

# Bernoulli-LoRA: A Theoretical Framework for Randomized Low-Rank Adaptation

Igor Sokolov<sup>1</sup>   Abdurakhmon Sadiev<sup>1</sup>   Yury Demidovich<sup>1</sup>  
Fawaz S Al-Qahtani<sup>2</sup>   Peter Richtárik<sup>1</sup>

<sup>1</sup> Center of Excellence for Generative AI,  
King Abdullah University of Science and Technology (KAUST), Saudi Arabia  
<sup>2</sup> Saudi Data & AI Authority (SDAIA) & National Center of AI (NCAI), Saudi Arabia

## Abstract

Parameter-efficient fine-tuning (PEFT) has emerged as a crucial approach for adapting large foundational models to specific tasks, particularly as model sizes continue to grow exponentially. Among PEFT methods, [Low-Rank Adaptation \(LoRA\)](#) [Hu et al., 2022] stands out for its effectiveness and simplicity, expressing adaptations as a product of two low-rank matrices. While extensive empirical studies demonstrate LoRA’s practical utility, theoretical understanding of such methods remains limited. Recent work on [RAC-LoRA](#) [Malinovsky et al., 2024] took initial steps toward rigorous analysis. In this work, we introduce [Bernoulli-LoRA](#), a novel theoretical framework that unifies and extends existing LoRA approaches. Our method introduces a probabilistic Bernoulli mechanism for selecting which matrix to update. This approach encompasses and generalizes various existing update strategies while maintaining theoretical tractability. Under standard assumptions from non-convex optimization literature, we analyze several variants of our framework: [Bernoulli-LoRA-GD](#), [Bernoulli-LoRA-SGD](#), [Bernoulli-LoRA-PAGE](#), and [Bernoulli-LoRA-MVR](#), [Bernoulli-LoRA-QGD](#), [Bernoulli-LoRA-MARINA](#), [Bernoulli-LoRA-EF21](#), establishing convergence guarantees for each variant. Additionally, we extend our analysis to convex non-smooth functions, providing convergence rates for both constant and adaptive (Polyak-type) stepsizes. Through extensive experiments on various tasks, we validate our theoretical findings and demonstrate the practical efficacy of our approach. This work is a step toward developing theoretically grounded yet practically effective PEFT methods.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Problem Statement</b>	<b>3</b>
<b>3</b>	<b>Motivation</b>	<b>4</b>
<b>4</b>	<b>Contributions</b>	<b>5</b>
<b>5</b>	<b>Notation</b>	<b>7</b>
<b>6</b>	<b>Bernoulli-LoRA Framework</b>	<b>8</b>
6.1	Reformulation as a Projected Gradient Step . . . . .	8
6.2	Core Algorithmic Variants . . . . .	9
6.3	Extensions for Federated Learning . . . . .	10

<b>7</b>	<b>Convergence Results</b>	<b>11</b>
<b>8</b>	<b>Experiments</b>	<b>15</b>
8.1	Linear Regression with Non-convex Regularization. . . . .	15
8.2	MLP on MNIST . . . . .	16
<b>A</b>	<b>Basic Facts and Useful Inequalities</b>	<b>24</b>
<b>B</b>	<b>Discussion on Positive Expected Projection (Assumption 1)</b>	<b>25</b>
<b>C</b>	<b>Proofs for Core Algorithmic Variants</b>	<b>27</b>
C.1	Analysis of Bernoulli-LoRA-GD . . . . .	27
C.1.1	Convergence for Smooth Non-Convex Functions . . . . .	29
C.1.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	31
C.1.3	Convergence for Non-Smooth Convex Functions . . . . .	32
C.2	Analysis of Bernoulli-LoRA-SGD . . . . .	39
C.2.1	Convergence for Smooth Non-Convex Functions . . . . .	39
C.2.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	42
C.3	Analysis of Bernoulli-LoRA-MVR . . . . .	43
C.3.1	Convergence for Smooth Non-Convex Functions . . . . .	45
C.3.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	46
C.4	Analysis of Bernoulli-LoRA-PAGE . . . . .	48
C.4.1	Convergence for Smooth Non-Convex Functions . . . . .	49
C.4.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	50
<b>D</b>	<b>Proofs for Federated Learning Extensions</b>	<b>51</b>
D.1	Analysis of Fed-Bernoulli-LoRA-QGD . . . . .	51
D.1.1	Convergence for Smooth Non-Convex Functions . . . . .	53
D.1.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	54
D.2	Analysis of Fed-Bernoulli-LoRA-MARINA . . . . .	55
D.2.1	Convergence for Smooth Non-Convex Functions . . . . .	56
D.2.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	57
D.3	Analysis of Fed-Bernoulli-LoRA-EF21 . . . . .	59
D.3.1	Convergence for Smooth Non-Convex Functions . . . . .	60
D.3.2	Convergence under Polyak-Łojasiewicz Condition . . . . .	61
<b>E</b>	<b>Experiments: Missing Details</b>	<b>63</b>
E.1	Linear Regression with Non-convex Regularization . . . . .	63

# 1 Introduction

Fine-tuning is a transfer learning method, where a pre-trained neural network is trained on a new dataset. In modern deep learning, adapting large models to specific tasks via fine-tuning has become central, especially in natural language processing [Peters et al., 2018, Devlin et al., 2019]. While full fine-tuning often yields strong results, it is computationally intensive for large models. Parameter-Efficient Fine-Tuning (PEFT) [He et al., 2021] addresses this by updating only a small subset of parameters [Richtárik and Takáč, 2016, Demidovich et al., 2023a], often with task-specific layers trained from scratch. PEFT offers performance close to full fine-tuning with reduced training time and resource use [Radford et al., 2019, Brown et al., 2020, Han et al., 2024], making it widely adopted in practice. Research on PEFT, especially for large foundation and language models, is rapidly growing.

Pre-trained models are known to have an inherently low intrinsic dimensionality [Li et al., 2018, Aghajanyan et al., 2020]. This means that fine-tuning can be effectively achieved within a reduced-dimensional subspace rather than the full parameter space. Among the various methods for utilizing this property, **Low-Rank Adaptation (LoRA)** [Hu et al., 2022] stands out as the most prominent reparameterization technique. **LoRA** minimizes the need to update an entire large, dense weight matrix by leveraging the product of two trainable low-rank matrices. This method significantly reduces the number of parameters required for fine-tuning. The low-rank matrices are optimized so that their scaled product serves as the update applied to the weight matrix:

$$W = W^0 + \frac{\alpha}{r}BA,$$

where  $W^0 \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{m \times r}$ , and  $A \in \mathbb{R}^{r \times n}$ . The pre-trained weight matrix  $W^0$  remains fixed, while  $A$  and  $B$  are the trainable matrices. Typically,  $A$  is initialized randomly using a Gaussian distribution, while  $B$  is set to zero to ensure that  $\Delta W = 0$  at the start. Various alternative initialization strategies have been investigated by [Zhu et al., 2024, Hayou et al., 2024, Meng et al., 2024, Wang et al., 2024]. The parameters of **LoRA** include the low rank  $r$  and the scaling factor  $\alpha$ . Since the dimensions  $m$  and  $n$  in deep learning models are usually large, selecting  $r \ll \min\{m, n\}$  drastically reduces the number of trainable parameters. The scaling factor  $\alpha$  acts as the stepsize. While **LoRA** may not always achieve the performance of full fine-tuning, it is more effective at mitigating forgetting when compared to traditional regularization techniques like weight decay and dropout. Additionally, it enhances diversity in generated outputs [Biderman et al., 2024]. Furthermore, **LoRA** is straightforward to implement and achieves performance comparable to full fine-tuning across a wide range of downstream tasks [Hu et al., 2022]. Complementing its algorithmic utility, research has also focused on enhancing the computational efficiency of **LoRA**; for instance, Cherniuk et al. [2023] demonstrated that by optimizing the computation graph of **LoRA** operations based on layer dimensions and rank, significant speedups and memory savings can be achieved without sacrificing accuracy. For a detailed summary of recent advancements in **LoRA**, refer to [Mao et al., 2025].

To bridge the gap between full fine-tuning and **LoRA**, Xia et al. [2024] introduced **Chain of LoRA (COLA)**, an iterative optimization framework that enhances model weights via higher-rank representations composed of multiple low-rank components—without added computational or memory cost. **COLA** incrementally refines low-rank approximations by training a sequence of **LoRA** modules through structured fine-tuning, merging, and extension. Each iteration adds a new low-rank component, forming a chain whose length reflects the number of optimized modules. The core idea involves applying **LoRA** updates iteratively over  $T$  steps: training a module, integrating its updates into fixed parameters, re-initializing, and repeating. This cyclic process builds higher-rank augmentations efficiently. In essence, **COLA** applies successive **LoRA** updates as:

$$W = W^0 + \frac{\alpha}{r} \sum_{t=0}^{T-1} B^t A^t.$$

Each pair  $(A^t, B^t)$  is initialized like standard **LoRA**. Unlike traditional **LoRA**, which may struggle with non-low-rank adaptations, **COLA** uses sequential low-rank decompositions to approximate updates of intermediate-to-high rank. This leads to more accurate and efficient adaptation while simplifying optimization by avoiding a direct high-rank fit.

## 2 Problem Statement

Supervised machine learning typically frames the training process as an optimization problem, aiming to minimize a loss function that quantifies the discrepancy between model predictions and true targets. This research focuses on the intricacies of this optimization challenge within the fine-tuning paradigm, where a pre-trained model is adapted to a new, specific task or dataset. Effective fine-tuning hinges on

making precise and efficient modifications to the model’s parameters to enhance its performance on the target task. We explore a generalized formulation of this problem, which is independent of the specific architecture of the underlying model:

$$\min_{\Delta W \in \mathbb{R}^{m \times n}} f(W^0 + \Delta W). \quad (1)$$

Here,  $W^0 \in \mathbb{R}^{m \times n}$  represents the parameters of the pre-trained model (or, for instance, the parameters of a single linear layer if other layers are kept constant), and  $\Delta W \in \mathbb{R}^{m \times n}$  is the adaptation term whose optimal value we seek. The function  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  denotes the empirical loss computed over the specific target dataset. Given that the dimensionality  $m \times n$  is typically very large in contemporary deep learning models, the adjustment  $\Delta W$  must possess a sufficiently simple structure to be practically trainable and applicable.

For the stochastic optimization methods developed and analyzed in this paper, we consider objective functions with one of the following specific structures:

- **Finite-Sum Setting:** The objective is an average of individual loss functions, a structure we address with methods like [Bernoulli-LoRA-PAGE](#):

$$f(W^0 + \Delta W) = \frac{1}{N} \sum_{i=1}^N f_i(W^0 + \Delta W), \quad (2)$$

where each  $f_i$  corresponds to the loss on a single data sample, and  $N$  is the total number of data points in the training set.

- **Expectation Setting:** The objective is an expectation over a data distribution  $\mathcal{D}$ , relevant for methods such as [Bernoulli-LoRA-MVR](#):

$$f(W^0 + \Delta W) = \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(W^0 + \Delta W)], \quad (3)$$

where  $f_\xi$  is the loss function associated with a data sample  $\xi$  drawn from  $\mathcal{D}$ .

Moreover, this paper extends its investigation to the **distributed optimization setting**, which is central to the Federated Learning (FL) algorithms we propose (e.g., [Fed-Bernoulli-LoRA-QGD](#), [Fed-Bernoulli-LoRA-MARINA](#), and [Fed-Bernoulli-LoRA-EF21](#)). In this context, we address problems formulated as:

$$f(W^0 + \Delta W) = \frac{1}{M} \sum_{l=1}^M f_l(W^0 + \Delta W), \quad (4)$$

where  $M$  is the total number of participating clients, and  $f_l$  represents the local loss function for client  $l$ , defined over its private dataset. The goal is to find a common adaptation  $\Delta W$  that minimizes this global, federated objective.

### 3 Motivation

Despite the widespread adoption and empirical success of [Low-Rank Adaptation \(LoRA\)](#) and its variants like [Chain of LoRA \(COLA\)](#), a comprehensive theoretical understanding underpinning these prevalent fine-tuning methods remains largely undeveloped. Several critical issues highlight this gap. Firstly, as pointed out by [Sun et al. \[2024\]](#), the [LoRA](#) re-parameterization inherently transforms a smooth Lipschitz loss into a non-smooth one. This alteration introduces significant theoretical complexities beyond those associated with managing the low-rank structure of updates, forming a key barrier to establishing robust theoretical frameworks. Secondly, the existing theoretical analysis of [COLA](#) by [Xia et al. \[2024\]](#) sidesteps

the core mechanism of low-rank updates by focusing on full-rank matrix optimization ( $\Delta W$ ). Such an approach is unsatisfactory as it fails to model or explain the very essence of LoRA’s efficiency.

Consequently, most methods based on LoRA are, in essence, heuristics, developed through empirical investigation without strong theoretical convergence guarantees. This is problematic, as these methods can be highly sensitive to hyperparameter choices [Khodak et al., 2021, Kuang et al., 2024], and their reliability beyond current empirical validation is not assured. In fact, Malinovsky et al. [2024] provided a concrete example of COLA’s potential divergence, further underscoring its heuristic nature. Their work introduced RAC-LoRA, the first comprehensive optimization framework designed to rigorously evaluate and establish convergence rates for methods utilizing LoRA-style updates, marking a significant step towards theoretically grounded PEFT.

