$$

Then we have an upper bound

$$\mathcal{L}^{(0)}(T) \leq \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)} e^{-BT} + \frac{A}{B} (1 - e^{-BT})}$$

where constants A, B are defined as

$$\begin{cases} A = \frac{1}{2} \min \left\{ \sum_{j \in \mathcal{V}_+} \left[\min_{i \in \mathcal{I}_+} \text{relu}((w_j^{(0)})^\top x_i) \right]^2, \sum_{j \in \mathcal{V}_-} \left[\min_{i \in \mathcal{I}_-} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right\} \\ B = \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^4 \right\}^{1/2} \end{cases}$$

□

We give the lower bound of approxiamte DP-LP below. We first give a loose lower bound as a warm-up. Then we improve the techniques and provide a tight lower bound.

Loose lower bound proof of Theorem 4.2. We rewrite the Langevin diffusion by asymptotic expansion (Freidlin et al., 2012, Equation 2.1, Chapter 2.2)

$$\begin{aligned} \dot{\mathcal{L}}^{(0)} &= - \|\nabla_v \mathcal{L}^{(0)}\|_2^2 + \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &= - \|\nabla_v \mathcal{L}^{(0)}\|_2^2 + \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &\geq - \|\nabla_v \mathcal{L}^{(0)}\|_2^2 + \left(\min_{i \in \mathcal{V}_+^{(0)}} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \sum_{i \in \mathcal{V}_+^{(0)}} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \end{aligned}$$

$$\begin{aligned}
& + \left(\min_{i \in \mathcal{V}_-^{(0)}} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \sum_{i \in \mathcal{V}_-^{(0)}} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
& = - \|\nabla_v \mathcal{L}^{(0)}\|_2^2 + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \sum_{i \in [n]} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
& = - \|\nabla_v \mathcal{L}^{(0)}\|_2^2 + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& = - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& \quad // \text{by trapping} \\
& = - \sum_{j \in \mathcal{V}_+^{(0)}} \left(\sum_{i \in \mathcal{I}_+} \exp(-f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\
& \quad - \sum_{j \in \mathcal{V}_-^{(0)}} \left(\sum_{i \in \mathcal{I}_-} \exp(f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\
& \quad + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& \geq - \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right) \sum_{j \in \mathcal{V}_+^{(0)}} \left(\sum_{i \in \mathcal{I}_+} \exp(-f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
& \quad - \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right) \sum_{j \in \mathcal{V}_-^{(0)}} \left(\sum_{i \in \mathcal{I}_-} \exp(f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
& \quad + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& \quad // a^2 + b^2 \leq (a + b)^2 \text{ when } a > 0, b > 0 \\
& \geq - \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right) \sum_{j \in [h]} \left(\sum_{i \in [n]} \exp(-f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
& \quad + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& \geq - h \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right) \left(\sum_{i \in [n]} \exp(-f(x_i; W^{(0)}, v^{(0)})) \right)^2 + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \\
& \geq - h \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right) (\mathcal{L}^{(0)})^2 + \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right) \cdot \frac{1}{2} \sigma^2 \mathcal{L}^{(0)}
\end{aligned}$$

In linear probing, the coefficients $h \left(\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2 \right)$ and $\frac{1}{2} \sigma^2 \left(\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right)$ are constants. We replace them with dummy notation A and B . We solve the first-order nonlinear ODE by turning it into a first-order linear ODE.

$$\begin{aligned}
\dot{\mathcal{L}}^{(0)} & \geq -A(\mathcal{L}^{(0)})^2 + B\mathcal{L}^{(0)} \\
\frac{1}{(\mathcal{L}^{(0)})^2} \dot{\mathcal{L}}^{(0)} & \geq -A + B \frac{1}{\mathcal{L}^{(0)}}
\end{aligned}$$

1134
1135
1136
1137
1138
1139
1140

$$-\frac{d}{dt} \left(\frac{1}{\mathcal{L}^{(0)}} \right) \geq -A + B \frac{1}{\mathcal{L}^{(0)}}$$

$$\mathcal{L}^{(0)}(T) \geq \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)}e^{-BT} + \frac{A}{B}(1 - e^{-BT})}$$

□

1141
1142
1143
1144
1145

Remark D.1 (On the qualitative properties of loose DP-LP lower bound). If we take the limit to initial point, then the lower bound degenerate to the initial loss value.

$$\lim_{t \rightarrow 0} \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)}e^{-BT} + \frac{A}{B}(1 - e^{-BT})} = \mathcal{L}^{(0)}(t=0) = \mathcal{L}(t=0) \quad (44)$$

1146
1147

If we take the limit to infinite time,

1148
1149
1150

$$\lim_{t \rightarrow \infty} \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)}e^{-BT} + \frac{A}{B}(1 - e^{-BT})} = \frac{B}{A} = \frac{\frac{1}{2}\sigma^2 (\min_{i \in [n]} \|\text{relu}((W^{(0)})^\top x_i)\|_2^2)}{h (\max_{j \in [h], i \in [n]} (\text{relu}((w_j^{(0)})^\top x_i))^2)} \quad (45)$$

1151
1152

the following interpretation holds:

1153
1154
1155
1156

1. For larger noise $\sigma \uparrow$, the lower bound is higher, i.e. worse performance.
2. For bad alignment between pretrained features $W^{(0)}$ and data points, both the denominator and the numerator could shrink. It is not obvious how the lower bound changes.

1157
1158

In the following result, we modify the proof, replace the $\min(\cdot)$, and provide a tighter bound.

1159
1160
1161

Tight lower bound proof of Theorem 4.2. This is an alternative construction of a lower bound for drift terms in the zeroth order approximation

1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

$$\begin{aligned} \|\nabla_v \mathcal{L}^{(0)}\|_2^2 &= \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &= \sum_{j \in \mathcal{V}_+^{(0)}} \left(\sum_{i \in \mathcal{I}_+} \exp(-f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &\quad + \sum_{j \in \mathcal{V}_-^{(0)}} \left(\sum_{i \in \mathcal{I}_-} \exp(f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &\quad // \text{by Lemma B.3} \\ &\leq \left(\sum_{j \in \mathcal{V}_+^{(0)}} \sum_{i \in \mathcal{I}_+} \exp(-f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &\quad + \left(\sum_{j \in \mathcal{V}_-^{(0)}} \sum_{i \in \mathcal{I}_-} \exp(f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &\leq \left(\sum_{j \in [h]} \sum_{i \in [n]} \exp(-f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &= \left(\sum_{i \in [n]} \sum_{j \in [h]} \exp(-f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \end{aligned}$$

$$\begin{aligned}
&\leq \left(\sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right] \exp(-f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
&\quad // \text{by Lemma B.1} \\
&\leq \left(\sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right) \left(\sum_{i \in [n]} \exp(-f(x_i; W^{(0)}, v^{(0)}))^2 \right) \\
&\quad // \text{by Lemma B.3} \\
&\leq \left(\sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right) \left(\sum_{i \in [n]} \exp(-f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
&\leq \left(\sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right) (\mathcal{L}^{(0)})^2
\end{aligned}$$

We replace the A constant by $\sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2$. This is an alternative construction of a lower bound for diffusion-resulted terms in the zeroth order approximation

$$\begin{aligned}
&\frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\
&= \frac{1}{2} \sigma^2 \sum_{i=1}^n \left\{ \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\} \cdot \left\{ \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \right\} \\
&\quad // \text{by Lemma B.2} \\
&\geq \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \ell^{1/2}(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\}^2 \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1} \\
&\quad // \text{by Lemma B.3} \\
&\geq \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right\} \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1} \\
&\geq \frac{1}{2} \sigma^2 \mathcal{L}^{(0)} \cdot \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1}
\end{aligned}$$

We replace the B constant by $\left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1}$ in the previous proof of loose lower bound of Theorem 4.2. Similarly,

$$\mathcal{L}^{(0)}(T) \geq \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)} e^{-BT} + \frac{A}{B} (1 - e^{-BT})}$$

where $A = \sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2$, $B = \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1}$. The limit of this lower bound is

$$\lim_{t \rightarrow \infty} \frac{1}{\frac{1}{\mathcal{L}^{(0)}(0)} e^{-BT} + \frac{A}{B} (1 - e^{-BT})} = \frac{B}{A} = \frac{1}{2} \sigma^2 \left\{ \sum_{i=1}^n \|\text{relu}((W^{(0)})^\top x_i)\|_2^{-2} \right\}^{-1} \left\{ \sum_{i \in [n]} \left[\max_{j \in [h]} \text{relu}((w_j^{(0)})^\top x_i) \right]^2 \right\}^{-1}$$

□

Example D.2 (On the downstream alignment of pretrained features (Theorem 4.2)). Here we provide an example on how the pretrained feature space affects the linear probing lower bound in Theorem 4.2 in the **overparametrized** regime. Consider one data point x_+ and two pretrained features $w_{+,1}, w_{+,2}$ with $\|x_+\|_2 = \|w_{+,1}\|_2 = \|w_{+,2}\|_2 = 1$, $\cos(x_+, w_{+,2}) = \frac{1}{3}\pi$.

- 1242 1. If we get lucky such that $w_{+,1} = x_+$, then the limit is $\frac{B}{A} = \frac{15}{24}\sigma^2$.
 1243
 1244 2. If the $w_{+,1}$ is not so good for the downstream task such that $\cos(x_+, w_{+,1}) = \frac{1}{6}\pi$, then the
 1245 limit becomes $\frac{B}{A} = \frac{16}{24}\sigma^2$.
 1246

1247 Since $\frac{16}{24} > \frac{15}{24}$, we can tell that when the pretrained features do not align well with the downstream
 1248 task, the lower bound gets higher, i.e. worse performance.
 1249

1250 D.2 APPROXIMATE DP-FT CONVERGENCE

1251
 1252 **Analysis of DP-FFT loss diffusion.** In the following 0th-order approximation of loss Langevin
 1253 diffusion, denote the drift term by W -gradient as T_1 , the drift term by v -gradient as T_2 , the diffusion
 1254 term by W -hessian as T_3 , the diffusion term by v -hessian as T_4 .

$$1255 \dot{\mathcal{L}}^{(0)} = - \underbrace{\left\| \nabla_W \mathcal{L}^{(0)} \right\|_F^2}_{T_1} - \underbrace{\left\| \nabla_v \mathcal{L}^{(0)} \right\|_2^2}_{T_2} \quad (46)$$

$$1256 + \frac{1}{2}\sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \left(\left\| \text{relu}((W^{(0)})^\top x_i) \right\|_2^2 + \sum_{j=1}^h (v_j^{(0)})^2 [\text{relu}'((w_j^{(0)})^\top x_i)]^2 \|x_i\|_2^2 \right) \quad (47)$$

$$1257 = - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \quad (48)$$

$$1258 - \sum_{j=1}^h \left\| \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i \right\|_2^2 \quad (49)$$

$$1259 + \frac{1}{2}\sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \left(\left\| \text{relu}((W^{(0)})^\top x_i) \right\|_2^2 + \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2 \right) \quad (50)$$

$$1260 = - \underbrace{\sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2}_{T_2} \quad (51)$$

$$1261 - \underbrace{\sum_{j=1}^h \left\| \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i \right\|_2^2}_{T_1} \quad (52)$$

$$1262 + \underbrace{\frac{1}{2}\sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \left\| \text{relu}((W^{(0)})^\top x_i) \right\|_2^2}_{T_4} \quad (53)$$

$$1263 + \underbrace{\frac{1}{2}\sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2}_{T_3} \quad (54)$$

1291 *Upper bound proof of Theorem 4.3. 1. Upper bounds for T_1, T_3 .* For T_1 , the key idea is $\|x\|_2^2 \geq$
 1292 $\langle x, z \rangle^2$ for any unit vector z .
 1293

$$1294 T_1 = - \sum_{j=1}^h \left\| \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i \right\|_2^2$$

$$\begin{aligned}
& // \text{since } \forall x \in \mathbb{R}^D, z \in \mathbb{S}^{D-1}, \|x\|_2^2 \geq \langle x, z \rangle^2 \\
& \leq - \sum_{j=1}^h \left\langle \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i, z \right\rangle^2 \\
& = - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle x_i, z \rangle \right)^2 \\
& = - \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle x_i, z \rangle \right)^2 \\
& // \text{pick } z = \frac{y_1 x_1}{\|x_1\|_2}, \text{ by Corollary B.5} \\
& \leq - \gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i=1}^n \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \right)^2 \\
& = - \gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
& = - \gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2
\end{aligned}$$

For T_3 , we align its form with T_1 .

$$\begin{aligned}
T_3 &= \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0}^2 \|x_i\|_2^2 \\
& // \text{since } \forall i \in [n], |y_i| = 1 \\
& = \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2 \\
& = \frac{1}{2} \sigma^2 \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i=1}^n \|x_i\|_2^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
& \leq \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i=1}^n \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
& = \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)}))
\end{aligned}$$

2. Upper bounds of T_2, T_4 . For T_2 , we use linear separability.

$$\begin{aligned}
T_2 &= - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\
& // \text{by Corollary B.5} \\
& \leq - \sum_{j=1}^h \left(\sum_{i \in [n]} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \gamma \|w_j^{(0)}\|_2 \right)^2 \\
& = - \gamma^2 \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2
\end{aligned}$$

$$= -\gamma^2 \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2$$

For T_4 , we align its form with T_3 .

$$\begin{aligned} T_4 &= \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &\quad // \text{since } \forall i \in [n], |y_i| = 1 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j \in [h]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle w_j^{(0)}, x_i \rangle^2 \\ &\leq \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j \in [h]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|w_j^{(0)}\|_2^2 \|x_i\|_2^2 \\ &\leq \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \sum_{i \in [n]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\ &= \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \end{aligned}$$

3. Combine upper bounds of T_1, T_2, T_3, T_4 .

$$\begin{aligned} \dot{\mathcal{L}}^{(0)} &= T_1 + T_2 + T_3 + T_4 \\ &\leq -\gamma^2 \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\ &\quad + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\ &\quad // \text{abbr. } \ell_i := \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\ &= -\gamma^2 \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 \\ &\quad + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \\ &= \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right) \right\} \\ &\because (v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \geq (v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \\ &\therefore \text{When the drift term (negative) still dominates the dynamics, we take } t = 0 \text{ for } (v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2. \\ \dot{\mathcal{L}}^{(0)} &\leq \sum_{j=1}^h \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right) \right\} \end{aligned}$$

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

4. Decompose loss by trapping. If the trapping condition holds, we can decompose the loss $\mathcal{L}^{(0)} = \mathcal{L}_+^{(0)} + \mathcal{L}_-^{(0)}$, where $\mathcal{L}_*^{(0)}$ is only controlled by w_j if $w_j^{(0)} \in \mathcal{S}_*$ ($*$ $\in \{+, -\}$).

$$\begin{aligned} \dot{\mathcal{L}}_*^{(0)} &\leq \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right) \right\} \\ &\leq \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\mathcal{L}_*^{(0)} \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \mathcal{L}_*^{(0)} \right\} \end{aligned}$$

Let $u = 1/\mathcal{L}_*^{(0)}$, $A = \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right]$, $B = \gamma^2$, $C = \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right)$. Then

$$\begin{aligned} -\frac{du}{dt} &\leq -AB + ACu \\ AB \exp(ACt) &\leq \frac{d}{dt} (ue^{ACt}) \\ \frac{B}{C} (\exp(ACt) - 1) &\leq ue^{ACt} - u_0 \\ \frac{B}{C} (\exp(ACt) - 1) + u_0 &\leq ue^{ACt} \\ \frac{B}{C} (1 - \exp(-ACt)) + u_0 e^{-ACt} &\leq u \\ \mathcal{L}_*^{(0)} &\leq \frac{1}{\frac{B}{C} (1 - e^{-ACt}) + \frac{1}{\mathcal{L}_{t=0,*}^{(0)}} e^{-ACt}} \end{aligned}$$

The time limit of the upper bound is

$$\lim_{t \rightarrow \infty} \mathcal{L}_*^{(0)} \leq \frac{C}{B} = \frac{\sigma^2}{2\gamma^2} \left(\max_{i \in [n]} \|x_i\|_2^2 \right) = \frac{1}{2} \frac{\max_{i \in [n]} \|x_i\|_2^2}{\min_{i \in [n]} \|x_i\|_2^2} \sigma^2 \frac{1}{\mu^2}$$

5. Combine clustered losses.

$$\begin{aligned} \mathcal{L}^{(0)} &= \mathcal{L}_-^{(0)} + \mathcal{L}_+^{(0)} \\ &\leq \frac{1}{\frac{B}{C} (1 - e^{-A_+ C t}) + \frac{1}{\mathcal{L}_{t=0,+}^{(0)}} e^{-A_+ C t}} + \frac{1}{\frac{B}{C} (1 - e^{-A_- C t}) + \frac{1}{\mathcal{L}_{t=0,-}^{(0)}} e^{-A_- C t}} \end{aligned}$$

□

Lower bound (type I) proof of Theorem 4.3. 1. Upper bounds for T_1, T_3 . For T_1 , the key idea is $\|x\|_2^2 \geq \langle x, z \rangle^2$ for any unit vector z .

$$\begin{aligned} T_1 &= - \sum_{j=1}^h \left\| \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i \right\|_2^2 \\ &\quad // \text{since } \forall x \in \mathbb{R}^D, z \in \mathbb{S}^{D-1}, \|x\|_2^2 \geq \langle x, z \rangle^2 \\ &\leq - \sum_{j=1}^h \left\langle \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i, z \right\rangle^2 \\ &= - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle x_i, z \rangle \right)^2 \\ &= - \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle x_i, z \rangle \right)^2 \end{aligned}$$

1458 //pick $z = \frac{y_1 x_1}{\|x_1\|_2}$, by Corollary B.5
 1459
 1460 $\leq -\gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i=1}^n \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \right)^2$
 1461
 1462
 1463 $= -\gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2$
 1464
 1465
 1466
 1467
 1468 $= -\gamma^2 \sum_{j=1}^h (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2$
 1469
 1470
 1471

1472 For T_3 , we align its form with T_1 .

1473
 1474 $T_3 = \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2$
 1475
 1476 //since $\forall i \in [n], |y_i| = 1$
 1477
 1478 $= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2$
 1479
 1480
 1481 $= \frac{1}{2} \sigma^2 \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i=1}^n \|x_i\|_2^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)}))$
 1482
 1483
 1484 $\leq \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i=1}^n \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)}))$
 1485
 1486
 1487 $= \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h (v_j^{(0)})^2 \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)}))$
 1488
 1489

1490 **2. Upper bounds of T_2, T_4 .** For T_2 , we use linear separability.

1491
 1492 $T_2 = -\sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2$
 1493
 1494 //by Corollary B.5
 1495
 1496 $\leq -\sum_{j=1}^h \left(\sum_{i \in [n]} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \gamma \|w_j^{(0)}\|_2 \right)^2$
 1497
 1498
 1499
 1500 $= -\gamma^2 \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \right)^2$
 1501
 1502
 1503
 1504 $= -\gamma^2 \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2$
 1505
 1506
 1507

1508 For T_4 , we align its form with T_3 .

1509
 1510 $T_4 = \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2$
 1511
 //since $\forall i \in [n], |y_i| = 1$

$$\begin{aligned}
&= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\
&= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j \in [h]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle w_j^{(0)}, x_i \rangle^2 \\
&\leq \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j \in [h]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|w_j^{(0)}\|_2^2 \|x_i\|_2^2 \\
&\leq \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \sum_{i \in [n]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
&= \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \|w_j^{(0)}\|_2^2 \sum_{i \in \mathcal{I}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)}))
\end{aligned}$$

3. Combine upper bounds of T_1, T_2, T_3, T_4 .

$$\begin{aligned}
\dot{\mathcal{L}}^{(0)} &= T_1 + T_2 + T_3 + T_4 \\
&\leq -\gamma^2 \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \right)^2 \\
&\quad + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \sum_{i \in \mathcal{I}(w_j^{(0)})} \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
&\quad // \text{abbr. } \ell_i := \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \\
&= -\gamma^2 \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right)^2 \\
&\quad + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \\
&= \sum_{j=1}^h \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right) \right\}
\end{aligned}$$

$$\because (v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \geq (v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2$$

\therefore When the drift term (negative) still dominates the dynamics, we take $t = 0$ for $(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2$.

$$\dot{\mathcal{L}}^{(0)} \leq \sum_{j=1}^h \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right) \right\}$$

4. Decompose loss by trapping.

If the trapping condition holds, we can decompose the loss $\mathcal{L}^{(0)} = \mathcal{L}_+^{(0)} + \mathcal{L}_-^{(0)}$, where $\mathcal{L}_*^{(0)}$ is only controlled by w_j if $w_j^{(0)} \in \mathcal{S}_*$ ($*$ \in $\{+, -\}$).

$$\dot{\mathcal{L}}_*^{(0)} \leq \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right)^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \right) \right\}$$

$$\leq \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right] \left\{ -\gamma^2 (\mathcal{L}_*^{(0)})^2 + \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right) \mathcal{L}_*^{(0)} \right\}$$

Let $u = 1/\mathcal{L}_*^{(0)}$, $A = \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \left[(v_{j,t=0}^{(0)})^2 + \|w_{j,t=0}^{(0)}\|_2^2 \right]$, $B = \gamma^2$, $C = \frac{1}{2} \sigma^2 \left(\max_{i \in [n]} \|x_i\|_2^2 \right)$. Then

$$\begin{aligned} -\frac{du}{dt} &\leq -AB + ACu \\ AB \exp(ACt) &\leq \frac{d}{dt}(ue^{ACt}) \\ \frac{B}{C}(\exp(ACt) - 1) &\leq ue^{ACt} - u_0 \\ \frac{B}{C}(\exp(ACt) - 1) + u_0 &\leq ue^{ACt} \\ \frac{B}{C}(1 - \exp(-ACt)) + u_0 e^{-ACt} &\leq u \\ \mathcal{L}_*^{(0)} &\leq \frac{1}{\frac{B}{C}(1 - e^{-ACt}) + \frac{1}{\mathcal{L}_{t=0,*}^{(0)}} e^{-ACt}} \end{aligned}$$

The time limit of the upper bound is

$$\lim_{t \rightarrow \infty} \mathcal{L}_*^{(0)} \leq \frac{C}{B} = \frac{\sigma^2}{2\gamma^2} \left(\max_{i \in [n]} \|x_i\|_2^2 \right) = \frac{1}{2} \frac{\max_{i \in [n]} \|x_i\|_2^2}{\min_{i \in [n]} \|x_i\|_2^2} \sigma^2 \frac{1}{\mu^2}$$

5. Combine clustered losses.

$$\begin{aligned} \mathcal{L}^{(0)} &= \mathcal{L}_-^{(0)} + \mathcal{L}_+^{(0)} \\ &\leq \frac{1}{\frac{B}{C}(1 - e^{-A_+ Ct}) + \frac{1}{\mathcal{L}_{t=0,+}^{(0)}} e^{-A_+ Ct}} + \frac{1}{\frac{B}{C}(1 - e^{-A_- Ct}) + \frac{1}{\mathcal{L}_{t=0,-}^{(0)}} e^{-A_- Ct}} \end{aligned}$$

□

Lower bound (type III) proof of Theorem 4.3. 1. Lower bounds for T_1, T_3 . For T_1 , we use $(\max_{k \in [n]} \|x_k\|_2^2)$.

$$\begin{aligned} T_1 &= - \sum_{j=1}^h \left\| \sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) v_j^{(0)} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} x_i \right\|_2^2 \\ &\quad // \text{abbr. } \ell_i := \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \\ &= - \sum_{j=1}^h \left\| \sum_{i \in \mathcal{I}(w_j^{(0)})} y_i \ell_i v_j^{(0)} x_i \right\|_2^2 \\ &= - \sum_{j=1}^h \left\| \sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i v_j^{(0)} x_i \right\|_2^2 \\ &= - \sum_{j \in [h]} (v_j^{(0)})^2 \left\| \sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i x_i \right\|_2^2 \\ &\geq - \sum_{j \in [h]} (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{I}(w_j^{(0)})} \ell_i \|x_i\|_2 \right)^2 \end{aligned}$$

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

$$\geq - \left(\max_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2$$

For T_3 , we align its form with T_1 .

$$\begin{aligned} T_3 &= \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell_i \sum_{j=1}^h (v_j^{(0)})^2 \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \|x_i\|_2^2 \\ &= \frac{1}{2} \sigma^2 \sum_{j \in [h]} (v_j^{(0)})^2 \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \|x_i\|_2^2 \\ &\geq \frac{1}{2} \sigma^2 \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} (v_j^{(0)})^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 \end{aligned}$$

2. Lower bounds for T_2, T_4 . For T_2 , we use $\langle x, y \rangle \leq \|x\|_2 \|y\|_2$.

$$\begin{aligned} T_2 &= - \sum_{j=1}^h \left(\sum_{i=1}^n y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) \text{relu}((w_j^{(0)})^\top x_i) \right)^2 \\ &= - \sum_{j=1}^h \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} y_i \exp(-y_i f(x_i; W^{(0)}, v^{(0)})) (w_j^{(0)})^\top x_i \right)^2 \\ &= - \sum_{j \in [h]} \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \langle w_j^{(0)}, x_i \rangle \right)^2 \\ &\geq - \sum_{j \in [h]} \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \|w_j^{(0)}\|_2 \|x_i\|_2 \right)^2 \\ &\geq - \left(\max_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 \end{aligned}$$

For T_4 , we align its form with T_2 .

$$\begin{aligned} T_4 &= \frac{1}{2} \sigma^2 \sum_{i=1}^n y_i^2 \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &\quad // \text{since } \forall i \in [n], |y_i| = 1 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \|\text{relu}((W^{(0)})^\top x_i)\|_2^2 \\ &= \frac{1}{2} \sigma^2 \sum_{i=1}^n \ell(y_i, f(x_i; W^{(0)}, v^{(0)})) \sum_{j \in [h]} \mathbb{1}_{(w_j^{(0)})^\top x_i > 0} \langle w_j^{(0)}, x_i \rangle^2 \\ &= \frac{1}{2} \sigma^2 \sum_{j \in [h]} \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \langle w_j^{(0)}, x_i \rangle^2 \\ &\quad // \text{by Lemma B.4} \end{aligned}$$

$$\begin{aligned}
&\geq \frac{1}{2}\sigma^2 \sum_{j \in [h]} \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \mu^2 \|w_j^{(0)}\|_2^2 \|x_i\|_2^2 \\
&= \frac{1}{2}\sigma^2 \mu^2 \sum_{j \in [h]} \|w_j^{(0)}\|_2^2 \sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \|x_i\|_2^2 \\
&\geq \frac{1}{2}\sigma^2 \mu^2 \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)
\end{aligned}$$