However, while RAC-LoRA provides a foundational theoretical lens, its scope does not encompass several critical aspects of modern optimization, particularly for non-convex problems and distributed settings. Specifically, the RAC-LoRA framework does not utilize optimal variance-reduced techniques for non-convex optimization, nor does it delve into sophisticated Federated Learning (FL) setting that incorporate crucial practical techniques such as communication compression [Alistarh et al., 2018, Wen et al., 2017, Horvóth et al., 2022, Panferov et al., 2024] and error feedback. Federated learning [Konečný et al., 2016, Konečný et al., 2016, McMahan et al., 2016, Kairouz et al., 2019] is a decentralized paradigm where multiple clients collaboratively train a model on their local, private data. The growing demand for training massive deep neural networks with billions of parameters on vast datasets [Brown et al., 2020, Kolesnikov et al., 2020] has intensified the ML community’s interest in distributed optimization. To achieve feasible training times [Li, 2020], distributing computation, especially stochastic gradient evaluations, is essential, driving the adoption of scalable algorithms [Goyal et al., 2017, You et al., 2019, Le Scao et al., 2023]. Our work is motivated by the need to bridge this gap by extending a theoretically sound LoRA framework to these advanced and practically vital optimization scenarios.

## 4 Contributions

The performance of LoRA-based methods is notably sensitive to the selection of hyperparameters [Khodak et al., 2021, Kuang et al., 2024], and a robust theoretical understanding to guide their application is still developing. While recent work by Malinovsky et al. [2024] on RAC-LoRA provided initial steps towards a rigorous analytical framework, we aim to further advance the theoretical foundations and practical versatility of low-rank adaptation techniques.

In PEFT approaches based on low-rank adaptation, two matrices,  $A$  and  $B$ , are typically updated. Existing methods may update only  $A$ , only  $B$ , or alternate between them deterministically [Malinovsky et al., 2024, Xia et al., 2024, Zhu et al., 2024]. Our primary contribution is the introduction of Bernoulli-LoRA, a novel and generic low-rank adaptation framework. Bernoulli-LoRA is characterized by its unique probabilistic update mechanism: at each step of the adaptation process, a Bernoulli trial (akin to a coin flip) determines which of the two matrices ( $A$  or  $B$ ) is selected for optimization, while the other matrix is sampled from a predefined distribution and remains fixed for that step. This randomized selection not only provides a flexible approach but also unifies and generalizes several existing update strategies within a single theoretical construct. Much like the iterative design of COLA [Xia et al., 2024], the Bernoulli-LoRA framework operates by applying a sequence of such probabilistically chosen low-rank updates.

Our theoretical analysis is grounded in standard assumptions common in non-convex optimization literature, such as  $L$ -smoothness of the objective function. We instantiate the Bernoulli-LoRA framework by developing and analyzing several distinct algorithmic variants. These variants span a range of optimization techniques, from foundational gradient-based methods to more advanced stochastic, variance-reduced, and federated learning algorithms, each designed to address specific challenges in modern machine learning. For every proposed method within the Bernoulli-LoRA framework, we establish rigorous convergence guarantees. Our key contributions, which advance the theoretical understanding

and practical applicability of LoRA-type methods, include:

- ◆ **Foundational Algorithmic Variants:** We begin by establishing the theoretical properties of **Bernoulli-LoRA** with two fundamental optimization schemes. These methods lay the groundwork for understanding how the randomized selection of  $A$  or  $B$  interacts with standard descent procedures in the context of low-rank updates.
  - **Bernoulli-LoRA-GD** (Algorithm 2) serves as the simplest instantiation, employing full gradient descent to update the trainable low-rank matrix. While computing the full gradient is often impractical for large-scale models, this variant provides crucial foundational understanding of the framework’s convergence behavior under idealized conditions, navigating the optimization landscape defined by the LoRA reparameterization.
  - **Bernoulli-LoRA-SGD** (Algorithm 4) offers a more practical and widely applicable alternative by utilizing stochastic gradients. This variant addresses the computational burden of full gradient methods and is a cornerstone for larger-scale learning tasks, providing insights into the interplay of stochasticity and randomized matrix adaptation.
- ◆ **Advanced Variance Reduction Techniques for Non-Convex Optimization:** Stochastic gradients, while efficient, introduce variance that can impede convergence. Integrating variance reduction (VR) into the **LoRA** structure, particularly with the additional Bernoulli randomization, presents unique analytical challenges. Our work addresses this by developing specific VR-enhanced variants for **Bernoulli-LoRA**. To the best of our knowledge, we provide the first theoretical analyses demonstrating provable benefits for LoRA-type methods incorporating advanced VR schemes in  $L$ -smooth non-convex settings. Specifically, we propose:
  - **Bernoulli-LoRA-PAGE** (Algorithm 6): Tailored for the finite-sum setting (2), this method integrates the **Probabilistic Gradient Estimator (PAGE)** [Li et al., 2021]. **PAGE** is recognized for achieving optimal non-convex convergence rates and implementation simplicity, and we successfully adapt it to the **Bernoulli-LoRA** context.
  - **Bernoulli-LoRA-MVR** (Algorithm 5): For the infinite-sum (expectation) setting, this variant employs **Momentum Variance Reduction** techniques inspired by **STORM** [Cutkosky and Orabona, 2019]. **MVR** offers an efficient batch-free approach to VR, and our work demonstrates its compatibility and effectiveness within the **Bernoulli-LoRA** paradigm.
- ◆ **Communication-Efficient Federated Learning Extensions:** The application of PEFT methods like **LoRA** in Federated Learning (FL) is promising but requires careful consideration of communication overhead and data heterogeneity. We extend **Bernoulli-LoRA** to FL by designing three specialized algorithms that combine our randomized adaptation with established FL communication-saving techniques. To the best of our knowledge, this constitutes the first comprehensive theoretical analysis of LoRA-type methods integrated with established communication-efficient FL techniques such as quantization, gradient difference compression, and error feedback. Our FL extensions include:
  - **Fed-Bernoulli-LoRA-QGD** (Algorithm 7): This method tackles high communication bandwidth by incorporating **QSGD**-style quantization [Alistarh et al., 2017, Wen et al., 2017, Horvóth et al., 2022, Panferov et al., 2024], enabling clients to transmit compressed gradient information, a crucial feature for practical FL deployments.
  - **Fed-Bernoulli-LoRA-MARINA** (Algorithm 8): We adapt the **MARINA** communication compression strategy [Gorbunov et al., 2021], which efficiently compresses gradient differences, to the **Bernoulli-LoRA** framework. This is particularly beneficial for non-convex distributed learning over potentially heterogeneous datasets.



Setting	Method	NC convergence rate	PL convergence rate
(1)	Bernoulli-LoRA-GD (Alg. 2)	$\frac{\Delta^0}{\gamma\lambda_{\min}T}$	$(1 - \gamma\mu\lambda_{\min})^T \Delta^0$
(1)	Bernoulli-LoRA-SGD (Alg. 4)	$\frac{\Delta^0}{\gamma\lambda_{\min}T} + \frac{\gamma LC_1 \lambda_{\max}}{\lambda_{\min}}$	$(1 - \gamma\mu\lambda_{\min})^T \Delta^0 + \frac{\gamma LC_1 \lambda_{\max}}{\mu\lambda_{\min}}$
(1)+(3)	Bernoulli-LoRA-MVR (Alg. 5)	$\frac{\Phi_1}{\gamma\lambda_{\min}T} + \frac{b\sigma^2 \lambda_{\max}}{(2-b)\lambda_{\min}} \text{ (1)}$	$(1 - \gamma\mu\lambda_{\min})^T \Phi_1 + \frac{b\sigma^2 \lambda_{\max}}{(2-b)\mu\lambda_{\min}} \text{ (1)}$
(1)+(2)	Bernoulli-LoRA-PAGE (Alg. 6)	$\frac{\Phi_2}{\gamma\lambda_{\min}T} \text{ (2)}$	$(1 - \gamma\mu\lambda_{\min})^T \Phi_2 \text{ (2)}$
(1)+(4)	Fed-Bernoulli-LoRA-QGD (Alg. 7)	$\frac{\Delta^0}{\gamma\lambda_{\min}T} + \frac{\gamma L\omega\Delta^* \lambda_{\max}}{M\lambda_{\min}}$	$(1 - \gamma\mu\lambda_{\min})^T \Delta^0 + \frac{\gamma L^2\omega \lambda_{\max}}{M\mu\lambda_{\min}}$
(1)+(4)	Fed-Bernoulli-LoRA-MARINA (Alg. 8)	$\frac{\Phi_2}{\gamma\lambda_{\min}T} \text{ (2)}$	$(1 - \gamma\mu\lambda_{\min})^T \Phi_2 \text{ (2)}$
(1)+(4)	Fed-Bernoulli-LoRA-EF21 (Alg. 9)	$\frac{\Phi_3}{\gamma\lambda_{\min}T} \text{ (3)}$	$(1 - \gamma\mu\lambda_{\min})^T \Phi_3 \text{ (3)}$

(1)  $\Phi_1 := \Delta^0 + \frac{\gamma}{b(2-b)}\mathcal{G}^0$ ;

(2)  $\Phi_2 := \Delta^0 + \frac{\gamma}{q}\mathcal{G}^0$ ;

(3)  $\Phi_3 := \Delta^0 + \frac{\gamma}{1-\sqrt{1-\beta}}\hat{\mathcal{G}}^0$ .

Table 1: Summary of the convergence rates for the proposed methods, presented for smooth non-convex functions (“NC”) and for functions satisfying the PL-condition (“PL”). Absolute constant factors are omitted. Notation:  $\Delta^0 := f(W^0) - f^*$ ;  $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_F^2$ ;  $\hat{\mathcal{G}}^0 := \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_F^2$ ;  $T$  is the chain length;  $\omega$  is the compression parameter;  $\Delta^* := f^* - \frac{1}{M} \sum_{l=1}^M f_l^*$ ;  $C_1$  is a constant from Asm. 4;  $q$  is the probability of a full gradient computation;  $\beta$  is the contractive compression parameter;  $b$  is the momentum parameter;  $\lambda_{\min} = \lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ , and  $\lambda_{\max} = \lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ .

- **Fed-Bernoulli-LoRA-EF21** (Algorithm 9): This algorithm integrates the modern EF21 error feedback mechanism [Richtárik et al., 2021]. Error feedback is vital for stabilizing training with contractive compressors, and we show how Bernoulli-LoRA can leverage this for robust distributed fine-tuning.

◆ **Analysis for Non-Smooth Convex Functions:** Recognizing that not all machine learning objectives are smooth, we broaden the applicability of our framework. To the best of our knowledge, we present the first theoretical analysis of LoRA-type methods specifically for the important class of non-smooth convex optimization problems. For this setting, we provide versions of Bernoulli-LoRA-GD (Algorithm 3) and establish their convergence rates for both constant stepsize policies and adaptive Polyak-type stepsizes, showcasing the versatility of the Bernoulli-LoRA approach beyond smooth, non-convex settings.

## 5 Notation

For matrices  $W \in \mathbb{R}^{m \times n}$ , where  $m$  and  $n$  denote the input and output dimensions respectively, we employ the Frobenius norm  $\|\cdot\|_F$ , defined as  $\|W\|_F = \sqrt{\text{Tr}(W^\top W)}$ , where  $\text{Tr}(\cdot)$  denotes the matrix trace. The inner product between two matrices  $A$  and  $B$  is denoted by  $\langle A, B \rangle = \text{Tr}(A^\top B)$ . In our low-rank adaptation framework,  $B \in \mathbb{R}^{m \times r}$  and  $A \in \mathbb{R}^{r \times n}$  represent the factors of rank  $r \ll \min\{m, n\}$ . We use  $\mathcal{O}(\cdot)$  to hide absolute constants. We denote  $\Delta^0 := f(W^0) - f^*$ ,  $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_F^2$  and  $\hat{\mathcal{G}}^0 := \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_F^2$ . For differentiable functions  $f$ , the gradient  $\nabla f(W) \in \mathbb{R}^{m \times n}$  is computed with respect to the trace inner product, while for non-smooth functions, the subgradient  $\partial f(W) \in \mathbb{R}^{m \times n}$  is similarly defined. The superscript  $\dagger$  denotes the Moore-Penrose pseudoinverse.

## 6 Bernoulli-LoRA Framework

In this section, we introduce the **Bernoulli-LoRA** framework, a novel and generic approach for low-rank adaptation. The core idea is to perform a sequence of low-rank updates, where at each step, a probabilistic choice determines which of the two factor matrices ( $A$  or  $B$ ) is trained. This randomized mechanism, formalized in Algorithm 1, not only provides a flexible and unifying theoretical construct for existing LoRA-style methods but also allows for a rigorous convergence analysis.

At each iteration, one of the two low-rank matrices is sampled from a fixed distribution and remains frozen, while the other is trained to minimize the objective. This strategy prevents optimization from being confined to a fixed subspace, reducing the risk of converging to a suboptimal point. We formalize these two configurations as Left and Right sketch updates.

**Definition 1** (Left Sketch). *The left sketch update rule is given by*

$$\Delta W = \frac{\alpha}{r} B_S \hat{A}, \quad (5)$$

where  $B_S \sim \mathcal{D}_B$  is sampled from a fixed distribution over  $\mathbb{R}^{m \times r}$  matrices, and only the matrix  $\hat{A} \in \mathbb{R}^{r \times n}$  is adjustable.

**Definition 2** (Right Sketch). *The right sketch update rule is given by*

$$\Delta W = \frac{\alpha}{r} \hat{B} A_S, \quad (6)$$

where  $A_S \sim \mathcal{D}_A$  is sampled from a fixed distribution over  $\mathbb{R}^{r \times n}$  matrices, and only the matrix  $\hat{B} \in \mathbb{R}^{m \times r}$  is adjustable.