3. Combine lower bounds of T_1, T_2, T_3, T_4 .

$$\begin{aligned}
\dot{\mathcal{L}}^{(0)} &= T_1 + T_2 + T_3 + T_4 \\
&\geq - \left(\max_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \left[(v_j^{(0)})^2 + \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 \\
&\quad + \frac{1}{2}\sigma^2 \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \left[(v_j^{(0)})^2 + \mu^2 \|w_j^{(0)}\|_2^2 \right] \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right) \\
&\quad // \text{by balancedness, } \|w_j^{(0)}\|_2^2 = (v_j^{(0)})^2 \\
&\geq -2 \left(\max_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)^2 + \frac{\sigma^2(1+\mu^2)}{2} \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h]} \|w_j^{(0)}\|_2^2 \left(\sum_{i \in \mathcal{S}(w_j^{(0)})} \ell_i \right)
\end{aligned}$$

4. Decompose loss by trapping. If the trapping condition holds, we can decompose the loss

$$\begin{aligned}
\mathcal{L}^{(0)} &= \mathcal{L}_+^{(0)} + \mathcal{L}_-^{(0)}, \text{ where } \mathcal{L}_*^{(0)} \text{ is only controlled by } w_j \text{ if } w_j^{(0)} \in \mathcal{S}_* (* \in \{+, -\}). \\
\dot{\mathcal{L}}_*^{(0)} &\geq -2 \left(\max_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \|w_j^{(0)}\|_2^2 (\mathcal{L}_*^{(0)})^2 + \frac{\sigma^2(1+\mu^2)}{2} \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \|w_j^{(0)}\|_2^2 \mathcal{L}_*^{(0)} \\
&= \left\{ \sum_{j \in [h], w_j^{(0)} \in \mathcal{S}_*} \|w_j^{(0)}\|_2^2 \right\} \cdot \left\{ -2 \left(\max_{k \in [n]} \|x_k\|_2^2 \right) (\mathcal{L}_*^{(0)})^2 + \frac{\sigma^2(1+\mu^2)}{2} \left(\min_{k \in [n]} \|x_k\|_2^2 \right) \mathcal{L}_*^{(0)} \right\}
\end{aligned}$$

The time limit of the loss lower bound is

$$\lim_{t \rightarrow \infty} \mathcal{L}_*^{(0)} \geq \frac{1}{2} \frac{\min_{k \in [n]} \|x_k\|_2^2}{\max_{k \in [n]} \|x_k\|_2^2} \sigma^2 \frac{1+\mu^2}{2}$$

By the previous lower bound proof,

$$\|W^{(0)}\|_F^2 \leq \|W_0^{(0)}\|_F^2 e^{2(\max_{k \in [n]} \|x_k\|_2) \mathcal{L}_0^{(0)} t}$$

Let $u = \frac{1}{\mathcal{L}_*^{(0)}}$, $A = \|W_0^{(0)}\|_F^2$, $\lambda_2 = 2(\max_{k \in [n]} \|x_k\|_2) \mathcal{L}_0^{(0)}$, $B = 2 \max_{k \in [n]} \|x_k\|_2^2$, $C = \frac{\sigma^2(1+\mu^2)}{2} \min_{k \in [n]} \|x_k\|_2^2$. Then consider integrating factor $\exp(AC/\lambda_2 \exp(\lambda_2 t))$.

$$-\frac{d}{dt} u \geq A e^{\lambda_2 t} (-B + C u)$$

$$A B e^{\lambda_2 t} \geq A C e^{\lambda_2 t} u + \frac{d}{dt} u$$

$$A B e^{\lambda_2 t} \exp(AC/\lambda_2 \exp(\lambda_2 t)) \geq A C \exp(AC/\lambda_2 \exp(\lambda_2 t)) e^{\lambda_2 t} u + \exp(AC/\lambda_2 \exp(\lambda_2 t)) \frac{d}{dt} u$$

$$\frac{B}{C} \frac{d}{dt} [\exp(AC/\lambda_2 \exp(\lambda_2 t))] \geq \frac{d}{dt} (u \cdot \exp(AC/\lambda_2 \exp(\lambda_2 t)))$$

$$\begin{aligned}
1728 \quad & \frac{B}{C} [\exp(AC/\lambda_2 \exp(\lambda_2 t)) - \exp(AC/\lambda_2)] \geq u \cdot \exp(AC/\lambda_2 \exp(\lambda_2 t)) - u_0 \cdot \exp(AC/\lambda_2) \\
1729 \quad & \\
1730 \quad & \frac{B}{C} [1 - \exp(AC/\lambda_2(1 - \exp(\lambda_2 t)))] \geq u - u_0 \cdot \exp(AC/\lambda_2(1 - \exp(\lambda_2 t))) \\
1731 \quad & \\
1732 \quad & \mathcal{L}_*^{(0)} \geq \frac{1}{\frac{1}{\mathcal{L}_{*,t=0}^{(0)}} e^{AC/\lambda_2(1-\exp(\lambda_2 t))} + \frac{B}{C} [1 - e^{AC/\lambda_2(1-\exp(\lambda_2 t))}]}
\end{aligned}$$

5. Combine clustered losses.

$$\begin{aligned}
1736 \quad & \mathcal{L}^{(0)} = \mathcal{L}_-^{(0)} + \mathcal{L}_+^{(0)} \\
1737 \quad & \\
1738 \quad & \geq \frac{1}{\frac{1}{\mathcal{L}_{+,t=0}^{(0)}} e^{AC/\lambda_2(1-\exp(\lambda_2 t))} + \frac{B}{C} [1 - e^{AC/\lambda_2(1-\exp(\lambda_2 t))}]} + \frac{1}{\frac{1}{\mathcal{L}_{-,t=0}^{(0)}} e^{AC/\lambda_2(1-\exp(\lambda_2 t))} + \frac{B}{C} [1 - e^{AC/\lambda_2(1-\exp(\lambda_2 t))}]} \\
1739 \quad & \\
1740 \quad & \square
\end{aligned}$$

D.3 PRIVACY BUDGET ALLOCATION

1742 *Proof of Theorem 5.1.* For any $j \in [h]$, with probability $1 - \rho$, its initial absolute value is bounded
1743 by

$$1744 \quad |v_j| \leq \sqrt{2\beta^2 \ln(2/\rho)} \quad (55)$$

1745 Then with probability $(1 - \rho)^h$, the maximum worst initial value is bounded by

$$1746 \quad \max_{j \in [h]} (c_j \cdot v_j) \leq \sqrt{\beta^2 \ln(2/\rho)} \quad (56)$$

1747 where we define c_j by $w_j \in S_{c_j}$. The approximate DP-LP dynamics is

$$1748 \quad \dot{v}_j = \sum_{i=1}^n y_i \ell_i \text{relu}(w_j^\top x_i) \quad (57)$$

1749 Say $w_j \in S_c$ for some $c \in \{-1, 1\}$, then during DP-LP, when $\text{sign}(v_j(T)) = \text{sign}(v_j(0))$,

$$1750 \quad |v_j(T) - v_j(0)| = \int_0^T \sum_{y_i=c} \ell_i \text{relu}(w_j^\top x_i) dt \quad (58)$$

$$1751 \quad \geq \min_{y_i=c} |\text{relu}(w_j^\top x_i)| \int_0^T \mathcal{L}_c(t) dt \quad (59)$$

$$1752 \quad // \text{by Theorem 4.2} \quad (60)$$

$$1753 \quad \geq \min_{y_i=c} \text{relu}(w_j^\top x_i) \frac{\frac{1}{2} \sigma^2 \left\{ \sum_{y_i=c} \|\text{relu}(W^\top x_i)\|_2^{-2} \right\}^{-1}}{\sum_{w_j \in S_c} [\max_{y_i=c} w_j^\top x_i]^2} \quad (61)$$

$$1754 \quad = \frac{1}{2} \sigma^2 \frac{\min_{y_i=c} \text{relu}(w_j^\top x_i)}{\sum_{w_j \in S_c} [\max_{y_i=c} w_j^\top x_i]^2} \left\{ \sum_{y_i=c} \|\text{relu}(W^\top x_i)\|_2^{-2} \right\}^{-1} \quad (62)$$

$$1755 \quad = \frac{1}{2} \sigma^2 Q \quad (63)$$

1756 where we define a constant Q to describe the pre-training quality. If the pre-trained features are
1757 better, Q becomes larger. To mitigate the feature distortion, we need $c \cdot v_j > 0$, then the necessary
1758 DP-LP run-time is

$$1759 \quad \Delta t \propto \frac{\sigma^2}{Q} \sqrt{\beta^2 \ln(2/\rho)} \propto \frac{\sigma^2}{Q} \sqrt{\ln(2/\rho)} \quad (64)$$

1760 where we ignore β as it is typically pre-determined in real implementations (e.g. the Linear layers
1761 in PyTorch). \square

1778 E APPENDIX: THEORY WITHOUT APPROXIMATION

1779 For convenience, we use different notations for the data input dimension $d = d_x$ and the backbone
1780 weight matrix $B = W^\top$ in the following proofs.

E.1 ITÔ'S FORMULA AND ITS CONSEQUENCES

We denote $M_{m,n}(\mathbb{R})$ as the space of m-by-n real matrices.

Theorem E.1 (Itô's formula). *Let X_t be a \mathbb{R}^n -valued Itô process satisfying the stochastic differential equation $\partial X_t = A_1(t, X_t)\partial t + A_2(t, X_t)\partial W_t$ with $A_1(t, X_t)$ being \mathbb{R}^n -valued, $A_2(t, X_t)$ being $M_{m,n}(\mathbb{R})$ -valued, and W_t being a standard n-dimensional brownian motion. Let $f : [0, \infty) \times \mathbb{R}^n \rightarrow \mathbb{R}$ be a function with continuous partial derivatives. Then $Y_t := f(t, X_t)$ is also an Itô process, and its stochastic differential equation is*

$$\partial Y_t = \frac{\partial f(t, X_t)}{\partial t} \partial t + \langle \nabla f(t, X_t), A_1(t, X_t)\partial t + A_2(t, X_t)\partial W_t \rangle + \frac{1}{2} \langle A_2(t, X_t)\partial W_t, H_f A_2(t, X_t)\partial W_t \rangle \quad (65)$$

where H_f is the Hessian matrix of f over X_t defined as $(H_f)_{ij} = \frac{\partial^2 f}{\partial (X_t)_i \partial (X_t)_j}$ and $(X_t)_i$ denotes the i -th entry of random vector X_t .

Corollary E.2 (Loss dynamics during linear probing). *During linear probing (Equation equation 90), the stochastic differential equation describing the loss dynamics is*

$$\partial \mathcal{L}_{\text{lp}} = -(B_0^T v - X^T Y)^T B_0^T B_0 (B_0^T v - X^T Y) \partial t + \sqrt{2\sigma^2} (B_0^T v - X^T Y)^T B_0^T \partial W_t + h\sigma^2 \partial t. \quad (66)$$

Proof of Corollary E.2. By Itô's formula (Equation equation E.1), the loss dynamics is

$$\partial \mathcal{L}_{\text{lp}} = \partial \frac{1}{2} \|X B_0^T v - Y\|^2 \quad (67)$$

$$= (X B_0^T v - Y)^T X B_0^T \partial v + \frac{1}{2} (\partial v)^T B_0 X^T X B_0^T (\partial v) \quad (68)$$

$$= (X B_0^T v - Y)^T X B_0^T \partial v + \frac{1}{2} (\partial v)^T (\partial v) \quad (69)$$

$$// \text{by Definition E.5} \quad (70)$$

$$= (X B_0^T v - Y)^T X B_0^T [-B_0 X^T (X B_0^T v - Y) \partial t + \sqrt{2\sigma^2} \partial W_t] + h\sigma^2 \partial t \quad (71)$$

$$= (B_0^T v - X^T Y)^T B_0^T [-B_0 (B_0^T v - X^T Y) \partial t + \sqrt{2\sigma^2} \partial W_t] + h\sigma^2 \partial t \quad (72)$$

$$= -(B_0^T v - X^T Y)^T B_0^T B_0 (B_0^T v - X^T Y) \partial t + \sqrt{2\sigma^2} (B_0^T v - X^T Y)^T B_0^T \partial W_t + h\sigma^2 \partial t \quad (73)$$

□

Corollary E.3 (Loss dynamics during fine-tuning). *During fine-tuning (Equation equation 91), the stochastic differential equation describing the loss dynamics is*

$$\begin{aligned} \partial \mathcal{L}_{\text{ft}} = & -(B^T v - X^T Y)^T B^T B (B^T v - X^T Y) \partial t + (B^T v - X^T Y)^T B^T \sqrt{2\sigma^2} \partial W_t \\ & - (B^T v - X^T Y)^T (B^T v - X^T Y) v^T v \partial t + (B^T v - X^T Y)^T (\sqrt{2\sigma^2} \partial W_t) v \\ & + \sigma^2 \|B\|_F^2 \partial t + \sigma^2 d \|v\|_2^2 \partial t. \end{aligned} \quad (74)$$

where we use ∂ as the differential sign and use d as the data input dimension.

Proof of Corollary E.3. Similar to Corollary E.2, we use Itô's formula (Equation E.1), the loss dynamics of fine-tuning is

$$\partial \mathcal{L}_{\text{ft}} = \partial \frac{1}{2} \|X B^T v - Y\|^2 \quad (75)$$

$$= \frac{1}{2} \langle \nabla_v \|X B^T v - Y\|^2, \partial v \rangle + \frac{1}{2} \langle \nabla_B \|X B^T v - Y\|^2, \text{vec}(\partial B) \rangle \quad (76)$$

$$+ \frac{1}{4} (\partial v)^T H_{\|X B^T v - Y\|^2} (\partial v) + \frac{1}{4} [\text{vec}(\partial B)]^T H_{\|X B^T v - Y\|^2} \text{vec}(\partial B) \quad (77)$$

$$= (X B^T v - Y)^T X B^T \partial v + (X B^T v - Y)^T X (\partial B)^T v \quad (78)$$

$$+ \frac{1}{2} (\partial v)^T B X^T X B^T (\partial v) + \frac{1}{2} [\text{vec}(\partial B)]^T \begin{bmatrix} v_1 \\ 0 \\ \vdots \\ v_h \end{bmatrix} \underbrace{[v_1 \quad 0 \quad \cdots \quad v_h]}_{d \times h} \text{vec}(\partial B) \quad (79)$$

$$1836 \quad = - (B^T v - X^T Y)^T B^T B (B^T v - X^T Y) \partial t + (B^T v - X^T Y)^T B^T \sqrt{2\sigma^2} \partial W_t \quad (80)$$

$$1837 \quad - (B^T v - X^T Y)^T (B^T v - X^T Y) v^T v \partial t + (B^T v - X^T Y)^T (\sqrt{2\sigma^2} \partial W_t') v \quad (81)$$

$$1838 \quad + \sigma^2 \text{trace}(B B^T) \partial t + \sigma^2 d \|v\|_2^2 \partial t \quad (82)$$

$$1840 \quad = - (B^T v - X^T Y)^T B^T B (B^T v - X^T Y) \partial t + (B^T v - X^T Y)^T B^T \sqrt{2\sigma^2} \partial W_t \quad (83)$$

$$1841 \quad - (B^T v - X^T Y)^T (B^T v - X^T Y) v^T v \partial t + (B^T v - X^T Y)^T (\sqrt{2\sigma^2} \partial W_t') v \quad (84)$$

$$1842 \quad + \sigma^2 \|B\|_F^2 \partial t + \sigma^2 d \|v\|_2^2 \partial t \quad (85)$$

1843
1844
1845 □

1846 *Remark E.4* (Noise effects on linear networks). In the loss dynamics of fine-tuning (Corollary E.3),
1847 the noise induced deterministic terms

$$1848 \quad \sigma^2 (\|B\|_F^2 + d \|v\|_2^2) \partial t$$

1849 does not explicitly depend on the linear head size h . We do a sanity check for this result in a
1850 discretized setting (so that we skip Itô's lemma and stochastic calculus). Say we inject noise ΔB
1851 to B , where ΔB is a $h \times d$ -matrix, and its entries are independent and follow Gaussian distribution
1852 $\mathcal{N}(0, \sigma)$. Then the expectation of the perturbed loss is:

$$1853 \quad \mathbb{E}[\mathcal{L}] = \frac{1}{2} \mathbb{E}[\|X(B + \Delta B)^T v - Y\|^2] \quad (86)$$

$$1854 \quad = \frac{1}{2} \|X B^T v - Y\|^2 + \mathbb{E}[(X B^T v - Y)^T X (\Delta B)^T v] + \frac{1}{2} \mathbb{E}[v^T \Delta B (\Delta B)^T v] \quad (87)$$

$$1855 \quad = \frac{1}{2} \|X B^T v - Y\|^2 + \frac{1}{2} \mathbb{E}[v^T \Delta B (\Delta B)^T v] \quad (88)$$

$$1856 \quad = \frac{1}{2} \|X B^T v - Y\|^2 + \frac{1}{2} \sigma^2 \cdot d \cdot \|v\|^2 \quad (89)$$

1857
1858
1859
1860
1861
1862 As a result, we find that, in the discrete updates, the noise induced deterministic terms does not
1863 explicitly depend on the linear head size h either. So our findings in the continuous case matches
1864 the discrete case.

1865 E.2 MODIFIED LANGEVIN DIFFUSION

1866
1867 **Definition E.5** (Langevin diffusion for linear probing). Let Q_t be the standard h -dimensional Brownian
1868 motion. Then the Langevin diffusion for linear probing is defined by the following stochastic
1869 differential equation:

$$1870 \quad \partial v = -\nabla_v \mathcal{L}(v, B_0) \partial t + \sqrt{2\sigma^2} \partial Q_t$$

$$1871 \quad = -B_0 X^T (X B_0^T v - Y) \partial t + \sqrt{2\sigma^2} \partial Q_t. \quad (90)$$

1872 Here we use “ ∂ ” as the differential notation.

1873
1874 **Definition E.6** (Langevin diffusion for fine-tuning). Let Q_t be the standard h -dimensional brownian
1875 motion and Q'_t be a matrix whose entries are standard and independent brownian motions. Then we
1876 define the Langevin diffusion for fine-tuning a two-layer linear network as

$$1877 \quad \partial v = -\nabla_v \mathcal{L}(v, B) \partial t + \sqrt{2\sigma^2} \partial Q_t$$

$$1878 \quad = -B X^T (X B^T v - Y) \partial t + \sqrt{2\sigma^2} \partial Q_t \quad (91)$$

$$1879 \quad \partial B = -\nabla_B \mathcal{L}(v, B) \partial t + \sqrt{2\sigma^2} \partial Q'_t$$

$$1880 \quad = -v (X B^T v - Y)^T X \partial t + \sqrt{2\sigma^2} \partial Q'_t.$$

1881
1882 Here we introduce an assumption based on random initialization. It describes a common phenom-
1883 enon in differential privacy deployment: the loss might not converge if the privacy mechanism
1884 perturbs the gradients too much (Ponomareva et al., 2023). To ensure that DP-SGD works for full
1885 fine-tuning, we assume that the noise scale (or variance) in the privacy mechanism is upper bounded
1886 by a constant.

Assumption E.7 (Upper bounded noise scale). Let $\beta > \frac{-\|X^T Y\| + \sqrt{\|X^T Y\|^2 + 4(1+d_x)\|X^T Y\| + 4d_x}}{2h}$. Then we assume that the noise scale $\sigma > 0$ we add for privacy in the fine-tuning process is upper-bounded by

$$\sigma^2 < \min \left\{ \frac{h\beta + \|B_0 X^T Y\|^2}{2h}, \frac{h\beta - 1}{\sqrt{2}(1+d)}, \frac{1}{1 + \sqrt{2}(1+d)} \left[\frac{h\beta(h\beta + \|X^T Y\|^2)}{(1+d)\|X^T Y\| + d} - 1 \right] \right\}. \quad (92)$$

Equation (18) upper monotonically decreases in time if Assumption E.7 also holds.

To understand the properties of a dynamics analysis problem, it can be useful to identify *invariants*, or functions whose output is conserved during optimization. Such conservation laws can be seen as a "weaker" form of implicit bias, helping to elucidate which properties (e.g., sparsity, low-rank) are preferred by the optimization dynamics among a potentially infinite set of minimizers (Marcotte et al., 2023). To prove the convergence of our optimization, we study the *imbalance matrix*, an invariant for multi-layer linear networks that has previously been studied in the context of gradient flows (but not Langevin dynamics, to the best of our knowledge).

Definition E.8 (Imbalance matrix). For a two-layer linear network, we define the imbalance matrix as

$$D := vv^T - BB^T. \quad (93)$$

Prior work on gradient flows has found that the imbalance matrix remains invariant over the evolution of gradient flows modeling gradient descent (Arora et al., 2018; Du et al., 2018; Marcotte et al., 2023). This property can be used to derive tight convergence bounds (Min et al., 2021; 2023a). However, a similar analysis has not materialized for Langevin diffusion models of DP-GD.

We observe that prior work on Langevin diffusion to analyze private optimization has implicitly assumed that the sensitivity of each layer in a neural network is the same (Ganesh et al., 2023b; Ye et al., 2023b). Hence, they fix a uniform noise scale for every parameter of the network. Under these conditions, we show that, when we ignore the sensitivity of each layer and use a uniform noise scale σ , the imbalance matrix is *not* invariant in expectation, unlike in (noise-free) gradient flow (Arora et al., 2018; Du et al., 2018; Marcotte et al., 2023); that is, its derivative over time is nonzero. This complicates the use of the imbalance matrix for theoretical analysis (Ye & Du, 2021).

Lemma E.9 (Imbalance matrix in fine-tuning). *During fine-tuning (Equation (91)), the derivative of the imbalance matrix D in Definition E.8 is*

$$\frac{\partial}{\partial t} \mathbb{E}[D] = (1-d)\sigma^2 I_{h \times h}, \quad (94)$$

where d is the dimension of data inputs ($B \in \mathbb{R}^{h \times d}$).

Our main observation is that by modeling differences in sensitivity of different layers, we can recover the invariance property of the imbalance matrix. The following proposition characterizes the sensitivity of the linear head and the feature extractor, and illustrates why they have differing sensitivities at initialization.

Proposition E.10. *We assume that the training dataset $\mathcal{D} = (X, Y)$ is normalized such that $X^T X = I_{d \times d}$, $\|Y\|_2 = 1$. We initialize the linear head by $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$ and $\beta = h/\sqrt{d}$. At the initialization of full fine-tuning, the linear head v has a greater layer sensitivity (Béthune et al., 2024) than the feature extractor B :*

$$\Delta(\nabla_v \mathcal{L}(v_0, B_0)) = \Theta \left(\sqrt{d} \cdot \Delta(\nabla_B \mathcal{L}(v_0, B_0)) \right) \quad (95)$$

Based on this observation, we propose a modified version of Langevin diffusion for full fine-tuning, which accounts for layer-wise sensitivity. With this modified definition, the imbalance matrix is again invariant in expectation.

Definition E.11 (Modified Langevin diffusion for fine-tuning). Let Q_t be the standard h -dimensional brownian motion. Let Q_t^i be a $h \times d$ matrix whose entries are standard and independent brownian motions. Then we define the modified Langevin diffusion for fine-tuning a two-layer

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

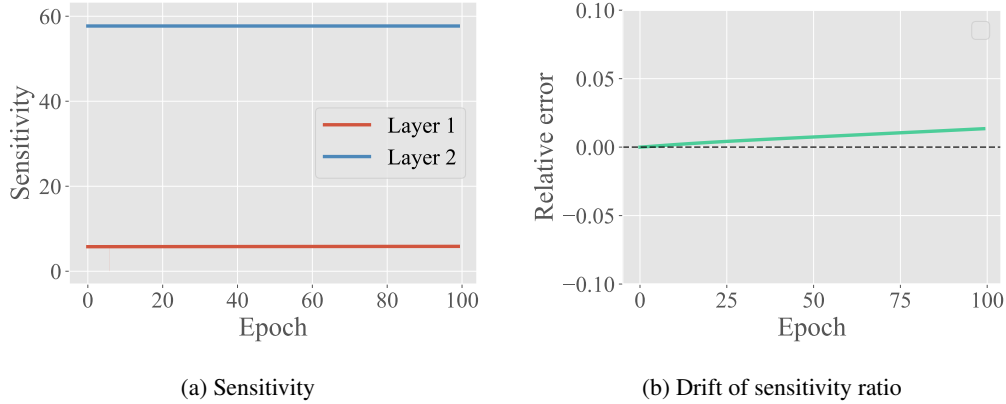


Figure 6: Evaluation of layer-wise sensitivity when running DP-GD on 2-layer linear networks and synthetic data (Béthune et al., 2024). We initialize the network parameter according to Proposition E.10. We take average on 10^4 random seeds with standard error smaller than 10^{-3} .

linear network as

$$\begin{aligned}
 \partial v &= -\nabla_v \mathcal{L}(v, B) \partial t + \sqrt{2\sigma^2 d} \partial Q_t \\
 &= -BX^T X(B^T v - X^T Y) \partial t + \sqrt{2\sigma^2 d} \partial Q_t \\
 \partial B &= -\nabla_B \mathcal{L}(v, B) \partial t + \sqrt{2\sigma^2} \partial Q'_t \\
 &= -v(XB^T v - Y)^T X \partial t + \sqrt{2\sigma^2} \partial Q'_t.
 \end{aligned} \tag{96}$$

The only difference between this diffusion and Equation (91) is the additional factor of \sqrt{d} , shown in red, reflecting the fact that the linear head has greater function sensitivity than the feature extractor.

E.3 LINEAR PROBING LOSS UPPER BOUND

The main idea of the proofs for convergence is to replace gradient terms with loss terms. By doing so, we obtain inequalities containing only loss terms and some other constants.

For the linear probing setting, we first show the strong convexity of the loss function. Then we can use the Lojasiewicz inequality to replace gradient terms with the loss terms.