---

### Algorithm 1 Bernoulli-LoRA Framework

---

- 1: **Parameters:** pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling factor  $\alpha > 0$ , chain length  $T$ , sketch distributions  $\mathcal{D}_S^B$  and  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ .
  - 2: **for**  $t = 0, 1, \dots, T - 1$  **do**
  - 3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
  - 4:   **if**  $c^t = 1$  **then**
  - 5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  (Left sketch)
  - 6:     Using a chosen optimizer, approximately solve  $\hat{A}^t \approx \arg \min_A f(W^t + \frac{\alpha}{r} B_S^t A)$ .
  - 7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ .
  - 8:   **else**
  - 9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  (Right sketch)
  - 10:    Using a chosen optimizer, approximately solve  $\hat{B}^t \approx \arg \min_B f(W^t + \frac{\alpha}{r} B A_S^t)$ .
  - 11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ .
  - 12:   **end if**
  - 13: **end for**
- 

### 6.1 Reformulation as a Projected Gradient Step

Building upon the work of Malinovsky et al. [2024] on their **RAC-LoRA** framework, the update steps in Algorithm 1 can be reformulated as projected gradient steps. The subproblems in lines 6 and 10 are typically solved approximately, for instance, by taking a single step of a suitable optimizer like Gradient Descent (GD) or its variants.



Following the approach of [Malinovsky et al. \[2024\]](#), let's consider the update for the trainable matrix  $\hat{A}^t$  in the Left Sketch case. Taking a single **GD** step on the subproblem corresponds to minimizing a quadratic approximation of the objective. This yields the solution for  $\hat{A}^t$ :

$$\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \nabla f(W^t),$$

where  $\eta$  is a learning rate for the subproblem and  $\dagger$  denotes the Moore-Penrose pseudoinverse. Substituting this into the update for  $W^{t+1}$  gives:

$$\begin{aligned} W^{t+1} &= W^t + \frac{\alpha}{r} B_S^t \hat{A}^t = W^t - \frac{\alpha\eta}{r} B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \nabla f(W^t) \\ &= W^t - \gamma H_B^t \nabla f(W^t), \end{aligned}$$

where we define the effective stepsize  $\gamma := \frac{\alpha\eta}{r}$  and the projection matrix  $H_B^t := B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top$ . A similar derivation for the Right Sketch case gives the update:

$$W^{t+1} = W^t - \gamma \nabla f(W^t) H_A^t,$$

where  $H_A^t := (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t$ . This reformulation reveals that both Left and Right sketch updates are equivalent to applying a standard gradient-based update, but projected onto a randomly chosen low-rank subspace.

While **RAC-LoRA** employs a deterministic choice for which matrix to update, our **Bernoulli-LoRA** framework generalizes this concept by introducing a probabilistic selection at each step. This allows us to express the update for any of our proposed methods in a single, unified form:

$$W^{t+1} = W^t - \gamma \hat{G}^t, \tag{7}$$

where  $\hat{G}^t$  is the *projected gradient estimator*. It is formed by taking a *base gradient estimator*  $G^t$  (e.g., a full gradient, a stochastic gradient, or a variance-reduced one) and projecting it based on the outcome of a Bernoulli trial:

$$\hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}. \tag{8}$$

The specific choice of the base estimator  $G^t$  defines the particular algorithm within the **Bernoulli-LoRA** family. We summarize our proposed methods in Table 2 and describe them next.

## 6.2 Core Algorithmic Variants

**Bernoulli-LoRA-GD.** The simplest instantiation of our framework is **Bernoulli-LoRA-GD** (Algorithm 2). This method serves as a foundational building block and a starting point for more elaborate variants. It uses the full gradient of the objective function as its base estimator, i.e.,  $G^t = \nabla f(W^t)$ . While impractical for large-scale deep learning, its analysis provides crucial insights into the convergence behavior of the Bernoulli-LoRA mechanism under idealized, deterministic conditions.

**Bernoulli-LoRA-SGD.** **Stochastic Gradient Descent (SGD)** [[Robbins and Monro, 1951](#)] is a highly effective and widely utilized algorithm for training a variety of machine learning models. The latest advancements in deep learning training methods are all based on different variations of **SGD** [[Sun, 2020](#)]. Its advantage over **GD** is that it uses stochastic gradients for updates, rather than relying on full gradients. Within our framework, we develop **Bernoulli-LoRA-SGD**, where the base estimator  $G^t$  is a general unbiased stochastic gradient of  $f$  at  $W^t$ .

Setting	Method	Base Gradient Estimator $G^t$	Thms. #
(1)	<b>Bernoulli-LoRA-GD</b> (Algs. 2 & 3)	$G^t = \nabla f(W^t)$	1 & 9 & 10
(1)	<b>Bernoulli-LoRA-SGD</b> (Alg. 4)	$G^t = g(W^t)$	11 & 12
(1)+(3)	<b>Bernoulli-LoRA-MVR</b> (Alg. 5)	$G^t = \nabla f_{\xi^t}(W^t) + (1 - b)(G^{t-1} - \nabla f_{\xi^t}(W^{t-1}))$	3 & 14
(1)+(2)	<b>Bernoulli-LoRA-PAGE</b> (Alg. 6)	$G^t = \begin{cases} \nabla f(W^t), & \text{w.p. } q \\ G^{t-1} + \nabla f_{i_t}(W^t) - \nabla f_{i_t}(W^{t-1}), & \text{w.p. } 1 - q \end{cases}$	4 & 16
(1)+(4)	<b>Fed-Bernoulli-LoRA-QGD</b> (Alg. 7)	$G^t = \frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t(\nabla f_l(W^t))$	17 & 18
(1)+(4)	<b>Fed-Bernoulli-LoRA-MARINA</b> (Alg. 8)	$\forall l : G_l^t = \begin{cases} \nabla f_l(W^t), & \text{w.p. } q \\ G_l^{t-1} + \mathcal{Q}_l^t(\nabla f_l(W^t) - \nabla f_l(W^{t-1})), & \text{w.p. } 1 - q \end{cases}$ $G^t = \frac{1}{M} \sum_{l=1}^M G_l^t$	6 & 20
(1)+(4)	<b>Fed-Bernoulli-LoRA-EF21</b> (Alg. 9)	$\forall l : G_l^t = G_l^{t-1} + \mathcal{C}_l^t(\nabla f_l(W^t) - G_l^{t-1})$ $G^t = \frac{1}{M} \sum_{l=1}^M G_l^t$	7 & 22

Table 2: Description of the methods developed and analyzed in this paper. All methods follow the general update rule  $W^{t+1} = W^t - \gamma \hat{G}^t$ , where the projected estimator  $\hat{G}^t$  is defined in (8). The table specifies the definition of the base gradient estimator  $G^t$  for each method. The projection matrices are  $H_A^t := (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t$  and  $H_B^t := B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top$ .

**Bernoulli-LoRA-PAGE.** Several optimal algorithms exist for addressing non-convex optimization problems, such as SPIDER [Fang et al., 2018] and SARAH [Pham et al., 2020]. However, their optimality is supported by a known lower bound that applies only in the small data setting. In contrast, ProbAbilistic Gradient Estimator (PAGE) [Li et al., 2021] stands out for its simplicity, ease of implementation, and ability to achieve optimal convergence in non-convex optimization. PAGE alternates between a full gradient update with probability  $q_t$  and a low-cost gradient adjustment with probability  $1 - q_t$ . Bernoulli-LoRA-PAGE is a new method based on PAGE within our Bernoulli-LoRA framework.

**Bernoulli-LoRA-MVR.** VR methods outperform SGD in reaching first-order critical points but often require finely tuned learning rates and large batch sizes to be effective. To overcome these challenges, Momentum Variance Reduction (MVR) [Cutkosky and Orabona, 2019] was introduced for server-only stochastic non-convex optimization. MVR uses a modified momentum technique to reduce variance without relying on large batch sizes. Several works employ this powerful approach [Tyurin and Richtárik, 2023, Karagulyan et al., 2024]. We propose Bernoulli-LoRA-MVR, where the base estimator  $G^t$  is updated using the MVR rule: a combination of the current stochastic gradient and a momentum term that incorporates the difference between past estimators and gradients.

### 6.3 Extensions for Federated Learning

Sun et al. [2024] identified instability in LoRA, arising from the mismatch between local clients simultaneously optimizing two low-rank matrices and the central server aggregating them independently. Factors such as data heterogeneity, multi-step local updates, and the amplification of additive noise applied to gradients for ensuring differential privacy (DP) significantly impact the process. Additionally, the final performance is highly sensitive to hyperparameter choices. Their proposed solution centers on keeping the randomly initialized non-zero matrices fixed while exclusively fine-tuning the zero-initialized ones. Based on this asymmetric approach, Malinovsky et al. [2024] proposed a distributed method Fed-RAC-LoRA.

We develop the theory further by incorporating compression, VR and EF techniques into FL methods for LoRA within the novel Bernoulli-LoRA framework.

The effectiveness of a distributed training method is primarily measured by its communication complexity, defined as the product of the required communication rounds and the communication volume per round. Following common practice, we assume client-to-server communication is the main bottleneck and exclude server-to-client communication from our analysis.

**Fed-Bernoulli-LoRA-QGD.** A key challenge for distributed methods lies in the high communication cost of gradient updates. Lossy compression techniques, such as QSGD [Alistarh et al., 2017], address this by enabling clients to send quantized gradients. We design Fed-Bernoulli-LoRA-QGD based on QSGD. The clients send compressed versions of their gradients. The base estimator  $G^t$  is formed by averaging the compressed local gradients received from all clients.

**Fed-Bernoulli-LoRA-MARINA.** MARINA [Gorbunov et al., 2021] is a communication-efficient method for non-convex distributed learning on heterogeneous datasets that uses a novel gradient difference compression strategy. Its biased gradient estimator underpins its strong theoretical and practical performance, with proven communication complexity bounds surpassing all prior first-order methods. We propose Fed-Bernoulli-LoRA-MARINA, where each client’s local estimator  $G_l^t$  is updated either with a full local gradient (with probability  $q$ ) or by adding a compressed gradient difference to its previous estimator. The server’s base estimator  $G^t$  is the average of these local estimators.

**Fed-Bernoulli-LoRA-EF21.** Error Feedback (EF) [Seide et al., 2014, Stich et al., 2018, Alistarh et al., 2018, Richtárik et al., 2021] is a widely adopted technique for stabilizing training with contractive compressors. We propose Fed-Bernoulli-LoRA-EF21, based on the modern EF21. Here, each client updates its local estimator  $G_l^t$  by adding a compressed version of the difference between the current local gradient and the previous local estimator. The server’s base estimator  $G^t$  is again the average of the clients’ estimators.

## 7 Convergence Results

The convergence properties of our framework hinge on the spectral properties of the expected projection matrix, which is introduced in Section 6.1. The magnitude of its eigenvalues, particularly the smallest (and in some cases, the largest), is a crucial factor that governs the optimization dynamics.

**Assumption 1.** (Positive Expected Projection) Consider a projection matrix  $H$  generated through either Left Sketch (Definition 1) or Right Sketch (Definition 2). For the sampling distributions  $\mathcal{D}_S^B$  and  $\mathcal{D}_S^A$ , the smallest eigenvalue of the expected projection matrix is strictly positive:

$$\lambda_{\min}^H = \lambda_{\min}[\mathbb{E}[H]] > 0.$$

**Assumption 2.** (Lower Bounded Function) The objective function  $f$  has a finite infimum  $f^* \in \mathbb{R}$ .

**Remark 1** (On the Practicality of Assumption 1). Assumption 1 is a mild and standard requirement, as it is satisfied by common practical choices for the sampling distributions  $\mathcal{D}_S^B$  and  $\mathcal{D}_S^A$ . For instance, a prevalent strategy [Xia et al., 2024, Mao et al., 2025] is to sample the entries of the fixed matrix from an i.i.d. Gaussian distribution. As shown in Appendix B (Lemma 2), this choice leads to an expected projection matrix  $\mathbb{E}[H] = \frac{r}{n} I_n$ , where  $r$  is the rank and  $n$  is the relevant dimension. Consequently,  $\lambda_{\min}^H = \frac{r}{n} > 0$ , readily satisfying the assumption.

Following classical optimization literature [Nemirovski et al., 2009, Beck, 2017, Duchi, 2018, Lan, 2020, Drusvyatskiy, 2020, Nesterov, 2018], we characterize convergence guarantees for two distinct

settings. In the non-smooth convex case, our objective is to find an  $\varepsilon$ -suboptimal solution: a random matrix  $\hat{W} \in \mathbb{R}^{m \times n}$  that satisfies

$$\mathbb{E} [f(\hat{W}) - f(W^*)] \leq \varepsilon, \quad (9)$$

where  $\mathbb{E} [\cdot]$  denotes the expectation with respect to the algorithm's randomness, and  $W^*$  is any minimizer of  $f$ . This same measure of performance—function value suboptimality—is also used to characterize convergence under the Polyak-Łojasiewicz condition, which we introduce later. For the smooth non-convex setting, where finding global minima is generally intractable, we instead aim to locate an  $\varepsilon$ -stationary point: a random matrix  $\hat{W} \in \mathbb{R}^{m \times n}$  satisfying

$$\mathbb{E} \left[ \left\| \nabla f(\hat{W}) \right\|_F^2 \right] \leq \varepsilon^2. \quad (10)$$

This condition guarantees that the expected squared norm of the gradient at our solution is sufficiently small, indicating proximity to a stationary point. To quantify the efficiency of our algorithms, we analyze their iteration complexity—the number of iterations required to achieve these criteria.

A fundamental assumption in the convergence analysis of gradient-based optimization is the Lipschitz continuity of the gradient [Bubeck, 2015, Nesterov, 2018, Beck, 2017, Demidovich et al., 2023b, Khaled and Richtárik, 2023]. This property, often referred to as Lipschitz smoothness, ensures the stability of the optimization trajectory and plays a crucial role in establishing convergence rates [Bottou et al., 2018, Sun, 2020].

**Assumption 3.** (*Lipschitz Smooth Gradient*) *A function  $f$  is differentiable, and there exists a constant  $L > 0$  such that*

$$\|\nabla f(W) - \nabla f(V)\|_F \leq L \|W - V\|_F,$$

for all  $W, V \in \mathbb{R}^{m \times n}$ .