Lemma E.12 ((Strong) convexity of linear probing phase). *The empirical risk $\mathcal{L} = \frac{1}{2} \sum_{i=1}^n \ell(f(x_i), y_i)$ is 1-strongly convex.*

Lemma E.13 (Initial loss before linear probing). *If we initialize the linear head by $v_{t=0} \sim \mathcal{N}(0, \beta I_{h \times h})$, then the expected empirical risk before linear probing is*

$$\mathbb{E}[\mathcal{L}_0] = \frac{1}{2}(h\beta + \|Y\|^2) \tag{97}$$

Proof of Lemma E.13. We initialize the linear head with a Gaussian distribution $\mathcal{N}(0, \beta I_{h \times h})$. So the expected initial loss is:

$$\mathbb{E}[\mathcal{L}_0] = \frac{1}{2} \mathbb{E}[\|XB_0^T v_0 - Y\|^2] \tag{98}$$

$$= \frac{1}{2} \mathbb{E}[v_0^T B_0 X^T X B_0^T v_0 + Y^T Y - 2Y^T X B_0^T v_0] \tag{99}$$

$$= \frac{1}{2} \mathbb{E}[v_0^T B_0 B_0^T v_0 + Y^T Y] \tag{100}$$

//we assumed in section 3.1 that B_0 has orthogonal rows $\tag{101}$

$$= \frac{1}{2} \mathbb{E}[v_0^T v_0 + Y^T Y] \tag{102}$$

$$\text{//by } v_{t=0} \sim \mathcal{N}(0, I_{h \times h}) \quad (103)$$

$$= \frac{1}{2}(h\beta + \|Y\|^2) \quad (104)$$

□

Theorem E.14 (Expected loss upper bound of linear probing). *The expected empirical risk in linear probing is upper bounded by*

$$\mathbb{E}[\mathcal{L}_{\text{lp}}(t)] \leq e^{-t}\mathbb{E}[\mathcal{L}_0] + (1 - e^{-t})(\gamma + h\sigma^2) \quad (105)$$

Proof of Theorem 4.4. By Lemma E.12, \mathcal{L} is 1-strongly convex, we have the Lojasiewicz inequality. Here we abuse the notation \mathcal{L} and consider it as a function of the linear head v because we fix B_0 in the linear probing process.

$$\mathcal{L}(v) - \{\min_v \mathcal{L}\} \leq \frac{1}{2}\|\nabla_v \mathcal{L}(v)\|_2^2 \quad (106)$$

For simplicity, we denote $\mathbb{E}[\mathcal{L}] := \hat{\mathcal{L}}$. Consider the Langevin diffusion in Equation equation 90 when $\mathcal{L}(v) - \{\min_v \mathcal{L}\} - h\sigma^2 > 0$, by Corollary E.2:

$$\partial \mathcal{L}(v) = \langle \nabla_v \mathcal{L}(v), -\nabla_v \mathcal{L}(v) \rangle dt + \sqrt{2\sigma^2} \partial W_t + h\sigma^2 \partial t \quad (107)$$

$$\partial \mathcal{L}(v) \leq -\|\nabla_v \mathcal{L}(v)\|_2^2 dt + \langle \nabla_v \mathcal{L}(v), \sqrt{2\sigma^2} \partial W_t \rangle + h\sigma^2 \partial t \quad (108)$$

$$\text{//By Lojasiewicz inequality} \quad (109)$$

$$\partial \mathcal{L}(v) \leq (-\mathcal{L}(v) + \{\min_v \mathcal{L}\}) \partial t + \langle \nabla_v \mathcal{L}(v), \sqrt{2\sigma^2} \partial W_t \rangle + h\sigma^2 \partial t \quad (110)$$

$$\partial(\mathbb{E}[\mathcal{L}(v)] - \{\min_v \mathcal{L}\} - h\sigma^2) \leq -(\mathbb{E}[\mathcal{L}(v)] - \{\min_v \mathcal{L}\}) \partial t + h\sigma^2 \partial t \quad (111)$$

$$\partial(\hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2) \leq -(\hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2) \partial t \quad (112)$$

$$\text{//When } \hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2 > 0 \quad (113)$$

$$\partial \ln |\hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2| \leq -1 \partial t \quad (114)$$

$$\ln |\hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2| \leq \ln |\widehat{\mathcal{L}(v_0)} - \{\min_v \mathcal{L}\} - h\sigma^2| - t \quad (115)$$

$$\hat{\mathcal{L}} - \{\min_v \mathcal{L}\} - h\sigma^2 \leq e^{-t}(\widehat{\mathcal{L}(v_0)} - \{\min_v \mathcal{L}\} - h\sigma^2) \quad (116)$$

$$\hat{\mathcal{L}} \leq e^{-t}(\widehat{\mathcal{L}(v_0)} - \{\min_v \mathcal{L}\} - h\sigma^2) + \{\min_v \mathcal{L}\} + h\sigma^2 \quad (117)$$

$$\hat{\mathcal{L}} \leq e^{-t}\widehat{\mathcal{L}(v_0)} + (1 - e^{-t})(\{\min_v \mathcal{L}\} + h\sigma^2) \quad (118)$$

$$\hat{\mathcal{L}} \leq e^{-t}\widehat{\mathcal{L}(v_0)} + (1 - e^{-t})(\gamma + h\sigma^2) \quad (119)$$

□

When we substitute the initial loss $\mathcal{L}(v_0)$ with the hyper-parameters we use in the random initialization, we obtain the following corollary.

Corollary E.15 (Expected loss upper bound of linear probing from random initialization). *If we initialize the linear head by $v_{t=0} \sim \mathcal{N}(0, I_{h \times h})$, then the expected loss is upper bounded by*

$$\mathbb{E}[\mathcal{L}_{\text{lp}}(t)] \leq \frac{1}{2}(h\beta + \|Y\|^2)e^{-t} + (1 - e^{-t})(\gamma + h\sigma^2) \quad (120)$$

Proof of Corollary E.15. The result is immediate when we combine Lemma E.13 and Theorem 4.4.

□

E.4 IMBALANCE MATRIX FROM LINEAR PROBING

In the convergence analysis of fine-tuning, we eliminate variables and simplify the Langevin dynamics by the imbalance matrix. In this part, we characterize how the imbalance matrix changes in the linear probing phase. The following results will later help us analyze LP-FT.

Lemma E.16 (Eigenvalues of imbalance matrix at the beginning of fine-tuning). *During the linear probing phase (Equation equation 90), for the imbalance matrix defined in Definition E.8,*

1. *the minimum eigenvalue of the imbalance matrix is always -1 ;*
2. *other eigenvalues evolve in this way:*

$$\mathbb{E}[\lambda] = \mathbb{E}[\|v\|_2^2] - 1 \geq -1 \quad (121)$$

Proof of Lemma E.16. Consider any eigenpair (λ, u) of matrix D , we have

$$Du = \lambda u \quad (122)$$

$$(vv^T - B_0 B_0^T)u = \lambda u \quad (123)$$

$$(vv^T - I_{h \times h})u = \lambda u \quad (124)$$

$$(v^T u)v = (\lambda + 1)u \quad (125)$$

$$(126)$$

We can take any $u \perp v$ and $(u, -1)$ is an eigenpair of D . So -1 is always an eigenvalue of D . We need to discuss two different cases here:

1. If $\lambda = -1$, we only know that $u \perp v$.
2. If $\lambda \neq -1$, then v and u are parallel. Say $u = \alpha v$, then

$$u = \frac{v^T u}{\lambda + 1} v \quad (127)$$

$$\alpha v = \frac{\alpha \|v\|_2^2}{\lambda + 1} v \quad (128)$$

$$\implies \lambda = \|v\|_2^2 - 1 \geq -1 \quad (129)$$

□

Proposition E.17 (Expected eigenvalue of imbalance matrix at the beginning of fine-tuning). *Say we run linear probing for time t . If we initialize the linear head by $v_{t=0} \sim \mathcal{N}(0, I_{h \times h})$, then for the imbalance matrix defined in Definition E.8, we have*

$$\mathbb{E}[\|v\|^2] = h\beta e^{-2t} + 2\|B_0 X^T Y\|^2 (e^{-t} - e^{-2t}) + (\|B_0 X^T Y\|^2 + h\sigma^2)(1 - e^{-2t}) \quad (130)$$

throughout the linear probing process. Then by Lemma E.16, for those eigenvalues not equal to -1 , we have

$$\mathbb{E}[\lambda] = \mathbb{E}[\|v\|_2^2] - 1 = h\beta e^{-2t} + 2\|B_0 X^T Y\|^2 (e^{-t} - e^{-2t}) + (\|B_0 X^T Y\|^2 + h\sigma^2)(1 - e^{-2t}) - 1 \quad (131)$$

at the beginning of fine-tuning after linear probing.

Proof of Proposition E.17. By Equation equation 90, the Langevin diffusion of linear probing is:

$$\partial v = -B_0 X^T (X B_0^T v - Y) \partial t + \sqrt{2\sigma^2} \partial W_t = -v \partial t + B_0 X^T Y \partial t + \sqrt{2\sigma^2} \partial W_t \quad (132)$$

We consider the evolution of $v^T v$: by Itô's formula (Equation equation E.1)

$$\partial v^T v = 2v^T \partial v + (\partial v)^T I_h (\partial v) \quad (133)$$

$$\partial v^T v = -2v^T (v - B_0 X^T Y) \partial t + 2v^T \sqrt{2\sigma^2} \partial W_t + 2h\sigma^2 \partial t \quad (134)$$

$$\partial v^T v = (-2v^T v + 2v^T B_0 X^T Y) \partial t + 2v^T \sqrt{2\sigma^2} \partial W_t + 2h\sigma^2 \partial t \quad (135)$$

To solve the above equation, we need to solve the dynamics of $v^T B_0 X^T Y$:

$$\partial Y^T X B_0^T v = -Y^T X B_0^T (v - B_0 X^T Y) \partial t + \sqrt{2\sigma^2} \partial W_t \quad (136)$$

$$\partial \mathbb{E}[Y^T X B_0^T v] = -\mathbb{E}[Y^T X B_0^T v] dt + \|B_0 X^T Y\|^2 \partial t \quad (137)$$

$$\frac{\partial}{\partial t} \mathbb{E}[Y^T X B_0^T v - \|B_0 X^T Y\|^2] = -\mathbb{E}[Y^T X B_0^T v - \|B_0 X^T Y\|^2] \quad (138)$$

$$\frac{\partial}{\partial t} \ln \mathbb{E}[Y^T X B_0^T v - \|B_0 X^T Y\|^2] = -1 \quad (139)$$

$$\mathbb{E}[Y^T X B_0^T v_t - \|B_0 X^T Y\|^2] = \mathbb{E}[Y^T X B_0^T v_0 - \|B_0 X^T Y\|^2] \cdot \exp(-t) \quad (140)$$

When we initialize the linear head by $v_{t=0} \sim \mathcal{N}(0, I_{h \times h})$, we have $\mathbb{E}[Y^T X B_0^T v_0] = 0$. Then

$$\mathbb{E}[Y^T X B_0^T v_t - \|B_0 X^T Y\|^2] = \mathbb{E}[Y^T X B_0^T v_0 - \|B_0 X^T Y\|^2] \cdot \exp(-t) \quad (141)$$

$$\mathbb{E}[\|B_0 X^T Y\|^2 - Y^T X B_0^T v_t] = \mathbb{E}[\|B_0 X^T Y\|^2 - Y^T X B_0^T v_0] \cdot \exp(-t) \quad (142)$$

So we can rewrite Equation equation 135 as:

$$\partial \mathbb{E}[\|v\|^2] = (-2\mathbb{E}[\|v\|^2] + 2\mathbb{E}[v^T B_0 X^T Y]) \partial t + 2h\sigma^2 \partial t \quad (143)$$

$$\partial \mathbb{E}[\|v\|^2] = (-2\mathbb{E}[\|v\|^2] + 2(\mathbb{E}[\|B_0 X^T Y\|^2 - Y^T X B_0^T v_0] \cdot \exp(-t) + \|B_0 X^T Y\|^2)) \partial t + 2h\sigma^2 \partial t \quad (144)$$

$$\frac{1}{2} \frac{\partial}{\partial t} \mathbb{E}[\|v\|^2] = -\mathbb{E}[\|v\|^2] + \mathbb{E}[\|B_0 X^T Y\|^2 - Y^T X B_0^T v_0] \cdot \exp(-t) + (\|B_0 X^T Y\|^2 + h\sigma^2) \quad (145)$$

Let $a_1 = \mathbb{E}[\|B_0 X^T Y\|^2 - Y^T X B_0^T v_0]$, $a_2 = \|B_0 X^T Y\|^2 + h\sigma^2$, $f(t) = \mathbb{E}[\|v\|^2]$ and rewrite the above equation:

$$\frac{1}{2} f'(t) + f(t) = a_1 e^{-t} + a_2 \quad (146)$$

$$f'(t) + 2f(t) = 2a_1 e^{-t} + 2a_2 \quad (147)$$

$$e^{2t} f'(t) + 2e^{2t} f(t) = 2a_1 e^t + 2a_2 e^{2t} \quad (148)$$

$$e^{2t} f(t) \Big|_0^t = (2a_1 e^t + a_2 e^{2t}) \Big|_0^t \quad (149)$$

$$e^{2t} f(t) = f(0) + 2a_1(e^t - 1) + a_2(e^{2t} - 1) \quad (150)$$

$$f(t) = f(0)e^{-2t} + 2a_1(e^{-t} - e^{-2t}) + a_2(1 - e^{-2t}) \quad (151)$$

Since we initialize the linear head by $v_{t=0} \sim \mathcal{N}(0, I_{h \times h})$, we have $f(0) = h\beta$ and $a_1 = \|B_0 X^T Y\|^2$. \square

Lemma E.18 (Imbalance matrix in fine-tuning). *During fine-tuning (Equation equation 91), the imbalance matrix D in Definition E.8 evolves as*

$$\frac{\partial}{\partial t} \mathbb{E}[D] = (1 - d)\sigma^2 I_{h \times h} \quad (152)$$

where d is the dimension of data inputs ($B \in \mathbb{R}^{h \times d}$).