A significant challenge arises when applying LoRA adaptation directly: the Lipschitz smoothness property is not preserved. Specifically, even if a function  $f(W)$  satisfies Assumption 3, its composition with the LoRA parameterization,  $f(W^0 + BA)$ , generally fails to maintain Lipschitz smoothness with respect to the variables  $\{B, A\}$ . This breakdown complicates the analysis of standard gradient-based methods when applied directly to the LoRA parameterization, as formally demonstrated by Sun et al. [2024]. Our framework, by reformulating the updates as projected steps on the full parameter space, circumvents this issue.

To unify our analysis, we define a probability-weighted eigenvalue  $\lambda_{\min(\max)}^p := p\lambda_{\min(\max)}^{H_B} + (1 - p)\lambda_{\min(\max)}^{H_A}$ . Let  $\widetilde{W}^T$  be an iterate drawn randomly from the sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ , with the specific sampling distribution depending on the method.

We begin by presenting the convergence result for the foundational Bernoulli-LoRA-GD method. The proof can be found in Appendix C.1.

**Theorem 1** (Smooth Non-Convex Setting). *Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L}$ . Then the iterates of Bernoulli-LoRA-GD (Algorithm 2), with matrices  $\hat{A}^t$  and  $\hat{B}^t$  computed according to Lemma 3, satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2\Delta^0}{\gamma\lambda_{\min}^p T},$$

where  $\Delta^0 := f(W^0) - f^*$ .

While insightful, full-gradient methods are often impractical for large-scale problems. We therefore extend our analysis to the stochastic setting, where the gradient is replaced by an unbiased estimator  $g(W)$ . For this, we use the general *expected smoothness* assumption.

**Assumption 4** (Expected Smoothness [Khaled and Richtárik, 2023]). *The stochastic gradient estimator  $g(W)$  satisfies*

$$\mathbb{E} \left[ \|g(W)\|_F^2 \right] \leq 2A_1 (f(W) - f^*) + B_1 \cdot \|\nabla f(W)\|_F^2 + C_1,$$

for some constants  $A_1, B_1, C_1 \geq 0$  and all  $W \in \mathbb{R}^{m \times n}$ .

The following theorem establishes the convergence for **Bernoulli-LoRA-SGD**. Its proof is in Appendix C.2.

**Theorem 2.** *Let Assumptions 2, 3, and 4 hold, and let the stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{1}{\sqrt{LA_1 \lambda_{\max}^p T}}, \frac{1}{LB_1} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}.$$

*Then the iterates generated by Bernoulli-LoRA-SGD (Algorithm 4) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{6\Delta^0}{\gamma \lambda_{\min}^p T} + \gamma LC_1 \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$ .

To analyze our variance-reduced methods, we also consider the more specific bounded variance assumption.

**Assumption 5** (Bounded Variance [Nemirovski et al., 2009]). *There exists a constant  $\sigma > 0$  such that, for all  $W \in \mathbb{R}^{m \times n}$ ,*

$$\begin{aligned} \mathbb{E} [\nabla f_\xi(W)] &= \nabla f(W), \\ \mathbb{E} \left[ \|\nabla f_\xi(W) - \nabla f(W)\|_F^2 \right] &\leq \sigma^2. \end{aligned}$$

The next result establishes convergence for **Bernoulli-LoRA-MVR**. The proof is in Appendix C.3.

**Theorem 3.** *Let Assumptions 1, 2, 3, and 5 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{2\lambda_{\max}^p(1-b)^2}{b}} \right)}$ .*

*Then the iterates of Bernoulli-LoRA-MVR (Algorithm 5) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2\Delta^0}{\gamma \lambda_{\min}^p T} + \left( \frac{\mathcal{G}^0}{bT} + \frac{2b\sigma^2}{2-b} \right) \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$  and  $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_F^2$ .

For the finite-sum setting, we analyze **Bernoulli-LoRA-PAGE**, with its convergence detailed in the following theorem and proven in Appendix C.4.

**Theorem 4.** *Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{1-q}{q} \lambda_{\max}^p} \right)}$ .*

*Then the iterates of Bernoulli-LoRA-PAGE (Algorithm 6) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2\Delta^0}{\gamma \lambda_{\min}^p T} + \frac{\mathcal{G}^0}{qT} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$  and  $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_F^2$ .

We now shift to our Federated Learning variants. The following theorem provides convergence guarantees for **Fed-Bernoulli-LoRA-QGD**, with the proof available in Appendix D.1.

**Theorem 5.** *Let Assumptions 1, 2, 3, and 11 hold, and let the stepsize satisfy*

*$0 < \gamma \leq \min \left\{ \frac{1}{L\sqrt{\frac{\omega}{M}\lambda_{\max}^p T}}, \frac{1}{L} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}$ . Then the iterates of **Fed-Bernoulli-LoRA-QGD** (Algorithm 7) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{6\Delta^0}{\gamma\lambda_{\min}^p T} + \frac{2\gamma L\omega\Delta^*}{M} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$ .

Next, we present the convergence result for **Fed-Bernoulli-LoRA-MARINA**. The proof can be found in Appendix D.2.

**Theorem 6.** *Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L(1+\sqrt{\lambda_{\max}^p \frac{1-q}{q} \cdot \frac{\omega}{M}})}$ . Then the iterates of **Fed-Bernoulli-LoRA-MARINA** (Algorithm 8) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2\Delta^0}{\gamma\lambda_{\min}^p T} + \frac{\mathcal{G}^0}{qT} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$  and  $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_F^2$ .

The convergence of **Fed-Bernoulli-LoRA-EF21** is established below, with a detailed proof in Appendix D.3.

**Theorem 7.** *Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L \left( 1 + \frac{\sqrt{\lambda_{\max}^p (1-\beta)}}{1-\sqrt{1-\beta}} \right)}$ . Then the iterates of **Fed-Bernoulli-LoRA-EF21** (Algorithm 9) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2\Delta^0}{\gamma\lambda_{\min}^p T} + \frac{2\hat{\mathcal{G}}^0}{\beta T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$  and  $\hat{\mathcal{G}}^0 := \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_F^2$ .

To obtain stronger, linear convergence rates, we introduce the Polyak–Łojasiewicz condition, a common generalization of strong convexity.

**Assumption 6** (Polyak–Łojasiewicz condition [Polyak, 1963, Łojasiewicz, 1963]). *There exists  $\mu > 0$  such that*

$$\frac{1}{2} \|\nabla f(W)\|_F^2 \geq \mu (f(W) - f^*).$$

The next theorem states the convergence of **Bernoulli-LoRA-SGD** under this condition. It is proven in Appendix C.2.

**Theorem 8.** *Let Assumptions 2, 3, 4, and 6 hold, and let the stepsize satisfy*

*$0 < \gamma \leq \min \left\{ \frac{\mu\lambda_{\min}^p}{2LA_1\lambda_{\max}^p}, \frac{2}{\mu\lambda_{\min}^p}, \frac{1}{LB_1} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}$ . Then the iterates of **Bernoulli-LoRA-SGD** (Algorithm 4) satisfy*

$$\mathbb{E} [f(W^T) - f^*] \leq \left( 1 - \frac{\gamma\mu\lambda_{\min}^p}{2} \right)^T \Delta^0 + \frac{\gamma LC_1}{\mu} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\Delta^0 := f(W^0) - f^*$ .

All other PL-condition results are relegated to the Appendix.



## 8 Experiments

To validate our theoretical findings, we conducted numerical experiments across multiple machine learning tasks.

### 8.1 Linear Regression with Non-convex Regularization.

We begin with a controlled linear regression problem with non-convex regularization, split into pre-training and fine-tuning phases. We use  $\widetilde{(\cdot)}$  for pre-training quantities and  $\hat{(\cdot)}$  for fine-tuning. During the **pre-training phase**, we solve

$$\min_{x \in \mathbb{R}^n} \left\{ \widetilde{f}(x) := \frac{1}{2\widetilde{m}} \left\| \widetilde{D}x - \widetilde{b} \right\|_2^2 + \widetilde{\lambda} \sum_{j=1}^d \frac{x_j^2}{1+x_j^2} \right\}, \quad (11)$$

where  $\widetilde{D} \in \mathbb{R}^{\widetilde{m} \times n}$ ,  $\widetilde{b} \in \mathbb{R}^{\widetilde{m}}$ ,  $\widetilde{m} = 9 \times 10^4$ , and  $n = 4096$ . We set  $\widetilde{\lambda} = \left\| \widetilde{D} \right\|_2 \approx 18.2$ , giving  $\widetilde{L} \approx 54.7$ . We optimize until  $\left\| \nabla \widetilde{f}(\widetilde{x}^*) \right\|^2 \leq 10^{-8}$  to obtain  $\widetilde{x}^*$ . For the **fine-tuning phase**, we use  $\widetilde{x}^*$  as the initialization and then solve

$$\min_{x \in \mathbb{R}^n} \left\{ \hat{f}(x) := \frac{1}{2\hat{m}} \left\| \hat{D}x - \hat{b} \right\|_2^2 + \hat{\lambda} \sum_{j=1}^d \frac{x_j^2}{1+x_j^2} \right\}, \quad (12)$$

where  $\hat{D} \in \mathbb{R}^{\hat{m} \times n}$ ,  $\hat{b} \in \mathbb{R}^{\hat{m}}$ , and  $\hat{m} = 10^4$ . We keep  $n = 4096$  and set  $\hat{\lambda} = \left\| \hat{D} \right\|_2 \approx 4101.7$ , yielding  $\hat{L} \approx 12305.3$ . This second phase uses a dataset with notably different characteristics to mirror realistic domain shifts.

**Stochastic setting.** We consider the stochastic setting, comparing **RAC-LoRA-SGD**, **Bernoulli-LoRA-SGD**, and **Bernoulli-LoRA-PAGE**. In all experiments, we use a batch size of 100, which corresponds to 1% of the data.

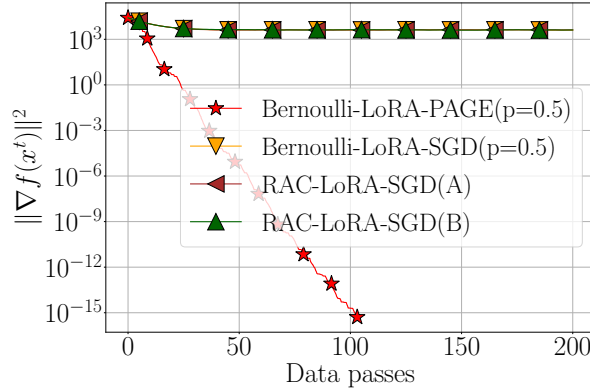


Figure 1: Comparison of **RAC-LoRA-SGD**, **Bernoulli-LoRA-SGD** and **Bernoulli-LoRA-PAGE** on linear regression fine-tuning. Curves with  $p = 0.01, 0.2, \dots$  indicate **Bernoulli-LoRA** sampling parameters. **RAC-LoRA-SGD(A)** trains  $B$  after resampling  $A$ , while **RAC-LoRA-SGD(B)** does the reverse. All methods use  $\gamma = c/\hat{L}$  with  $c$  tuned individually.

Figure 1 shows that **Bernoulli-LoRA-PAGE** successfully reduces variance and converges to the target tolerance, whereas all **SGD** variants stall at a certain accuracy. This underscores the practical advantage of **Bernoulli-LoRA-PAGE** over the baseline **RAC-LoRA-SGD** in the stochastic setting from an optimization standpoint.

## 8.2 MLP on MNIST

In this section, we evaluate **Bernoulli-LoRA** against established baselines in parameter-efficient fine-tuning, following the setup of [Malinovsky et al. \[2024\]](#). Source code of our experiments is available at [https://github.com/IgorSokoloff/Bernoulli-LoRA\\_experiments](https://github.com/IgorSokoloff/Bernoulli-LoRA_experiments).

**Methodology.** We first pre-train a three-layer MLP on MNIST digits 0–4 [[LeCun et al., 1998](#)], then adapt it with various **LoRA**-type methods to classify digits 5–9. Only unseen classes are used for evaluation. All adaptations use rank  $r = 1$  and train for 50 epochs with **AdamW** [[Loshchilov, 2017](#)] ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ ), a fixed learning rate of  $2 \times 10^{-4}$ , and batch size 128. Each method is run 20 times using different seeds, and Table 3 reports the median accuracy (with standard deviation). For **Bernoulli-LoRA**, we show the best median accuracy among all tested settings.

Method	$\mathcal{D}_A$	$\mathcal{D}_B$	Acc. (test)	Train Params
FPFT	-	-	99.5	54,700
LoRA	Gaussian	Zero	$85.69 \pm 1.60$	1K
LoRA	Zero	Gaussian	$89.82 \pm 0.90$	1K
COLA	Gaussian	Zero	$93.32 \pm 0.50$	1K
COLA	Zero	Gaussian	$96.55 \pm 0.20$	1K
AsymmLoRA	Gaussian	Zero	$64.04 \pm 6.90$	133
AsymmLoRA	Zero	Gaussian	$74.52 \pm 7.20$	912
RAC-LoRA	Gaussian	Zero	$93.02 \pm 0.50$	133
RAC-LoRA	Zero	Gaussian	$96.49 \pm 0.20$	912
Bernoulli-LoRA <sup>2</sup>	Zero <sup>1</sup>	Gaussian	$96.46 \pm 0.17$	$\approx 904$

<sup>1</sup> Although **Bernoulli-LoRA** prescribes probabilistic selection from the first iteration, a deterministic assignment of fixed and trainable matrices at initialization yielded better performance.

<sup>2</sup> Achieved with  $p = 0.99$ , giving an expected trainable parameter count  $p \cdot 912 + (1 - p) \cdot 133 \approx 904$ . Here, 912 and 133 are the parameter counts for matrices  $A$  and  $B$ , respectively.

Table 3: Performance on MNIST classification using an MLP with rank  $r$  and scaling  $\alpha = 1$ . For **AsymmLoRA** and **RAC-LoRA**, only the zero-initialized matrix is trained.

**Discussion.** From Table 3, standard **LoRA** attains roughly 86% of full-parameter fine-tuning (FPFT) accuracy, indicating room for improvements via chaining. **COLA** improves upon vanilla **LoRA**, though both lack formal convergence guarantees. **AsymmLoRA** approximates **LoRA** in practice [[Sun et al., 2024](#)] but similarly lacks convergence analysis, whereas **RAC-LoRA** and **Bernoulli-LoRA** both boost accuracy and have theoretical backing. Notably, **Bernoulli-LoRA** matches **RAC-LoRA** in generalization and also guarantees convergence. An additional benefit is that **RAC-LoRA** and **Bernoulli-LoRA** each train only one matrix per **LoRA** block, whereas **COLA** needs two. In **RAC-LoRA**, either  $A$  or  $B$  is trained deterministically; in **Bernoulli-LoRA**, the choice is probabilistic, yielding an expected  $pmr + (1 - p)rn$  trainable parameters. This advantage is especially valuable in resource-constrained settings such as Federated Learning.

Detailed configurations, hardware specs, and dataset descriptions are provided in Appendix E.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## **Acknowledgements**

The research reported in this publication was supported by funding from King Abdullah University of Science and Technology (KAUST): i) KAUST Baseline Research Scheme, ii) CRG Grant ORFS-CRG12-2024-6460, iii) Center of Excellence for Generative AI, under award number 5940, and iv) SDAIA-KAUST Center of Excellence in Artificial Intelligence and Data Science.

## References

- Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. *arXiv preprint arXiv:2012.13255*, 2020.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in Neural Information Processing Systems*, 30, 2017.
- Dan Alistarh, Torsten Hoefler, Mikael Johansson, Nikola Konstantinov, Sarit Khirirat, and Cédric Renggli. The convergence of sparsified gradient methods. *Advances in Neural Information Processing Systems*, 31, 2018.
- Amir Beck. *First-order methods in optimization*. SIAM, 2017.
- Dan Biderman, Jacob Portes, Jose Javier Gonzalez Ortiz, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. LoRA learns less and forgets less. *arXiv preprint arXiv:2405.09673*, 2024.
- Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM review*, 60(2):223–311, 2018.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Daria Cherniuk, Aleksandr Mikhalev, and Ivan Oseledets. Run lora run: Faster and lighter lora implementations. *arXiv preprint arXiv:2312.03415*, 2023.
- Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex SGD. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yury Demidovich, Grigory Malinovsky, Egor Shulgin, and Peter Richtárik. MAST: Model-agnostic sparsified training. *arXiv preprint arXiv:2311.16086*, 2023a.
- Yury Demidovich, Grigory Malinovsky, Igor Sokolov, and Peter Richtárik. A guide through the zoo of biased SGD. *Advances in Neural Information Processing Systems*, 36:23158–23171, 2023b.
- Yury Demidovich, Grigory Malinovsky, and Peter Richtárik. Streamlining in the riemannian realm: Efficient riemannian optimization with loopless variance reduction. *arXiv preprint arXiv:2403.06677*, 2024a.
- Yury Demidovich, Petr Ostroukhov, Grigory Malinovsky, Samuel Horváth, Martin Takáč, Peter Richtárik, and Eduard Gorbunov. Methods with local steps and random reshuffling for generally smooth non-convex federated optimization. *arXiv preprint arXiv:2412.02781*, 2024b.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-[]1423. URL [https://doi.org/10.18653/v1/n19-\[\]1423](https://doi.org/10.18653/v1/n19-[]1423).

- Dmitriy Drusvyatskiy. Convex analysis and nonsmooth optimization. *University Lecture*, 2020.
- John C Duchi. Introductory lectures on stochastic optimization. *The mathematics of data*, 25:99–186, 2018.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. SPIDER: near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in Neural Information Processing Systems*, 31, 2018.
- Ilyas Fatkhullin, Igor Sokolov, Eduard Gorbunov, Zhize Li, and Peter Richtárik. EF21 with bells & whistles: Practical algorithmic extensions of modern error feedback. *arXiv preprint arXiv:2110.03294*, 2021.
- Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. MARINA: Faster non-convex distributed learning with compression. In *International Conference on Machine Learning*, pages 3788–3798. PMLR, 2021.
- Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. SGD: General analysis and improved rates. In *International conference on machine learning*, pages 5200–5209. PMLR, 2019.
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training ImageNet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Brian C Hall. Lie groups, lie algebras, and representations. In *Quantum Theory for Mathematicians*, pages 333–366. Springer, 2013.
- Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. The impact of initialization on LoRA finetuning dynamics. *Advances in Neural Information Processing Systems*, 37:117015–117040, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.
- Samuel Horváth, Chen-Yu Ho, Ludovít Horváth, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. In *Mathematical and Scientific Machine Learning*, pages 129–141. PMLR, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. LoRA: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Póczos, and Alexander J Smola. On variance reduction in stochastic gradient descent and its asynchronous variants. *Advances in Neural Information Processing Systems*, 28, 2015.
- Peter Kairouz, H. B. McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary B. Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Salim Y. El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badi Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaïd Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Oluwasanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek

- Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, R. Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Xiaodong Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning. *Found. Trends Mach. Learn.*, 14:1–210, 2019.
- Avetik Karagulyan, Egor Shulgin, Abdurakhmon Sadiev, and Peter Richtárik. SPAM: Stochastic proximal point method with momentum variance reduction for non-convex cross-device federated learning. *arXiv preprint arXiv:2405.20127*, 2024.
- Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=AU4qHN2VkS>. Survey Certification.
- Ahmed Khaled, Othmane Sebbouh, Nicolas Loizou, Robert M Gower, and Peter Richtárik. Unified analysis of stochastic gradient methods for composite convex and smooth optimization. *Journal of Optimization Theory and Applications*, 199(2):499–540, 2023.
- Sarit Khirirat, Abdurakhmon Sadiev, Artem Riabinin, Eduard Gorbunov, and Peter Richtárik. Error feedback under  $(L_0, L_1)$ -smoothness: Normalization and momentum. *arXiv preprint arXiv:2410.16871*, 2024.
- Mikhail Khodak, Renbo Tu, Tian Li, Liam Li, Maria-Florina F Balcan, Virginia Smith, and Ameet Talwalkar. Federated hyperparameter tuning: Challenges, baselines, and connections to weight-sharing. *Advances in Neural Information Processing Systems*, 34:19184–19197, 2021.
- Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 491–507. Springer, 2020.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: distributed machine learning for on-device intelligence. *arXiv:1610.02527*, 2016.
- Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. FederatedScope-LLM: A comprehensive package for fine-tuning large language models in federated learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5260–5271, 2024.
- Guanghui Lan. *First-order and stochastic optimization methods for machine learning*, volume 1. Springer, 2020.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.



- Chuan Li. Demystifying gpt-3 language model: A technical overview, 2020. URL <https://lambdalabs.com/blog/demystifying-gpt-3>.
- Chunyuan Li, Heerad Farkhor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension of objective landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
- Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. PAGE: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International conference on machine learning*, pages 6286–6295. PMLR, 2021.
- Stanislaw Lojasiewicz. A topological property of real analytic subsets. *Coll. du CNRS, Les équations aux dérivées partielles*, 117(87-89):2, 1963.
- I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Grigory Malinovsky, Kai Yi, and Peter Richtárik. Variance reduced ProxSkip: Algorithm, theory and application to federated learning. *Advances in Neural Information Processing Systems*, 35:15176–15189, 2022.
- Grigory Malinovsky, Umberto Michieli, Hasan Abed Al Kader Hammoud, Taha Ceritli, Hayder Elesedy, Mete Ozay, and Peter Richtárik. Randomized asymmetric chain of LoRA: The first meaningful theoretical framework for low-rank adaptation. *arXiv preprint arXiv:2410.08305*, 2024.
- Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. A survey on lora of large language models. *Frontiers of Computer Science*, 19(7):197605, 2025. doi: 10.1007/s11704-024-040663-9. URL [https://journal.hep.com.cn/fcs/EN/abstract/article\\_47717.shtml](https://journal.hep.com.cn/fcs/EN/abstract/article_47717.shtml).
- H. B. McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics*, 2016.
- Fanxu Meng, Zhaohui Wang, and Muhan Zhang. PiSSA: Principal singular values and singular vectors adaptation of large language models. *Advances in Neural Information Processing Systems*, 37:121038–121072, 2024.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Andrei Panferov, Yury Demidovich, Ahmad Rammal, and Peter Richtárik. Correlated quantization for faster nonconvex distributed optimization. *arXiv preprint arXiv:2401.05518*, 2024.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In Marilyn A. Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 2227–2237. Association for Computational Linguistics, 2018. doi: 10.18653/v1/N18-11202. URL <https://doi.org/10.18653/v1/n18-11202>.

- Nhan H Pham, Lam M Nguyen, Dzung T Phan, and Quoc Tran-Dinh. ProxSARAH: An efficient algorithmic framework for stochastic composite nonconvex optimization. *Journal of Machine Learning Research*, 21(110):1–48, 2020.
- Boris Polyak. Gradient methods for the minimisation of functionals. *Ussr Computational Mathematics and Mathematical Physics*, 3:864–878, 1963.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Peter Richtárik and Martin Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 156:433–484, 2016.
- Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. EF21: A new, simpler, theoretically better, and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34:4384–4396, 2021.
- Peter Richtárik, Igor Sokolov, Elnur Gasanov, Ilyas Fatkhullin, Zhize Li, and Eduard Gorbunov. 3PC: Three point compressors for communication-efficient distributed training and a better theory for lazy aggregation. In *International Conference on Machine Learning*, pages 18596–18648. PMLR, 2022.
- Peter Richtárik, Abdurakhmon Sadiev, and Yury Demidovich. A unified theory of stochastic proximal point methods without smoothness. *arXiv preprint arXiv:2405.15941*, 2024.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- Abdurakhmon Sadiev, Grigory Malinovsky, Eduard Gorbunov, Igor Sokolov, Ahmed Khaled, Konstantin Pavlovich Burlachenko, and Peter Richtárik. Don’t compress gradients in random reshuffling: Compress gradient differences. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=CzPtBzgfae>.
- Issai Schur. Neue begründung der theorie der gruppencharaktere. 2024.
- Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- Fanhua Shang, Kaiwen Zhou, Hongying Liu, James Cheng, Ivor W Tsang, Lijun Zhang, Dacheng Tao, and Licheng Jiao. Vr-sgd: A simple stochastic variance reduction method for machine learning. *IEEE Transactions on Knowledge and Data Engineering*, 32(1):188–202, 2018.
- Igor Sokolov and Peter Richtárik. MARINA-P: Superior performance in non-smooth federated optimization with adaptive stepsizes. *arXiv preprint arXiv:2412.17082*, 2024.
- Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. *Advances in Neural Information Processing Systems*, 31, 2018.
- Ruo-Yu Sun. Optimization for deep learning: An overview. *Journal of the Operations Research Society of China*, 8(2):249–294, 2020.
- Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. Improving LoRA in privacy-preserving federated learning. *arXiv preprint arXiv:2403.12313*, 2024.

- Alexander Tyurin and Peter Richtárik. DASHA: Distributed nonconvex optimization with communication compression and optimal oracle complexity. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=VA1YpcNr7ul>.
- Roman Vershynin. High-dimensional probability, 2009.
- Evgeniya Vorontsova, Roland Hildebrand, Alexander Gasnikov, and Fedor Stonyakin. Convex optimization. *arXiv preprint arXiv:2106.01946*, 2021.
- Hanqing Wang, Yixia Li, Shuo Wang, Guanhua Chen, and Yun Chen. MiLoRA: Harnessing minor singular components for parameter-efficient LLM finetuning. *arXiv preprint arXiv:2406.09044*, 2024.
- Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. TernGrad: Ternary gradients to reduce communication in distributed deep learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of LoRA: efficient fine-tuning of language models via residual learning. *arXiv preprint arXiv:2401.04151*, 2024.
- Kai Yi, Timur Kharisov, Igor Sokolov, and Peter Richtárik. Cohort squeeze: Beyond a single communication round per cohort in cross-device federated learning. *arXiv preprint arXiv:2406.01115*, 2024.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training BERT in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Jiacheng Zhu, Kristjan Greenewald, Kimia Nadjahi, Haitz Saez De Ocariz Borde, Rickard Brühl Gabrielson, Leshem Choshen, Marzyeh Ghassemi, Mikhail Yurochkin, and Justin Solomon. Asymmetry in low-rank adapters of foundation models. *arXiv preprint arXiv:2402.16842*, 2024.

# APPENDIX

## A Basic Facts and Useful Inequalities

**Tower property.** For any random variables  $X$  and  $Y$ , we have

$$\mathbb{E}[\mathbb{E}[X \mid Y]] = \mathbb{E}[X]. \quad (13)$$

**Cauchy-Bunyakovsky-Schwarz inequality.** For any random variables  $X$  and  $Y$ , we have

$$|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]}. \quad (14)$$

**Variance decomposition.** For any random vector  $X \in \mathbb{R}^d$  and any non-random  $c \in \mathbb{R}^d$ , we have

$$\mathbb{E}[\|X - c\|_2^2] = \mathbb{E}[\|X - \mathbb{E}[X]\|_2^2] + \|\mathbb{E}[X] - c\|_2^2. \quad (15)$$

**Jensen's inequality.** For any random vector  $X \in \mathbb{R}^d$  and any convex function  $g : \mathbb{R}^d \mapsto \mathbb{R}$ , we have

$$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]. \quad (16)$$

## B Discussion on Positive Expected Projection (Assumption 1)

Assumption 1 merits further discussion. While any single projection matrix has eigenvalues that are either 0 or 1 (with the smallest being 0), the expected value of a *random* projection matrix can have all its eigenvalues strictly greater than zero. This property is crucial for ensuring stable convergence behavior in our framework.