Proof of Lemma E.9. We prove this lemma by analyzing the infinitesimal generator A of imbalance matrix D at any time:

$$A(D)_{ij} := \lim_{t \downarrow 0} \frac{\mathbb{E}^D[(D(t))_{ij}] - (D)_{ij}}{t} \quad (153)$$

$$= 0 + \sigma^2 \sum_{i' \in [h]} \sum_{j' \in [h]} \mathbf{1}[i' = j' = i = j] \quad (154)$$

$$- \sigma^2 \sum_{i' \in [h], j' \in [d]} \sum_{i'' \in [h], j'' \in [d]} \mathbf{1}[i' = i'' = i = j \text{ and } j' = j''] \quad (155)$$

2160 the generator is zero for $i \neq j$. So we can just consider the case where $i = j$.

$$2161 \quad A(D)_{ii} = \sigma^2 \sum_{i' \in [h]} \sum_{j' \in [h]} \mathbf{1}[i' = j' = i] \quad (156)$$

$$2162 \quad - \sigma^2 \sum_{i' \in [h], j' \in [d]} \sum_{i'' \in [h], j'' \in [d]} \mathbf{1}[i' = i'' = i \text{ and } j' = j''] \quad (157)$$

$$2163 \quad = (1 - d)\sigma^2 \quad (158)$$

2164 □

2165 **Lemma E.19** (Monotonic eigenvalue of imbalance matrix in fine-tuning). *Denote D_{lp} as the imbalance matrix right after linear probing phase. All eigenvalues of the imbalance matrix are decreasing in expectation during fine-tuning. Specifically,*

$$2166 \quad \mathbb{E}[\lambda(D)] = \mathbb{E}[\lambda(D_{\text{lp}})] + (1 - d)\sigma^2 t \quad (159)$$

2167 where t is the time-span of fine-tuning process.

2168 *Proof of Lemma E.19.* Pick any eigenpair (λ, u) of imbalance matrix D (Definition E.8) such that $\|u\|_2 = 1$. By Itô's lemma (Equation equation E.1):

$$2169 \quad \partial \lambda = u^T (\partial D) u + u^T (\partial D) (\lambda I - D)^\dagger (\partial D) u^T \quad (160)$$

$$2170 \quad = (1 - d)\sigma^2 \|u\|_2^2 \partial t + \partial M_t + (1 - d)^2 \sigma^4 u^T (\lambda I - D)^\dagger u^T \quad (161)$$

$$2171 \quad = (1 - d)\sigma^2 \partial t + \partial M_t + (1 - d)^2 \sigma^4 u^T (\lambda I - D)^\dagger u^T \quad (162)$$

2172 where M_t is the martingale induced by the Brownian noise and $(\cdot)^\dagger$ denotes the pseudo inverse of a certain matrix. Say the the singular value decomposition (SVD) of D is

$$2173 \quad D = U \Sigma U^T = U \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \lambda_2 & \\ \mathbf{0} & & \ddots \end{bmatrix} U^T \quad (163)$$

2174 where we have $\lambda \in \text{diag} \Sigma$ and u being a column vector in U . So we can write the SVD of $(\lambda I - D)$ as:

$$2175 \quad \lambda I - D = V \Sigma' V^T = V \begin{bmatrix} \lambda - \lambda_1 & & \mathbf{0} \\ & \lambda - \lambda_2 & \\ \mathbf{0} & & \ddots \end{bmatrix} V^T \quad (164)$$

2176 where we obtain V by removing u in the columns of U and we obtain Σ' by removing λ in Σ . Then the pseudo inverse of $(\lambda I - D)$ is

$$2177 \quad (\lambda I - D)^\dagger = V \Sigma' V^T = V \begin{bmatrix} \frac{1}{\lambda - \lambda_1} & & \mathbf{0} \\ & \frac{1}{\lambda - \lambda_2} & \\ \mathbf{0} & & \ddots \end{bmatrix} V^T \quad (165)$$

2178 Since U is orthogonal, we shall have $V^T u = \mathbf{0}$. Then we can rewrite the stochastic dynamics of D as:

$$2179 \quad \frac{\partial}{\partial t} \mathbb{E}[\lambda] = (1 - d)\sigma^2 \quad (166)$$

2180 □

2181 E.5 FINE-TUNING LOSS

2182 **Lemma E.20** (Bounding the norm of linear head $\|v\|_2^2$). *During fine-tuning (Equation equation 91), we can bound the norm of $\|v\|_2^2$ with the imbalance matrix D in Definition E.8 as*

$$2183 \quad \frac{\underline{\lambda} + \sqrt{\underline{\lambda}^2 + 4\|w\|^2}}{2} \leq \|v\|_2^2 \leq \frac{\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4\|w\|^2}}{2} \quad (167)$$

2184 where we denote $\underline{\lambda} = \lambda_{\min}(\hat{D})$, $\bar{\lambda} = \lambda_{\max}(\hat{D})$.

2214 *Proof of Lemma E.20.* Given the information of imbalance matrix, we can bound the linear head
 2215 norm. Denote $\underline{\lambda} = \lambda_{\min}(D)$, $\bar{\lambda} = \lambda_{\max}(D)$. Denote $w = B^T v$ and multiply D with v on both
 2216 sides:

$$2217 \quad v^T D v = (v^T v)^2 - (v^T B)(B^T v) \quad (168)$$

$$2218 \quad v^T D v = \|v\|_2^4 - \|w\|_2^2 \quad (169)$$

2219 We have a range for the Rayleigh quotient: $\frac{x^T D x}{x^T x} \in [\underline{\lambda}, \bar{\lambda}]$. So we obtain two inequalities:

$$2220 \quad \begin{cases} \|v\|_2^4 - \|w\|_2^2 \geq \underline{\lambda} \|v\|_2^2 \\ \|v\|_2^4 - \|w\|_2^2 \leq \bar{\lambda} \|v\|_2^2 \end{cases} \quad (170)$$

$$2221 \quad = \begin{cases} \|v\|_2^4 - \underline{\lambda} \|v\|_2^2 - \|w\|_2^2 \geq 0 \\ \|v\|_2^4 - \bar{\lambda} \|v\|_2^2 - \|w\|_2^2 \leq 0 \end{cases} \quad (171)$$

2222 To get a lower bound of v , we can solve two quadratic inequalities. For the first quadratic equation,
 2223 since the smaller root is non-positive, $\underline{\lambda} - \sqrt{\underline{\lambda}^2 + 4\|w\|_2^2} \leq 0$, we just bound $\|v\|_2^2$ with the larger
 2224 root:

$$2225 \quad \|v\|_2^2 \geq \frac{\underline{\lambda} + \sqrt{\underline{\lambda}^2 + 4\|w\|_2^2}}{2} \quad (172)$$

2226 similarly, for the second quadratic equation, we obtain an upper bound for $\|v\|_2^2$ with the right-side
 2227 zero point:

$$2228 \quad \|v\|_2^2 \leq \frac{\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4\|w\|_2^2}}{2} \quad (173)$$

2229 \square

2230 **Lemma E.21** (Bounding eigenvalues of $B^T B$ (re-stated from Min et al. (2023b))). *During fine-*
 2231 *tuning (Equation equation 91), we can bound any nonzero eigenvalue λ_i of $B^T B$ as*

$$2232 \quad \lambda_i \in \left[\frac{-\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4(z_i^T w)^2}}{2}, \frac{-\underline{\lambda} + \sqrt{\underline{\lambda}^2 + 4(z_i^T w)^2}}{2} \right] \quad (174)$$

2233 where we use the imbalance matrix D in Definition E.8 and denote

$$2234 \quad \begin{cases} \bar{\lambda} = \lambda_{\max}(D) \\ \underline{\lambda} = \lambda_{\min}(D) \end{cases} \quad (175)$$

2235 *Proof of Lemma E.21.* The proof of this lemma follows the proof of Lemma 3 in Min et al. (2023b).
 2236 $B^T B$ is symmetric and positive semidefinite ($x^T B^T B x = \|Bx\|_2^2 \geq 0$). So every eigenvalue of
 2237 $B^T B$ is non-negative.

2238 D has at most one positive eigenvalue: if D has more than one eigenvalues, then the subspace of
 2239 \mathbb{R}^h spanned by the all positive eigenvectors has dimension at least 2, which must have non-trivial
 2240 intersection with $\ker(v^T)$ as $\dim(\ker(v^T)) = h - 1$. Then there exists a nonzero vector $z \in \ker(v^T)$
 2241 such that $z^T D z > 0$, which would imply $-z^T B B^T z = z^T D z > 0$, a contradiction.

2242 For any eigenvalue-eigenvector pair (λ_i, z_i) of $B^T B$ where $\lambda_i \neq 0$ and $z_i \in \mathbb{S}^{d-1}$,

$$2243 \quad \lambda_i^2 = z_i^T (B^T B)^2 z_i \quad (176)$$

$$2244 \quad // \text{replace something with imbalance matrix} \quad (177)$$

$$2245 \quad \lambda_i^2 = (z_i^T w)^2 - z_i^T B^T D B z_i \quad (178)$$

$$2246 \quad \lambda_i^2 - (z_i^T w)^2 = -z_i^T B^T D B z_i \quad (179)$$

$$2247 \quad \lambda_i^2 - (z_i^T w)^2 \in (z_i^T (B^T B) z_i) \cdot [-\lambda_{\max}, -\lambda_{\min}] \quad (180)$$

$$2248 \quad \lambda_i^2 - (z_i^T w)^2 \in \lambda_i \cdot [-\lambda_{\max}, -\lambda_{\min}] \quad (181)$$

again, we can rewrite this as two quadratic inequalities

$$\begin{cases} \lambda_i^2 + \lambda_{\max} \lambda_i - (z_i^T w)^2 \geq 0 \\ \lambda_i^2 + \lambda_{\min} \lambda_i - (z_i^T w)^2 \leq 0 \end{cases} \quad (182)$$

from them we know that there are two possible intervals:

$$\begin{cases} \lambda_i \in \left[-\infty, \frac{-\lambda_{\max} - \sqrt{\lambda_{\max}^2 + 4(z_i^T w)^2}}{2} \right] \cup \left[\frac{-\lambda_{\max} + \sqrt{\lambda_{\max}^2 + 4(z_i^T w)^2}}{2}, +\infty \right] \\ \lambda_i \in \left[\frac{-\lambda_{\min} - \sqrt{\lambda_{\min}^2 + 4(z_i^T w)^2}}{2}, \frac{-\lambda_{\min} + \sqrt{\lambda_{\min}^2 + 4(z_i^T w)^2}}{2} \right] \end{cases} \quad (183)$$

Note that we must have $\lambda_i \geq 0$ since $B^T B$ is positive semidefinite. So we can rewrite the bounds:

$$\lambda_i \in \left[\frac{-\lambda_{\max} + \sqrt{\lambda_{\max}^2 + 4(z_i^T w)^2}}{2}, \frac{-\lambda_{\min} + \sqrt{\lambda_{\min}^2 + 4(z_i^T w)^2}}{2} \right] \quad (184)$$

since the function $f(x) = -x + \sqrt{x + c^2}$ is monotonically decreasing, we have $f(\lambda_{\max}) \leq f(\lambda_{\min})$, i.e. the lower bound is no greater than the upper bound, i.e. the above interval is always non-empty. \square

E.6 NUMERICAL CONJECTURE ON THE EIGENVALUES

Conjecture E.22 (Small relative error induced by Jensen gap (Equation 216)). We denote the minimum eigenvalue of the imbalance matrix D as $\underline{\lambda}$. The relative error $\frac{\mathbb{E}[\max(0, -\underline{\lambda})^{1/2}]^2 - \mathbb{E}[\underline{\lambda}]}{\mathbb{E}[\max(0, -\underline{\lambda})^{1/2}]^2}$ increases slowly in time and is smaller than 1% under reasonable number of training epochs. Here we provide an empirical example with huge noise scale (much greater than the common noise scale in real-world applications). We observe that the relative approximation error is insignificant even with huge noise scale.

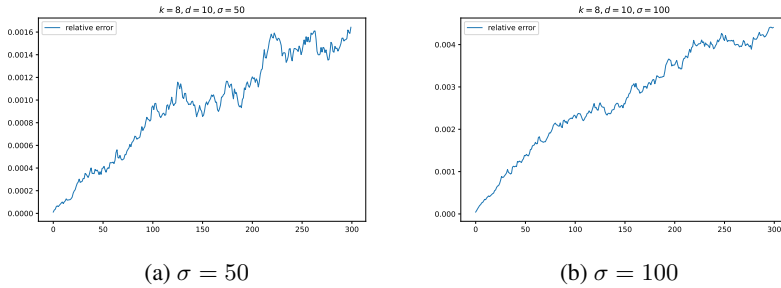


Figure 7: Growth of the relative error $\frac{\mathbb{E}[\max(0, -\underline{\lambda})^{1/2}]^2 - \mathbb{E}[\underline{\lambda}]}{\mathbb{E}[\max(0, -\underline{\lambda})^{1/2}]^2}$ in the experiment setting: (1) we use a two-layer linear network with a linear head of size $h = 8$ and a feature extractor of size $h \times d = 8 \times 10$; (2) we train the linear network with DP-SGD; (3) we repeat the experiment with large noise multipliers $\sigma = 50$ and $\sigma = 100$.

E.7 FINE-TUNING LOSS UPPER BOUND

Lemma E.23 (Imbalance matrix in fine-tuning under layerwise noise). *During fine-tuning (Equation (96)), the imbalance matrix D in Definition E.8 evolves as*

$$\mathbb{E} \left[\frac{dD}{dt} \right] = 0 \quad (185)$$

Proof of Lemma E.23. We prove this lemma by analyzing the infinitesimal generator A of imbalance matrix D :

$$A(D_0(v, B))_{ij} := \lim_{t \downarrow 0} \frac{\mathbb{E}^{D_0}[D_{ij}] - (D_0)_{ij}}{t} \quad (186)$$

$$= 0 + \sigma^2 \sum_{i' \in [h]} \sum_{j' \in [h]} \mathbf{1}[i' = j' = i = j] \quad (187)$$

$$- \sigma^2 \sum_{i' \in [h], j' \in [d]} \sum_{i'' \in [h], j'' \in [d]} \mathbf{1}[i' = i'' = i = j \text{ and } j' = j''] \quad (188)$$

the generator is zero for $i \neq j$. So we can just consider the case where $i = j$.

$$A(D_0(v, B))_{ii} = \sigma^2 \sum_{i' \in [h]} \sum_{j' \in [h]} \mathbf{1}[i' = j' = i] \quad (189)$$

$$- \sigma^2 \sum_{i' \in [h], j' \in [d]} \sum_{i'' \in [h], j'' \in [d]} \mathbf{1}[i' = i'' = i \text{ and } j' = j''] \quad (190)$$

$$= (d - d)\sigma^2 \quad (191)$$

$$= 0 \quad (192)$$

□

Theorem E.24 (Loss upper bound of fine-tuning). *In fine-tuning under layerwise noise (Equation equation 96), we have*

$$\mathbb{E}[\mathcal{L}] \lesssim \mathbb{E}[\mathcal{L}]e^{(-\bar{\lambda} + \sqrt{2}\sigma^2(1+d))t} + L^\square (1 - e^{(-\bar{\lambda} + \sqrt{2}\sigma^2(1+d))t}) \quad (193)$$

where $L^\square = \sigma^2 \frac{(1+d)\|X^T Y\| - d\bar{\lambda}}{\bar{\lambda} - \sqrt{2}\sigma^2(1+d)}$.