Later in this section, we will utilize the following lemma, which is a classical result from linear algebra, often known as a direct consequence of Schur's Lemma [Hall, 2013, Schur, 2024].

**Lemma 1** (Rotational Invariance Implies Scalar Matrix). *Let  $M \in \mathbb{R}^{n \times n}$  be a matrix satisfying*

$$M = QMQ^\top \quad \text{for all orthonormal matrices } Q \in \mathbb{R}^{n \times n}. \quad (17)$$

*Then  $M = \alpha I_n$  for some scalar  $\alpha \in \mathbb{R}$ .*

*Proof.* The condition  $M = QMQ^\top$  is equivalent to  $MQ = QM$ , which means that  $M$  commutes with every orthonormal matrix  $Q$ . Since  $M$  is a real symmetric matrix, it is guaranteed to have at least one real eigenvector. Let  $v$  be such an eigenvector with corresponding eigenvalue  $\lambda$ . We can normalize this eigenvector to create a unit vector  $u_1 = v/\|v\|$ , which is also an eigenvector with the same eigenvalue:

$$Mu_1 = M\left(\frac{v}{\|v\|}\right) = \frac{1}{\|v\|}Mv = \frac{1}{\|v\|}(\lambda v) = \lambda\left(\frac{v}{\|v\|}\right) = \lambda u_1.$$

Now, let  $u$  be any other arbitrary unit vector in  $\mathbb{R}^n$ . Because both  $u_1$  and  $u$  are unit vectors (i.e., they lie on the unit sphere), there always exists an orthonormal matrix  $Q$  (specifically, a rotation) that maps  $u_1$  to  $u$ . That is,  $u = Qu_1$ .

We now examine the action of  $M$  on this arbitrary unit vector  $u$ :

$$Mu = M(Qu_1) = (MQ)u_1 = (QM)u_1 = Q(Mu_1) = Q(\lambda u_1) = \lambda(Qu_1) = \lambda u.$$

We have shown that any arbitrary unit vector  $u$  is an eigenvector of  $M$  with the same eigenvalue  $\lambda$ . If every unit vector is an eigenvector with eigenvalue  $\lambda$ , then for any non-zero vector  $x \in \mathbb{R}^n$ , we can write  $x = \|x\| \cdot \frac{x}{\|x\|}$ . Let  $u_x := x/\|x\|$  be the corresponding unit vector. Then:

$$Mx = M(\|x\|u_x) = \|x\|(Mu_x) = \|x\|(\lambda u_x) = \lambda(\|x\|u_x) = \lambda x.$$

Since  $Mx = \lambda x$  for all vectors  $x \in \mathbb{R}^n$ , the matrix  $M$  must be a scalar multiple of the identity matrix, i.e.,  $M = \lambda I_n$ .  $\square$

In practice, LoRA-type methods often employ Gaussian sampling for the matrices  $A_S$  or  $B_S$  [Xia et al., 2024, Mao et al., 2025]. The following lemma, a standard result in multivariate statistics, demonstrates that under such Gaussian sampling, Assumption 1 is naturally satisfied.

**Lemma 2** (Expected Eigenvalues of Random Projection Matrices). *Consider a projection matrix  $H_B$  generated by a random matrix  $B \in \mathbb{R}^{n \times r}$  whose entries are i.i.d.  $\mathcal{N}(0,1)$  with  $r \leq n$ , defined as:*

$$H_B = B(B^\top B)^\dagger B^\top,$$

*where  $\dagger$  denotes the Moore-Penrose pseudoinverse. Similarly, for a random matrix  $A \in \mathbb{R}^{r \times n}$  with i.i.d.  $\mathcal{N}(0,1)$  entries, we define:*

$$H_A = A^\top (AA^\top)^\dagger A.$$

For these matrices, we have:

$$\mathbb{E}[H_B] = \mathbb{E}[H_A] = \frac{r}{n} I_n,$$

which implies:

$$\lambda_{\min}(\mathbb{E}[H_B]) = \lambda_{\min}(\mathbb{E}[H_A]) = \frac{r}{n}.$$

*Proof.* The proof leverages the rotational invariance property of the standard Gaussian distribution. We will prove the result for  $H_B$ ; the argument for  $H_A$  is analogous.

First, we establish that  $\mathbb{E}[H_B]$  must be a scalar multiple of the identity matrix. Let  $Q \in \mathbb{R}^{n \times n}$  be an arbitrary orthonormal matrix. Due to the rotational invariance of the multivariate standard normal distribution, the random matrix  $QB$  has the same distribution as  $B$ .

Consider the projection matrix  $H_{QB}$  generated by  $QB$ :

$$\begin{aligned} H_{QB} &= (QB) \left( (QB)^\top QB \right)^\dagger (QB)^\top \\ &= QB \left( B^\top Q^\top QB \right)^\dagger B^\top Q^\top \\ &= QB \left( B^\top B \right)^\dagger B^\top Q^\top \\ &= Q \left( B(B^\top B)^\dagger B^\top \right) Q^\top = QH_B Q^\top. \end{aligned}$$

Since  $QB$  and  $B$  are identically distributed, their expectations must be equal:  $\mathbb{E}[H_{QB}] = \mathbb{E}[H_B]$ . This implies:

$$\mathbb{E}[H_B] = Q\mathbb{E}[H_B]Q^\top,$$

for every orthonormal matrix  $Q$ . By Lemma 1,  $\mathbb{E}[H_B]$  must be a scalar multiple of the identity matrix, so  $\mathbb{E}[H_B] = \alpha I_n$  for some scalar  $\alpha \in \mathbb{R}$ .

To determine this scalar, we use the property that the trace of a projection matrix is equal to its rank. Since the columns of  $B$  are drawn from a continuous distribution, they are linearly independent almost surely (as  $r \leq n$ ). Thus, the rank of  $H_B$  is  $r$ .

$$\mathbb{E}[\text{Tr}(H_B)] = \mathbb{E}[\text{rank}(H_B)] = r.$$

By linearity of expectation and trace, we also have:

$$\mathbb{E}[\text{Tr}(H_B)] = \text{Tr}(\mathbb{E}[H_B]) = \text{Tr}(\alpha I_n) = \alpha n.$$

Equating the two expressions gives  $\alpha n = r$ , which implies  $\alpha = \frac{r}{n}$ . Therefore,

$$\mathbb{E}[H_B] = \frac{r}{n} I_n.$$

The same argument applies to  $H_A$  by observing that  $A^\top$  is an  $n \times r$  matrix with i.i.d.  $\mathcal{N}(0,1)$  entries, which completes the proof.  $\square$

**Remark 2.** This result is foundational in the study of random projections and can be found in standard textbooks on multivariate statistics; for example, see Lemma 5.3.2 in [Vershynin, 2009].



## C Proofs for Core Algorithmic Variants

### C.1 Analysis of Bernoulli-LoRA-GD

---

**Algorithm 2** Bernoulli-LoRA-GD

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling factor  $\alpha > 0$ , stepsize  $\gamma_t$ 
   chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \nabla f(W^t)$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta \nabla f(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13: end for

```

---

The following lemma establishes that the Bernoulli-LoRA update can be reformulated as a standard projected gradient descent step, providing a crucial foundation for our subsequent convergence analysis.

**Lemma 3.** Consider the updates  $\hat{A}^t$  and  $\hat{B}^t$  from Algorithm 2 computed as solutions to the following optimization problems:

$$\begin{aligned}
\hat{A}^t &:= \arg \min_A \left\{ f(W^t) + \frac{\alpha}{r} \langle \nabla f(W^t), B_S^t A \rangle_F + \frac{\alpha^2}{2\gamma r^2} \|B_S^t A\|_F^2 \right\}, \\
\hat{B}^t &:= \arg \min_B \left\{ f(W^t) + \frac{\alpha}{r} \langle \nabla f(W^t), B A_S^t \rangle_F + \frac{\alpha^2}{2\gamma r^2} \|B A_S^t\|_F^2 \right\}.
\end{aligned} \tag{18}$$

Then the Left and Right sketch updates can be expressed as a gradient descent step:

$$W^{t+1} = W^t - \gamma G^t, \tag{19}$$

where  $G^t$  is defined by

$$G^t = \begin{cases} H_B^t \nabla f(W^t), & \text{with probability } p \\ \nabla f(W^t) H_A^t, & \text{with probability } 1 - p \end{cases} \tag{20}$$

with projection matrices  $H_A^t$  and  $H_B^t$  given by:

$$H_A^t := (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t \quad \text{and} \quad H_B^t := B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top, \tag{21}$$

where  $^\dagger$  denotes the Moore-Penrose pseudoinverse.

*Proof.* Following Algorithm 2, at each iteration we randomly select either the Left sketch (with probability  $p$ ) or the Right sketch (with probability  $1 - p$ ). We analyze both cases separately and then combine them into a unified update rule.

**Left Sketch Analysis.** When the Left sketch is selected, the update takes the form:

$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t. \quad (22)$$

Minimizing the right-hand side with respect to  $\hat{A}^t$  yields:

$$\begin{aligned} \frac{\alpha}{r} (B_S^t)^\top \nabla f(W^t) + \frac{\alpha^2}{\gamma r^2} (B_S^t)^\top B_S^t \hat{A}^t &= 0; \\ (B_S^t)^\top B_S^t \hat{A}^t &= -\frac{\gamma r}{\alpha} (B_S^t)^\top \nabla f(W^t); \\ \hat{A}^t &= -\frac{\gamma r}{\alpha} \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \nabla f(W^t). \end{aligned} \quad (23)$$

This leads to the Left sketch update:

$$\begin{aligned} W^{t+1} &= W^t + \frac{\alpha}{r} B_S^t \hat{A}^t \\ &= W^t - \gamma B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \nabla f(W^t) \\ &= W^t - \gamma H_B^t \nabla f(W^t), \end{aligned} \quad (24)$$

where  $H_B^t := B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top$  is a projection matrix.

**Right Sketch Analysis.** For the Right sketch, we follow a similar approach. The update rule is:

$$W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t. \quad (25)$$

First, observe that:

$$\left\| \hat{B}^t A_S^t \right\|_F^2 = \left\langle \hat{B}^t A_S^t, \hat{B}^t A_S^t \right\rangle_F = \left\langle A_S^t, \left( \hat{B}^t \right)^\top \hat{B}^t A_S^t \right\rangle_F. \quad (26)$$

For the linear term from (18):

$$\frac{\alpha}{r} \left\langle \nabla f(W^t), \hat{B}^t A_S^t \right\rangle_F = \frac{\alpha}{r} \text{Tr} \left( \left( \nabla f(W^t) \right)^\top \hat{B}^t A_S^t \right), \quad (27)$$

with gradient  $\nabla f(W^t) (A_S^t)^\top$  with respect to  $\hat{B}^t$ . Using the matrix calculus identity  $\nabla_X \|X\|_F^2 = 2X$ , the gradient of the quadratic term is:

$$\frac{\alpha^2}{\gamma r^2} \hat{B}^t A_S^t (A_S^t)^\top. \quad (28)$$

Setting the total gradient to zero and solving for  $\hat{B}^t$ :

$$\hat{B}^t = -\frac{\gamma r}{\alpha} \nabla f(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger, \quad (29)$$

which yields the Right sketch update:

$$\begin{aligned} W^{t+1} &= W^t + \frac{\alpha}{r} \hat{B}^t A_S^t \\ &= W^t - \gamma \nabla f(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t \\ &= W^t - \gamma \nabla f(W^t) H_A^t, \end{aligned} \quad (30)$$

where  $H_A^t := (A_S^t)^\top (A_S^t (A_S^t)^\top)^\dagger A_S^t$  is a projection matrix.

**Combined Update Rule.** Combining equations (24) and (30), we obtain the unified update:

$$W^{t+1} = W^t - \gamma G^t, \quad (31)$$

where  $G^t$  takes the form given in the lemma statement, completing the proof.  $\square$

With these assumptions in place, we can now state our main convergence result for RAC-LoRA with Gradient Descent updates.

### C.1.1 Convergence for Smooth Non-Convex Functions

**Theorem 1.** Let Assumptions 1, 3, and 2 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L}$ . Then the iterates of Bernoulli-LoRA-GD (Algorithm 2), with matrices  $\hat{A}^t$  and  $\hat{B}^t$  computed according to Lemma 3, satisfy

$$\mathbb{E} \left[ \left\| \nabla f(\tilde{W}^T) \right\|_F^2 \right] \leq \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T}, \quad (32)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$  and  $\tilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .

*Proof.* From Lemma 3, we know that Bernoulli-LoRA updates can be expressed as

$$W^{t+1} = W^t - \gamma G^t, \quad (33)$$

where  $G^t$  takes the form

$$G^t = \begin{cases} H_B^t \nabla f(W^t), & \text{with probability } p \\ \nabla f(W^t) H_A^t, & \text{with probability } 1 - p \end{cases} \quad (34)$$

with projection matrices  $H_A^t$  and  $H_B^t$  as defined in the lemma.

To analyze the convergence, we first compute the conditional expectation and second moment of  $G^t$ :

$$\begin{aligned} \mathbb{E} [G^t \mid W^t, H^t] &= p H_B^t \nabla f(W^t) + (1-p) \nabla f(W^t) H_A^t, \\ \mathbb{E} [\|G^t\|_F^2 \mid W^t, H^t] &= p \|H_B^t \nabla f(W^t)\|_F^2 + (1-p) \|\nabla f(W^t) H_A^t\|_F^2, \end{aligned} \quad (35)$$

where we defined  $H^t := \{H_A^t, H_B^t\}$ .

We begin by establishing several key auxiliary bounds. For the Left sketch term:

$$\begin{aligned} & -\gamma p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F + \frac{L\gamma^2}{2} p \|H_B^t \nabla f(W^t)\|_F^2 \\ &= -\gamma p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F + \frac{L\gamma^2}{2} p \langle H_B^t \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F \\ &= -\gamma p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F + \frac{L\gamma^2}{2} p \langle \nabla f(W^t), (H_B^t)^\top H_B^t \nabla f(W^t) \rangle_F \\ &= p \left( -\gamma \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F + \frac{L\gamma^2}{2} \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F \right) \\ &\stackrel{\gamma \leq 1/L}{\leq} -\frac{\gamma}{2} p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F. \end{aligned} \quad (36)$$

For any projection matrix  $H_A^t$ , we have:

$$\begin{aligned}
\langle \nabla f(W^t) H_A^t, \nabla f(W^t) H_A^t \rangle_F &= \text{Tr} \left( (H_A^t)^\top (\nabla f(W^t))^\top \nabla f(W^t) H_A^t \right) \\
&= \text{Tr} \left( (\nabla f(W^t))^\top \nabla f(W^t) H_A^t (H_A^t)^\top \right) \\
&= \text{Tr} \left( (\nabla f(W^t))^\top \nabla f(W^t) H_A^t \right) \\
&= \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F.
\end{aligned} \tag{37}$$

Therefore:

$$\begin{aligned}
& -\gamma(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F + \frac{L\gamma^2}{2}(1-p) \|\nabla f(W^t) H_A^t\|_F^2 \\
&= -\gamma(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F + \frac{L\gamma^2}{2}(1-p) \langle \nabla f(W^t) H_A^t, \nabla f(W^t) H_A^t \rangle_F \\
&= -\gamma(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F + \frac{L\gamma^2}{2}(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F \\
&\stackrel{\gamma \leq 1/L}{\leq} -\frac{\gamma}{2}(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F.
\end{aligned} \tag{38}$$

Using the Lipschitz gradient condition and the above bounds:

$$\begin{aligned}
\mathbb{E} [f(W^{t+1}) \mid W^t, H^t] &\leq f(W^t) + \mathbb{E} [\langle \nabla f(W^t), W^{t+1} - W^t \rangle_F \mid W^t, H^t] \\
&+ \frac{L}{2} \mathbb{E} [\|W^{t+1} - W^t\|_F^2 \mid W^t, H^t] \\
&= f(W^t) - \gamma \langle \nabla f(W^t), \mathbb{E} [G^t \mid W^t, H^t] \rangle_F + \frac{L\gamma^2}{2} \mathbb{E} [\|G^t\|_F^2 \mid W^t, H^t] \\
&= f(W^t) - \gamma p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F - \gamma(1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F \\
&+ \frac{L\gamma^2}{2} p \|H_B^t \nabla f(W^t)\|_F^2 + \frac{L\gamma^2}{2} (1-p) \|\nabla f(W^t) H_A^t\|_F^2 \\
&\stackrel{(36),(38)}{\leq} f(W^t) - \frac{\gamma}{2} (p \langle \nabla f(W^t), H_B^t \nabla f(W^t) \rangle_F + (1-p) \langle \nabla f(W^t), \nabla f(W^t) H_A^t \rangle_F).
\end{aligned} \tag{39}$$

For the first term:

$$\begin{aligned}
-\langle \nabla f(W^t), \mathbb{E} [H_B^t] \nabla f(W^t) \rangle_F &= -\text{Tr} \left( (\nabla f(W^t))^\top \mathbb{E} [H_B^t] \nabla f(W^t) \right) \\
&\leq -\lambda_{\min} (\mathbb{E} [H_B^t]) \text{Tr} \left( (\nabla f(W^t))^\top \nabla f(W^t) \right) \\
&= -\lambda_{\min}^{H_B} \|\nabla f(W^t)\|_F^2.
\end{aligned} \tag{40}$$

Similarly, for the second term:

$$\begin{aligned}
-\langle \nabla f(W^t), \nabla f(W^t) \mathbb{E} [H_A^t] \rangle_F &= -\text{Tr} \left( (\nabla f(W^t))^\top \nabla f(W^t) \mathbb{E} [H_A^t] \right) \\
&= -\text{Tr} \left( \mathbb{E} [H_A^t] (\nabla f(W^t))^\top \nabla f(W^t) \right) \\
&\leq -\lambda_{\min}^{H_A} \|\nabla f(W^t)\|_F^2.
\end{aligned} \tag{41}$$

Therefore:

$$\begin{aligned}
\mathbb{E} [f(W^{t+1}) | W^t] &= \mathbb{E} [\mathbb{E} [f(W^{t+1}) | W^t, H^t] | W^t] \\
&\leq f(W^t) - \frac{\gamma}{2} (p \langle \nabla f(W^t), \mathbb{E} [H_B^t] \nabla f(W^t) \rangle_F + (1-p) \langle \nabla f(W^t), \nabla f(W^t) \mathbb{E} [H_A^t] \rangle_F) \\
&\leq f(W^t) - \frac{\gamma}{2} (p \lambda_{\min}^{H_B} + (1-p) \lambda_{\min}^{H_A}) \|\nabla f(W^t)\|_F^2 \\
&= f(W^t) - \frac{\gamma}{2} \lambda_{\min}^p \|\nabla f(W^t)\|_F^2,
\end{aligned} \tag{42}$$

where  $\lambda_{\min}^p := p \lambda_{\min}^{H_B} + (1-p) \lambda_{\min}^{H_A}$ . Further,

$$\mathbb{E} [\mathbb{E} [f(W^{t+1}) | W^t, H^t] | W^t] - f^* \leq f(W^t) - f^* - \frac{\gamma}{2} \lambda_{\min}^p \|\nabla f(W^t)\|_F^2. \tag{43}$$

Taking the sum over  $t = 0, \dots, T-1$  and using the tower property of expectation:

$$\mathbb{E} [f(W^T) - f^*] \leq f(W^0) - f^* - \frac{\gamma}{2} \lambda_{\min}^p \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(W^t)\|_F^2]. \tag{44}$$

By rearranging terms, we get:

$$\frac{\gamma}{2} \lambda_{\min}^p \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(W^t)\|_F^2] \leq f(W^0) - f^*. \tag{45}$$

Finally, dividing both sides by  $\frac{\gamma T}{2} \lambda_{\min}^p$  yields:

$$\mathbb{E} [\|\nabla f(\widetilde{W}^T)\|_F^2] \leq \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T}, \tag{46}$$

where  $\widetilde{W}^T$  is chosen uniformly at random from  $\{W^0, W^1, \dots, W^{T-1}\}$ , completing the proof.  $\square$

### C.1.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 9.** Let Assumptions 1, 2, 3, and 6 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L}$ . Then the iterates of Bernoulli-LoRA-GD (Algorithm 2), with matrices  $\hat{A}^t$  and  $\hat{B}^t$  computed according to Lemma 3, satisfy

$$\mathbb{E} [f(W^T) - f^*] \leq (1 - \gamma \mu \lambda_{\min}^p)^T (f(W^0) - f^*),$$

where  $\lambda_{\min}^p := p \lambda_{\min}^{H_B} + (1-p) \lambda_{\min}^{H_A}$ .

*Proof.* We begin our analysis from a key inequality derived in the proof of Theorem 1:

$$\mathbb{E} [f(W^{t+1}) | W^t] \leq f(W^t) - \frac{\gamma}{2} \lambda_{\min}^p \|\nabla f(W^t)\|_F^2. \tag{47}$$

By invoking the Polyak-Łojasiewicz condition (Assumption 6), which states that  $\frac{1}{2} \|\nabla f(W)\|_F^2 \geq \mu (f(W) - f^*)$ , we can further bound the right-hand side of the inequality (47):

$$\mathbb{E} [f(W^{t+1}) | W^t] \leq f(W^t) - \gamma \lambda_{\min}^p (\mu (f(W^t) - f^*)).$$

Subtracting the optimal function value  $f^*$  from both sides, we get a recursive relationship for the expected suboptimality gap:

$$\begin{aligned}\mathbb{E} [f(W^{t+1}) - f^* \mid W^t] &\leq (f(W^t) - f^*) - \gamma\mu\lambda_{\min}^p (f(W^t) - f^*) \\ &= (1 - \gamma\mu\lambda_{\min}^p) (f(W^t) - f^*).\end{aligned}$$

By taking the full expectation over all randomness up to iteration  $t$  and applying the tower property, we obtain:

$$\mathbb{E} [f(W^{t+1}) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E} [f(W^t) - f^*].$$

Unrolling this recursion from  $t = T - 1$  down to  $t = 0$  yields the final linear convergence result:

$$\mathbb{E} [f(W^T) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p)^T (f(W^0) - f^*).$$

This completes the proof.  $\square$

### C.1.3 Convergence for Non-Smooth Convex Functions

---

#### Algorithm 3 Bernoulli-LoRA-GD (Non-smooth setting)

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling factor  $\alpha > 0$ , stepsize  $\gamma_t$ 
   chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = \arg \min_A \left\{ f(W^t) + \frac{\alpha}{r} \langle \partial f(W^t), B_S^t A \rangle_F + \frac{\alpha^2}{2\gamma_t r^2} \|B_S^t A\|_F^2 \right\}$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = \arg \min_B \left\{ f(W^t) + \frac{\alpha}{r} \langle \partial f(W^t), B A_S^t \rangle_F + \frac{\alpha^2}{2\gamma_t r^2} \|B A_S^t\|_F^2 \right\}$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13: end for

```

---

Our analysis relies on the following standard assumptions that are widely used in non-smooth optimization theory:

**Assumption 7.** *The function  $f$  has at least one minimizer, denoted by  $W^*$ .*

**Assumption 8.** *The function  $f$  is convex.*

**Assumption 9** (Lipschitz continuity). *The function  $f$  is  $L_0$ -Lipschitz continuous. That is, there exists  $L_0 > 0$  such that*

$$|f(W) - f(V)| \leq L_0 \|W - V\|_F, \quad \forall W, V \in \mathbb{R}^{m \times n}. \quad (48)$$

The combination of convexity and Lipschitz continuity represents a standard framework in non-smooth optimization [Vorontsova et al., 2021, Nesterov, 2013, Bubeck, 2015, Beck, 2017, Duchi, 2018,



Lan, 2020, Drusvyatskiy, 2020]. Notably, the  $L_0$ -Lipschitz continuity implies uniformly bounded subgradients [Beck, 2017], a property that plays a crucial role in our analysis:

$$\|\partial f(W)\|_F \leq L_0, \quad \forall W \in \mathbb{R}^{m \times n}. \quad (49)$$

This boundedness of subgradients ensures the stability of our optimization process and enables us to establish rigorous convergence guarantees.

The following lemma establishes that the **Bernoulli-LoRA** update in the non-smooth case can also be reformulated as a subgradient descent step, which plays a central role in our convergence analysis for non-smooth objectives.

**Lemma 4.** *Consider the updates  $\hat{A}^t$  and  $\hat{B}^t$  from Algorithm 3 computed as solutions to the following optimization problems:*

$$\begin{aligned} \hat{A}^t &:= \arg \min_A \left\{ f(W^t) + \frac{\alpha}{r} \langle \partial f(W^t), B_S^t A \rangle_F + \frac{\alpha^2}{2\gamma_t r^2} \|B_S^t A\|_F^2 \right\}, \\ \hat{B}^t &:= \arg \min_B \left\{ f(W^t) + \frac{\alpha}{r} \langle \partial f(W^t), B A_S^t \rangle_F + \frac{\alpha^2}{2\gamma_t r^2} \|B A_S^t\|_F^2 \right\}. \end{aligned} \quad (50)$$

Then the Left and Right sketch updates can be expressed as a subgradient descent step:

$$W^{t+1} = W^t - \gamma_t G^t, \quad (51)$$

where  $G^t$  is defined by

$$G^t = \begin{cases} H_B^t \partial f(W^t), & \text{with probability } p \\ \partial f(W^t) H_A^t, & \text{with probability } 1 - p \end{cases} \quad (52)$$

with projection matrices  $H_A^t$  and  $H_B^t$  given by:

$$H_A^t := (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t \quad \text{and} \quad H_B^t := B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top, \quad (53)$$

where  $\dagger$  denotes the Moore-Penrose pseudoinverse.

*Proof.* The proof follows a similar structure to that of Lemma 3, with subgradients replacing gradients throughout the analysis. We examine both sketch types separately before combining them into a unified update rule.

**Left Sketch Analysis.** When the Left sketch is selected, the update takes the form:

$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t. \quad (54)$$

The matrix  $\hat{A}^t$  is defined as the solution to the optimization problem:

$$\hat{A}^t := \arg \min_A \left\{ f(W^t) + \frac{\alpha}{r} \langle \partial f(W^t), B_S^t A \rangle_F + \frac{\alpha^2}{2\gamma_t r^2} \|B_S^t A\|_F^2 \right\}. \quad (55)$$

By computing the gradient of the objective with respect to  $A$  and setting it to zero, we obtain:

$$\begin{aligned} \frac{\alpha}{r} (B_S^t)^\top \partial f(W^t) + \frac{\alpha^2}{\gamma_t r^2} (B_S^t)^\top B_S^t \hat{A}^t &= 0; \\ \hat{A}^t &= -\frac{\gamma_t r}{\alpha} \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \partial f(W^t). \end{aligned} \quad (56)$$

Substituting this expression back into the update equation yields the Left sketch update:

$$\begin{aligned}
W^{t+1} &= W^t + \frac{\alpha}{r} B_S^t \hat{A}^t \\
&= W^t - \gamma_t B_S^t \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top \partial f(W^t) \\
&= W^t - \gamma_t H_B^t \partial f(W^t).
\end{aligned} \tag{57}$$

**Right Sketch Analysis.** For the Right sketch, we follow an analogous approach. The update rule takes the form:

$$W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t. \tag{58}$$

Applying similar optimization steps but now with respect to matrix  $B$ , we obtain:

$$\hat{B}^t = -\frac{\gamma_t r}{\alpha} \partial f(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger, \tag{59}$$

which leads to the Right sketch update:

$$\begin{aligned}
W^{t+1} &= W^t + \frac{\alpha}{r} \hat{B}^t A_S^t \\
&= W^t - \gamma_t \partial f(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t \\
&= W^t - \gamma_t \partial f(W^t) H_A^t.
\end{aligned} \tag{60}$$

**Combined Update Rule.** By combining equations (57) and (60), we arrive at the unified update rule:

$$W^{t+1} = W^t - \gamma_t G^t, \tag{61}$$

where  $G^t$  takes the form specified in the lemma statement, thus completing the proof.  $\square$

**Assumption 10.** Consider a projection matrix  $H$  generated through either Left Sketch (Definition 1) or Right Sketch (Definition 2). For the sampling distributions  $\mathcal{D}_S^B$  and  $\mathcal{D}_S^A$ , the expected projection matrix  $H$  satisfies

$$\mathbb{E}[H] = \alpha I, \tag{62}$$

where a constant  $\alpha > 0$ .