Proof of Theorem 4.5. We first simplify the loss dynamics:

$$\partial \mathcal{L} = \partial \frac{1}{2} \|XB^T v - Y\|^2 \quad (194)$$

$$= \frac{1}{2} \langle \nabla_v \|XB^T v - Y\|^2, \partial v \rangle + \frac{1}{2} \langle \nabla_B \|XB^T v - Y\|^2, \text{vec}(\partial B) \rangle \quad (195)$$

$$+ \frac{1}{4} (\partial v)^T H_{\|XB^T v - Y\|^2} (\partial v) + \frac{1}{4} [\text{vec}(\partial B)]^T H_{\|XB^T v - Y\|^2} \text{vec}(\partial B) \quad (196)$$

$$= (XB^T v - Y)^T XB^T \partial v + (XB^T v - Y)^T X (\partial B)^T v \quad (197)$$

$$+ \frac{1}{2} (\partial v)^T BB^T (\partial v) + \frac{1}{2} [\text{vec}(\partial B)]^T H_{\|XB^T v - Y\|^2} \text{vec}(\partial B) \quad (198)$$

$$= - (XB^T v - Y)^T XB^T BX^T (XB^T v - Y) \partial t + (XB^T v - Y)^T XB^T \sqrt{2\sigma^2 d} \partial W_t \quad (199)$$

$$- (XB^T v - Y)^T XX^T (XB^T v - Y) v^T v \partial t + (XB^T v - Y)^T X (\sqrt{2\sigma^2} \partial W_t') v \quad (200)$$

$$+ \sigma^2 \text{trace}(BB^T) \partial t + \sigma^2 d \|v\|^2 \partial t \quad (201)$$

$$= - (B^T v - X^T Y)^T B^T B (B^T v - X^T Y) \partial t + (B^T v - X^T Y)^T B^T \sqrt{2\sigma^2} \partial W_t \quad (202)$$

$$- (B^T v - X^T Y)^T (B^T v - X^T Y) v^T v \partial t + (B^T v - X^T Y)^T (\sqrt{2\sigma^2} \partial W_t') v \quad (203)$$

$$+ \sigma^2 \text{trace}(B^T B) \partial t + \sigma^2 d \|v\|^2 \partial t \quad (204)$$

By Lemma E.20 and Lemma E.21, we have

$$\partial \mathbb{E} \mathcal{L} = - \mathbb{E}[(w - X^T Y)^T (B^T B + v^T v I_{d \times d}) (w - X^T Y)] \partial t + \sigma^2 \mathbb{E}[\|B\|_F^2 + d \|v\|_2^2] \partial t \quad (205)$$

$$\leq \mathbb{E} \left\{ - \|w - X^T Y\|_2^2 \frac{\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4\|w\|^2}}{2} \partial t - \|w - X^T Y\|_2^2 \frac{-\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4(z_{\min}^T w)^2}}{2} \partial t \right\} \quad (206)$$

$$+ \mathbb{E} \left\{ \sigma^2 d \frac{-\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4(z_{\min}^T w)^2}}{2} \partial t + \sigma^2 d \frac{\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4\|w\|^2}}{2} \partial t \right\} \quad (207)$$

$$\leq - \frac{1}{2} \mathbb{E}[\|w - X^T Y\|_2^2 (\Lambda_{\min} + \Lambda_{\max})] \partial t + \frac{1}{2} \sigma^2 \mathbb{E}[d \Gamma_{\min} + \Gamma_{\max}] \partial t \quad (208)$$

where we define

$$\begin{cases} \Lambda_{\min} = \underline{\lambda} + \sqrt{\underline{\lambda}^2 + 4\|w\|^2} \geq \max(0, 2\underline{\lambda}) \\ \Lambda_{\max} = -\bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4(z_{\min}^T w)^2} \geq \max(0, -2\bar{\lambda}) \\ \Gamma_{\min} = -\underline{\lambda} + \sqrt{\underline{\lambda}^2 + 4(z_{\min}^T w)^2} \leq \max(2\|w\|, 2\|w\| - 2\underline{\lambda}) = 2\|w\| + 2\max(0, -\underline{\lambda}) \\ \Gamma_{\max} = \bar{\lambda} + \sqrt{\bar{\lambda}^2 + 4\|w\|^2} \leq \max(2\|w\|, 2\|w\| + 2\bar{\lambda}) = 2\|w\| + 2\max(0, \bar{\lambda}) \end{cases} \quad (209)$$

Denote the probability measure of the state at time t as ν_t . Then by using Jensen’s inequality, reverse Hölder’s inequality, etc., we can bound the first term:

$$\mathbb{E}[\|w - w_*\|_2^2(\Lambda_{\min} + \Lambda_{\max})] = \int \|w - w_*\|_2^2(\Lambda_{\min} + \Lambda_{\max})d\nu_t \quad (210)$$

$$\geq \left(\int \|w - w_*\|_2^{-1} d\nu_t \right)^{-2} \left(\int (\Lambda_{\min} + \Lambda_{\max})^{1/2} d\nu_t \right)^2 \quad (211)$$

$$= \mathbb{E}[\|w - w_*\|_2^{-1}]^{-2} \mathbb{E}[(\Lambda_{\min} + \Lambda_{\max})^{1/2}]^2 \quad (212)$$

$$\geq \mathbb{E}[\|w - w_*\|_2^2] \mathbb{E}[(\Lambda_{\min} + \Lambda_{\max})^{1/2}]^2 \quad (213)$$

$$\text{according our empirical observation (Conjecture E.22)} \quad (214)$$

$$\text{we ignore the Jensen gap for the second multiplier} \quad (215)$$

$$\gtrsim -\frac{1}{2} \mathbb{E}[\|w - w_*\|_2^2] \mathbb{E}[\bar{\lambda}] \quad (216)$$

$$\text{By Lemma E.19} \quad (217)$$

$$= \mathbb{E}[\|w - w_*\|_2^2] (-\mathbb{E}[\bar{\lambda}(D_0)] + (d-1)\sigma^2 t) \quad (218)$$

$$= 2(-\mathbb{E}[\bar{\lambda}(D_0)] + (d-1)\sigma^2 t) \cdot \mathbb{E}[\mathcal{L}] \quad (219)$$

Then we rewrite the upper bound:

$$\partial \mathbb{E}[\mathcal{L}] \leq -\frac{1}{2} \mathbb{E}[\|w - X^T Y\|_2^2(\Lambda_{\min} + \Lambda_{\max})] \partial t + \frac{1}{2} \sigma^2 \mathbb{E}[d\Gamma_{\min} + \Gamma_{\max}] \partial t \quad (220)$$

$$\partial \mathbb{E}[\mathcal{L}] \lesssim -\bar{\lambda} \mathbb{E}[\mathcal{L}] \partial t + \sigma^2 (\sqrt{2}(1+d) \mathbb{E}[\mathcal{L}]^{1/2} + (1+d) \|X^T Y\| - d\underline{\lambda}) \partial t \quad (221)$$

$$\partial \mathbb{E}[\mathcal{L}] \lesssim (-\bar{\lambda} + \sqrt{2}\sigma^2(1+d)) \mathbb{E}[\mathcal{L}] \partial t + \sigma^2 ((1+d) \|X^T Y\| - d\underline{\lambda}) \partial t \quad (222)$$

$$\mathbb{E}[\mathcal{L}] \lesssim \mathbb{E}[\mathcal{L}] e^{(-\bar{\lambda} + \sqrt{2}\sigma^2(1+d))t} + L^\square (1 - e^{(-\bar{\lambda} + \sqrt{2}\sigma^2(1+d))t}) \quad (223)$$

where $L^\square = \sigma^2 \frac{(1+d) \|X^T Y\| - d\underline{\lambda}}{\bar{\lambda} - \sqrt{2}\sigma^2(1+d)}$. □

F THEORY WITH CLIPPING

In this section, we present the **first** theoretical investigation on Langevin diffusion **with clipping**. We believe that our contribution is significant for the Langevin diffusion and private optimization research community. We summarize our findings and contributions in the following list:

- A new definition for Langevin diffusion with clipping (Definition F.1).
- Zeroth order approximation error for the clipped Langevin diffusion (Theorem F.3).
- Privacy guarantee for the clipped Langevin diffusion (Theorem F.4).
- The exact “discrete vs. continuous” algebraic correspondence between the clipped Langevin diffusion and vanilla DP-SGD (Remark F.2).
- Feature distortion analysis for the clipped Langevin diffusion (Theorem F.5).
- The existence proof of a unique strong solution for the clipped Langevin diffusion (Corollary F.7).

Definition F.1 (Clipped Langevin diffusion). Say we work on parameter $\theta \in \mathbb{R}^p$ to minimize a group of loss functions $\{\ell_i\}_{i \in [n]}$. The parameter evolve according to the following stochastic differential

2430 equation.

$$2431 \quad \partial\theta = - \sum_{i \in [n]} \text{clip}_C(\nabla\ell_i(\theta))\partial t + \sigma\partial\xi_t \quad (224)$$

2432 This equation is the clipped Langevin diffusion. ξ_t is a vector containing p independent 1-
2433 dimensional Brownian motion. The clipping function is defined by a constant $C > 0$ and

$$2434 \quad \text{clip}_C(\nabla\ell_i(\theta)) := \min\left(1, \frac{C}{\|\nabla\ell_i(\theta)\|_2}\right) \nabla\ell_i(\theta).$$

2435 This definition allows us to establish the first exact "discrete vs. continuous" algebraic correspon-
2436 dence between clipped Langevin diffusion and vanilla DP-SGD, creating a continuous analytical
2437 framework that closely mirrors real DP-SGD implementations.

2438 *Remark F.2* (Algebraic correspondence between the clipped Langevin diffusion and DP-SGD). The
2439 update rule of the vanilla DP-SGD with step-size $\eta > 0$ can be written as (Abadi et al., 2016):

$$2440 \quad \theta_{k+1} = \theta_k - \eta \frac{1}{|B|} \sum_{i \in B_k} (\text{clip}_C(\nabla\ell_i(\theta)) + \sigma\mathcal{N}(0, C^2\mathbf{I})) \quad (225)$$

2441 where B is the batch size and B_k is the batch of data points sampled at step k . We can rewrite the
2442 update rule by assuming full sampling, $\tilde{\eta} = \eta \frac{1}{|B|}$ and $\tilde{\sigma} = \sigma C$:

$$2443 \quad \theta_{k+1} = \theta_k - \tilde{\eta} \sum_{i \in [n]} (\text{clip}_C(\nabla\ell_i(\theta)) + \tilde{\sigma}\mathcal{N}(0, \mathbf{I})) \quad (226)$$

2444 One can compare this update rule with the clipped Langevin diffusion (Equation (224)):

$$2445 \quad \partial\theta = - \sum_{i \in [n]} \text{clip}_C(\nabla\ell_i(\theta))\partial t + \sigma\partial\xi_t \quad (227)$$

2446 It is easy to see the algebraic correspondence between the above two equations. We provide a
2447 rigorous derivation of DP-SGD update by discretizing the clipped Langevin diffusion with the Euler-
2448 Maruyama method.

2449 Suppose that we want to solve the clipped Langevin diffusion on some interval of time $[0, T]$. Then
2450 the Euler-Maruyama approximation to the true solution θ is the Markov chain $\tilde{\theta}$ defined as follows:

- 2451 • Partition the interval $[0, T]$ into K equal subintervals of width $\tilde{\eta} > 0$:

$$2452 \quad 0 = \tau_0 < \tau_1 < \dots < \tau_K = T \text{ and } \tilde{\eta} = \frac{T}{K} \quad (228)$$

- 2453 • Let $\tilde{\theta}_0 = \theta_0$ at the initialization.
- 2454 • Iteratively compute $\tilde{\theta}_k$ for $1 \leq k \leq K$ by

$$2455 \quad \tilde{\theta}_k = \tilde{\theta}_{k-1} - \tilde{\eta} \sum_{i \in [n]} (\text{clip}_C(\nabla\ell_i(\tilde{\theta}_{k-1})) + \tilde{\sigma}\mathcal{N}(0, \mathbf{I})) \quad (229)$$

2456 In this way, we rediscover the update rules for DP-SGD by discretizing the clipped Langevin diffu-
2457 sion.

2458 We give an approximation error bound following (Freidlin et al., 2012, Theorem 1.2, Chapter 2.1).

2459 **Theorem F.3** (Zeroth order approximation error). *For all $t > 0, \delta > 0$, we have*

$$2460 \quad \mathbb{E} \left[\left\| \theta_t - \theta_t^{(0)} \right\|^2 \right] \leq \left(\sigma(2p)^{\frac{1}{2}} t^{\frac{1}{2}} + 2nCt \right)^2 \quad (230)$$

2461 *Proof of Theorem F.3.*

$$2462 \quad \mathbb{E}[\partial\|\theta_t - \theta_t^{(0)}\|^2] = \mathbb{E}[\langle \theta_t - \theta_t^{(0)}, \partial\theta_t - \partial\theta_t^{(0)} \rangle + 2p\sigma^2\partial t] \quad (231)$$

$$\partial \mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2] \leq \mathbb{E}[4nC\|\theta_t - \theta_t^{(0)}\|\partial t + 2p\sigma^2\partial t] \quad (232)$$

$$\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2] \leq \int_0^T (4nC \cdot \mathbb{E}[\|\theta_t - \theta_t^{(0)}\|] + 2p\sigma^2)\partial t \quad (233)$$

$$\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2] \leq \int_0^T (4nC \cdot \sqrt{\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2]} + 2p\sigma^2)\partial t \quad (234)$$

$$\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2] \leq 2p\sigma^2T + 4nC \int_0^T \cdot \sqrt{\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2]}\partial t \quad (235)$$

By Lemma F.9, we have

$$\mathbb{E}[\|\theta_t - \theta_t^{(0)}\|^2] \leq \left(\sigma(2p)^{\frac{1}{2}}t^{\frac{1}{2}} + 2nCt\right)^2 \quad (236)$$

□

Note that this approximation error significantly improves upon the $O(\exp(T))$ error found under standard regularity assumptions (Freidlin et al., 2012, Theorem 1.2, Chapter 2.1).

We present a privacy guarantee for the clipped Langevin diffusion by deriving an upper bound on the KL divergence.

Theorem F.4 (KL Divergence Bound for Clipped Langevin Diffusion). *Let θ_0, θ'_0 have the same distribution Θ_0, Θ'_0 , θ_T be the solution to Equation (224) given initial condition θ_0 and database D , θ'_T be the solution to Equation (224) given initial condition θ'_0 and database D' , such that $D \sim D'$. Let $\Theta_{[0,T]}$ be the distribution of the trajectory $\theta_{t \in [0,T]}$. Then for any $T > 0$:*

$$\text{KL}(\Theta_{[0,T]} \|\Theta'_{[0,T]}) \leq \frac{2n^2C^2}{\sigma^2}T \quad (237)$$

Proof of Theorem F.4. By Theorem B.1 & 3.1 of Ye et al. (2023a),

$$\begin{aligned} \text{KL}(\Theta_{[0,T]} \|\Theta'_{[0,T]}) &= \frac{1}{2\sigma^2} \int_0^T \mathbb{E} \left[\left\| \sum_{i \in [n]} \text{clip}_C(\nabla \ell_i(\theta; D)) - \sum_{i \in [n]} \text{clip}_C(\nabla \ell_i(\theta; D')) \right\|_2^2 \right] dt \\ &\leq \frac{1}{2\sigma^2} \int_0^T 4n^2C^2 dt \\ &= \frac{2n^2C^2}{\sigma^2}T \end{aligned}$$

□

We demonstrate that our main result on feature distortion holds for clipped Langevin diffusion, reinforcing our paper’s key insight. Here, our approximation technique is essential, as the stochastic analysis of Langevin diffusion with nonlinear & nonconvex coefficients would be extremely challenging without it.

Theorem F.5 (Random initialization causes feature distortion). *If Assumption 3.1 and Assumption 3.2 hold, and the linear head is randomly initialized by $v_0 \sim \mathcal{N}(0, \beta I_{h \times h})$, then with probability $1 - 2^{-h}$, $\forall \beta > 0, \exists j \in [h], \Delta t > 0$ such that during the time interval $(0, \Delta t)$, DP-FFT distorts w_j reducing its alignment with the data cluster. The cosine similarity between w_j and the data cluster mean $\bar{x}_{c(j)}$ decreases monotonically:*

$$\left. \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) \right|_t < 0, \quad \forall t \in (0, \Delta t) \quad (238)$$

Proof of Theorem F.5. The per-sample gradient for the i -th data point (before clipping) is

$$\nabla_{(v,W)} \ell_i = \begin{bmatrix} \nabla_v \ell_i \\ \text{vec}(\nabla_W \ell_i) \end{bmatrix} = \begin{bmatrix} y_i \ell_i \text{relu}(W^\top x_i) \\ y_i \ell_i v_1 \text{relu}'(w_1^\top x_i) x_i \\ y_i \ell_i v_2 \text{relu}'(w_2^\top x_i) x_i \\ \vdots \\ y_i \ell_i v_h \text{relu}'(w_h^\top x_i) x_i \end{bmatrix} = y_i \ell_i \begin{bmatrix} \text{relu}(W^\top x_i) \\ v_1 \text{relu}'(w_1^\top x_i) x_i \\ v_2 \text{relu}'(w_2^\top x_i) x_i \\ \vdots \\ v_h \text{relu}'(w_h^\top x_i) x_i \end{bmatrix} \quad (239)$$

2537

where the $\text{vec}(\cdot)$ operator is defined as an operation that converts a tensor to a vector (Magnus & Neudecker, 1999, Chapter 2.4). We use $\text{vec}(\cdot)$ to collect the gradients of v and W into one vector. Then we can write the clipped per-sample gradient for the i -th data point as:

$$\text{clip}_C(\nabla_{(v,W)}\ell_i) = \min\left(1, \frac{C}{\|\nabla_{(v,W)}\ell_i\|_2}\right) \cdot y_i \ell_i \begin{bmatrix} \text{relu}(W^\top x_i) \\ v_1 \text{relu}'(w_1^\top x_i)x_i \\ v_2 \text{relu}'(w_2^\top x_i)x_i \\ \vdots \\ v_h \text{relu}'(w_h^\top x_i)x_i \end{bmatrix}. \quad (240)$$

Therefore, the dynamics of the parameter w_j for any $j \in [h]$ under gradient clipping is,

$$\frac{\partial w_j}{\partial t} = \min\left(1, \frac{C}{\|\nabla_{(v,W)}\ell_i\|_2}\right) \cdot y_i \ell_i \cdot v_j \text{relu}'(w_j^\top x_i)x_i \quad (241)$$

Note that the clipping operation only multiplies the gradient with a normalization term $\min\left(1, \frac{C}{\|\nabla_{(v,W)}\ell_i\|_2}\right)$. As a result, it does not change the signs of the gradient entries. Then we are ready to analyze the cosine similarity between w_j and the mean data direction:

$$\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) = \frac{2(w_j^\top \bar{x}_{c(j)})}{\|w_j\|_2^2} \left[\|w_j\|_2^2 \bar{x}_{c(j)}^\top \frac{\partial w_j}{\partial t} - \bar{x}_{c(j)}^\top w_j w_j^\top \frac{\partial w_j}{\partial t} \right] \quad (242)$$

$$= \frac{2(w_j^\top \bar{x}_{c(j)})}{\|w_j\|_2^2} \left[\|w_j\|_2^2 \bar{x}_{c(j)} - (\bar{x}_{c(j)}^\top w_j) w_j \right]^\top \frac{\partial w_j}{\partial t} \quad (243)$$

$$\text{//by Assumption 3.2} \quad (244)$$

$$\text{sign}\left(\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)})\right) = \text{sign}\left(\left[\|w_j\|_2^2 \bar{x}_{c(j)} - (\bar{x}_{c(j)}^\top w_j) w_j \right]^\top \frac{\partial w_j}{\partial t}\right) \quad (245)$$

$$\text{//the clipping operation preserves the sign} \quad (246)$$

$$= \text{sign}\left(v_j(\|w_j\|_2^2 - (\bar{x}_{c(j)}^\top w_j)^2)\right) \quad (247)$$

$$= \text{sign}(v_j) \quad (248)$$

Since we initialize $v \sim \mathcal{N}(0, \beta I_{h \times h})$, with probability $1 - 2^{-h}$, there exists j such that $v_j < 0$ at $t = 0 \implies \frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) < 0$ at $t = 0$. By the continuity of the approximated Langevin diffusion, there exists $\Delta t > 0$ such that for any $t \in (0, \Delta t)$,

$$\frac{\partial}{\partial t} \cos(w_j, \bar{x}_{c(j)}) < 0. \quad (249)$$

□

We establish that a unique and strong solution exists for the clipped Langevin diffusion. This result is particularly noteworthy because it bypasses the standard regularity assumptions typically required in existence proofs for stochastic differential equations (Mao, 1997; Øksendal, 2014). Standard conditions demand that both the drift and diffusion coefficients exhibit linear growth in their parameters and are Lipschitz continuous. However, such assumptions are often impractical for the loss functions prevalent in modern machine learning. Additionally, deep learning architectures frequently introduce non-differentiability (as seen in the discontinuities of ReLU activation functions, for instance). In response, we propose relaxed regularity criteria to address these challenges.

Theorem F.6 (Criteria of unique strong solution for SDE with irregular drift (Veretennikov, 1981, Theorem 1)). *Consider the following stochastic differential equation:*

$$dx_t = a(x_t, t)dt + b(x_t, t)dX_t \quad (250)$$

where

- X_t denotes the standard Wiener process.
- a is a bounded, d -dimensional vector-valued, measurable function.

- b is a bounded, matrix-valued, continuous measurable function of size $d \times d$. b satisfies the following properties:

- (Uniform elliptic condition): For any $x \in \mathbb{R}^d, v \in \mathbb{R}^d, t \geq 0$, there exists a constant $\lambda > 0$ such that

$$v^T b(x, t) b^T(x, t) v \geq \lambda v^T v \quad (251)$$

- (Fixed time uniform continuity): For every $T > 0$ and any $t \in [0, T]$, $b(\cdot, t)$ is uniformly continuous on any compact metric subspace $U \subset \mathbb{R}^d$.

Then a unique strong solution X_t exists for the stochastic differential equation.

Corollary F.7. If the per-sample loss function ℓ has a discontinuity set with Lebesgue measure 0, then the clipped Langevin diffusion (Equation (224)) has a unique strong solution.

Remark F.8 (Toy-case example of Corollary F.7). Consider a 2-layer ReLU network f parametrized by $v \in \mathbb{R}^h, W \in \mathbb{R}^{d \times h}$:

$$f(x) := v^T \text{relu}(W^T x), \quad (252)$$

a singleton training dataset $D := \{(x_0, y_0)\}$:

$$x_0 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad y_0 = 1 \quad (253)$$

and exponential loss $\ell(y, \hat{y}) := \exp(-y\hat{y})$. Then the drift coefficient (e.g. $a(x_t, t)$ in Theorem F.6) of the loss Langevin diffusion is

$$-\text{clip}_C(\nabla \ell_0(y_0, f(x_0))) = -\text{clip}_C(\nabla \ell_0(y_0, f(x_0))) \quad (254)$$

$$= -\min\left(1, \frac{C}{\|\nabla_{(v,W)} \ell_0\|_2}\right) \cdot y_i \ell_i \begin{bmatrix} \text{relu}(W^T x_i) \\ v_1 \text{relu}'(w_1^T x_i) x_i \\ v_2 \text{relu}'(w_2^T x_i) x_i \\ \vdots \\ v_h \text{relu}'(w_h^T x_i) x_i \end{bmatrix} \quad (255)$$

The set of all discontinuities of this drift coefficient has Lebesgue measure zero in the parameter space $\mathbb{R}^h \times \mathbb{R}^{d \times h}$. This drift coefficient is a measurable function. So we can apply Theorem F.6 in this example.

F.1 TECHNICAL RESULTS

Lemma F.9 (Gronwall type inequality IV). Let $x : [a, b] \rightarrow \mathbb{R}_+$ be a continuous function that satisfies the inequality:

$$x(t) \leq M + \int_a^t \Psi(s) \omega(x(s)) ds, \quad t \in [a, b]$$

where $M \geq 0, \Psi : [a, b] \rightarrow \mathbb{R}_+$ is continuous and $\omega : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is continuous and monotone-increasing. Then the estimation

$$x(t) \leq \Phi^{-1}\left(\Phi(M) + \int_a^t \Psi(s) ds\right), \quad t \in [a, b]$$

holds, where $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ is give by

$$\Phi(u) := \int_{u_0}^u \frac{1}{\omega(s)} ds, \quad u \in \mathbb{R}$$

Proof of Lemma F.9. This proof is done by Sever Silvestru Dragomir.

We just copy the proof here for completeness.

2646 Denote $y(t)$ as

$$2647 \quad y(t) := \int_a^t \omega(x(s))\Psi(s)ds, \quad t \in [a, b]$$

2648 we have $y(a) = 0$, and by the recursive integral condition of x , we obtain:

$$2651 \quad y'(t) = x(t)\Psi(t), \quad t \in [a, b]$$

$$2652 \quad y'(t) \leq \omega(M + y(t))\Psi(t)$$

$$2653 \quad \frac{1}{\omega(M + y(t))}d(y(t)) \leq \Psi(t)dt$$

2656 By integration on $[a, t]$, we have

$$2658 \quad \left(\int_0^{y(t)} \frac{1}{\omega(M + s)} ds \right) - \Phi(M) \leq \int_a^t \Psi(s)ds$$

$$2660 \quad \int_0^{y(t)} \frac{1}{\omega(M + s)} ds \leq \int_a^t \Psi(s)ds + \Phi(M)$$

2663 that is,

$$2664 \quad \Phi(y(t) + M) \leq \int_a^t \Psi(s)ds + \Phi(M)$$

2666 By taking the inverse mapping of Φ on both sides, we finish the proof. \square

2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699