**Theorem 10.** Let Assumptions 1, 7, 8, 9, and 10 hold. Let us define the following quantities:  $\bar{W}^T := \frac{1}{T} \sum_{t=0}^{T-1} W^t$  as the averaged iterate,  $R_0^2 := \|W^0 - W^*\|_F^2$  as the initial distance to optimum. Consider the sequence  $\{W^t\}$  produced by Bernoulli-LoRA (Algorithm 3) with updates of  $\hat{A}^t$  and  $\hat{B}^t$  computed according to Lemma 4.

**1. (Constant stepsize).** If the stepsize is constant, i.e.,  $\gamma_t := \gamma > 0$ , then

$$\mathbb{E} \left[ f(\bar{W}^T) - f(W^*) \right] \leq \frac{R_0^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2}. \tag{63}$$

Moreover, with the optimal stepsize  $\gamma_* = \sqrt{\frac{(R_0^2)^2}{T\alpha L_0^2}}$ , we obtain:

$$\mathbb{E} \left[ f(\bar{W}^T) - f(W^*) \right] \leq \frac{R_0^2 L_0}{\sqrt{\alpha T}}. \tag{64}$$

2. (**Polyak stepsize**). If the stepsize is chosen adaptively as

$$\gamma_t = \frac{(f(W^t) - f(W^*))}{\|\partial f(W^t)\|_F^2}, \quad (65)$$

then

$$\mathbb{E} [f(\bar{W}^T) - f(W^*)] \leq \frac{R^0 L_0}{\sqrt{\alpha T}}. \quad (66)$$

*Proof.* From Lemma 4, we know that Bernoulli-LoRA updates in the non-smooth setting can be expressed as

$$W^{t+1} = W^t - \gamma_t G^t, \quad (67)$$

where  $G^t$  takes the form

$$G^t = \begin{cases} H_B^t \partial f(W^t), & \text{with probability } p \\ \partial f(W^t) H_A^t, & \text{with probability } 1 - p \end{cases} \quad (68)$$

with projection matrices  $H_A^t$  and  $H_B^t$  as defined in the lemma.

To analyze the convergence, we first compute the conditional expectation and second moment of  $G^t$ :

$$\mathbb{E} [G^t \mid W^t, H^t] = p H_B^t \partial f(W^t) + (1 - p) \partial f(W^t) H_A^t, \quad (69)$$

$$\mathbb{E} [\|G^t\|_F^2 \mid W^t, H^t] = p \|H_B^t \partial f(W^t)\|_F^2 + (1 - p) \|\partial f(W^t) H_A^t\|_F^2, \quad (70)$$

where we defined  $H^t := \{H_A^t, H_B^t\}$ .

By the definition of subgradient, we have:

$$f(W^*) \geq f(W^t) + \langle \partial f(W^t), W^* - W^t \rangle_F, \quad (71)$$

which implies:

$$\langle \partial f(W^t), W^t - W^* \rangle_F \geq f(W^t) - f(W^*). \quad (72)$$

Let us establish key auxiliary bounds. First, for the inner product terms:

$$\begin{aligned} -2\gamma_t \mathbb{E} [\langle G^t, W^t - W^* \rangle_F \mid W^t, H^t] &\stackrel{(69)}{=} -2\gamma_t p \langle H_B^t \partial f(W^t), W^t - W^* \rangle_F \\ &\quad -2\gamma_t (1 - p) \langle \partial f(W^t) H_A^t, W^t - W^* \rangle_F. \end{aligned} \quad (73)$$

For projection matrices, we have the following properties:

$$\begin{aligned} \|\partial f(W^t) H_A^t\|_F^2 &= \langle \partial f(W^t) H_A^t, \partial f(W^t) H_A^t \rangle_F \\ &= \text{Tr} \left( (H_A^t)^\top (\partial f(W^t))^\top \partial f(W^t) H_A^t \right) \\ &= \text{Tr} \left( (\nabla f(W^t))^\top \nabla f(W^t) H_A^t (H_A^t)^\top \right) \\ &= \text{Tr} \left( (\partial f(W^t))^\top \partial f(W^t) H_A^t \right) \\ &= \langle \partial f(W^t), \partial f(W^t) H_A^t \rangle_F, \end{aligned} \quad (74)$$

and similarly, one can show that

$$\|H_B^t \partial f(W^t)\|_F^2 = \langle \partial f(W^t), H_B^t \partial f(W^t) \rangle_F. \quad (75)$$

This allows us to express the second moment term as:

$$\begin{aligned} \gamma_t^2 \mathbb{E} \left[ \|G^t\|_F^2 \mid W^t, H^t \right] &\stackrel{(70)}{=} \gamma_t^2 p \|H_B^t \partial f(W^t)\|_F^2 + \gamma_t^2 (1-p) \|\partial f(W^t) H_A^t\|_F^2 \\ &\stackrel{(74), (75)}{=} \gamma_t^2 p \langle \partial f(W^t), H_B^t \partial f(W^t) \rangle_F + \gamma_t^2 (1-p) \langle \partial f(W^t), \partial f(W^t) H_A^t \rangle_F. \end{aligned} \quad (76)$$

Combining these bounds, we can analyze the distance to the optimal solution:

$$\begin{aligned} \mathbb{E} \left[ \|W^{t+1} - W^*\|_F^2 \mid W^t, H^t \right] &= \mathbb{E} \left[ \|W^t - \gamma_t G^t - W^*\|_F^2 \mid W^t, H^t \right] \\ &= \|W^t - W^*\|_F^2 - 2\gamma_t \mathbb{E} \left[ \langle G^t, W^t - W^* \rangle_F \mid W^t, H^t \right] \\ &\quad + \gamma_t^2 \mathbb{E} \left[ \|G^t\|_F^2 \mid W^t, H^t \right] \\ &\stackrel{(73), (76)}{=} \|W^t - W^*\|_F^2 - 2\gamma_t p \langle H_B^t \partial f(W^t), W^t - W^* \rangle_F \\ &\quad - 2\gamma_t (1-p) \langle \partial f(W^t) H_A^t, W^t - W^* \rangle_F + \gamma_t^2 p \langle \partial f(W^t), H_B^t \partial f(W^t) \rangle_F \\ &\quad + \gamma_t^2 (1-p) \langle \partial f(W^t), \partial f(W^t) H_A^t \rangle_F. \end{aligned} \quad (77)$$

For the expected projection matrices (see Assumption 10), we have:

$$\begin{aligned} \langle \partial f(W^t), \mathbb{E} [H_B^t] \partial f(W^t) \rangle_F &= \text{Tr} \left( (\partial f(W^t))^\top \mathbb{E} [H_B^t] \partial f(W^t) \right) \\ &= \alpha \text{Tr} \left( (\partial f(W^t))^\top \partial f(W^t) \right) \\ &= \alpha \|\partial f(W^t)\|_F^2, \end{aligned} \quad (78)$$

and similarly,

$$\langle \partial f(W^t), \partial f(W^t) \mathbb{E} [H_A^t] \rangle_F = \alpha \|\partial f(W^t)\|_F^2. \quad (79)$$

Taking expectation of both sides of (77) again, we get

$$\mathbb{E} \left[ \|W^{t+1} - W^*\|_F^2 \mid W^t \right] = \mathbb{E} \left[ \mathbb{E} \left[ \|W^{t+1} - W^*\|_F^2 \mid W^t, H^t \right] \mid W^t \right] \quad (80)$$

$$\begin{aligned} &= \|W^t - W^*\|_F^2 - 2\gamma_t p \langle \mathbb{E} [H_B^t] \partial f(W^t), W^t - W^* \rangle_F \\ &\quad - 2\gamma_t (1-p) \langle \partial f(W^t) \mathbb{E} [H_A^t], W^t - W^* \rangle_F \\ &\quad + \gamma_t^2 p \langle \partial f(W^t), \mathbb{E} [H_B^t] \partial f(W^t) \rangle_F + \gamma_t^2 (1-p) \langle \partial f(W^t), \partial f(W^t) \mathbb{E} [H_A^t] \rangle_F \end{aligned} \quad (81)$$

$$\begin{aligned} &\stackrel{(78), (79)}{=} \|W^t - W^*\|_F^2 - 2\gamma_t p \alpha \langle \partial f(W^t), W^t - W^* \rangle_F \\ &\quad - 2\gamma_t (1-p) \alpha \langle \partial f(W^t), W^t - W^* \rangle_F + \gamma_t^2 \alpha \|\partial f(W^t)\|_F^2 \end{aligned} \quad (82)$$

$$\begin{aligned} &= \|W^t - W^*\|_F^2 - 2\gamma_t \alpha \langle \partial f(W^t), W^t - W^* \rangle_F + \gamma_t^2 \alpha \|\partial f(W^t)\|_F^2 \\ &\stackrel{(72)}{=} \|W^t - W^*\|_F^2 - 2\gamma_t \alpha (f(W^t) - f(W^*)) + \gamma_t^2 \alpha \|\partial f(W^t)\|_F^2. \end{aligned} \quad (83)$$

By Assumption 9, subgradients are uniformly bounded (see [Beck, 2017]):

$$\|\partial f(W)\|_F \leq L_0 \quad \forall W \in \mathbb{R}^{m \times n}. \quad (84)$$

Now we analyze both stepsize strategies separately.

**1. (Constant stepsize).** Let us first consider using a fixed stepsize  $\gamma_t := \gamma > 0$ . Taking expectation of both sides of (80) again, applying tower property (13) and using the bound (84), we obtain:

$$\mathbb{E} \left[ \|W^{t+1} - W^*\|_F^2 \right] \leq \mathbb{E} \left[ \|W^t - W^*\|_F^2 \right] - 2\gamma \alpha \mathbb{E} [f(W^t) - f(W^*)] + \gamma^2 \alpha L_0^2. \quad (85)$$

Rearranging terms in (85):

$$2\gamma\alpha\mathbb{E}[f(W^t) - f(W^*)] \leq \mathbb{E}[\|W^t - W^*\|_F^2] - \mathbb{E}[\|W^{t+1} - W^*\|_F^2] + \gamma^2\alpha L_0^2. \quad (86)$$

Summing inequality (86) for  $t = 0, \dots, T-1$ :

$$\begin{aligned} 2\gamma\alpha \sum_{t=0}^{T-1} \mathbb{E}[f(W^t) - f(W^*)] &\leq \sum_{t=0}^{T-1} \left( \mathbb{E}[\|W^t - W^*\|_F^2] - \mathbb{E}[\|W^{t+1} - W^*\|_F^2] \right) \\ &\quad + T\gamma^2\alpha L_0^2 \\ &= \mathbb{E}[\|W^0 - W^*\|_F^2] - \mathbb{E}[\|W^T - W^*\|_F^2] + T\gamma^2\alpha L_0^2 \\ &\leq \|W^0 - W^*\|_F^2 + T\gamma^2\alpha L_0^2, \end{aligned} \quad (87)$$

where the last inequality follows from the non-negativity of  $\|W^T - W^*\|_F^2$ .

For the averaged iterate  $\bar{W}^T := \frac{1}{T} \sum_{t=0}^{T-1} W^t$ , by convexity of  $f$  we have:

$$\begin{aligned} \mathbb{E}[f(\bar{W}^T) - f(W^*)] &\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(W^t) - f(W^*)] \\ &\stackrel{(87)}{\leq} \frac{\|W^0 - W^*\|_F^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2} \\ &= \frac{(R^0)^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2}, \end{aligned} \quad (88)$$

where we denoted  $(R^0)^2 := \|W^0 - W^*\|_F^2$ .

To optimize this bound, we minimize it with respect to  $\gamma$ . The optimal stepsize  $\gamma_*$  solves:

$$\begin{aligned} \gamma_* &= \arg \min_{\gamma > 0} \left( \frac{(R^0)^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2} \right) \\ &= \sqrt{\frac{(R^0)^2}{T\alpha L_0^2}}. \end{aligned} \quad (89)$$

Substituting  $\gamma_*$  back into (88), we obtain the optimal convergence rate:

$$\mathbb{E}[f(\bar{W}^T) - f(W^*)] \leq \frac{R^0 L_0}{\sqrt{\alpha T}}. \quad (90)$$

**2. (Polyak stepsize).** For this strategy, we choose the stepsize adaptively based on the current function value:

$$\begin{aligned} \gamma_t &= \arg \min_{\gamma > 0} \left\{ \|W^t - W^*\|_F^2 - 2\gamma\alpha (f(W^t) - f(W^*)) + \gamma^2\alpha \|\partial f(W^t)\|_F^2 \right\} \\ &= \frac{(f(W^t) - f(W^*))}{\|\partial f(W^t)\|_F^2}. \end{aligned} \quad (91)$$

Substituting this stepsize into inequality (80):

$$\begin{aligned} \mathbb{E}[\|W^{t+1} - W^*\|_F^2 | W^t] &= \mathbb{E}[\mathbb{E}[\|W^{t+1} - W^*\|_F^2 | W^t, H^t] | W^t] \\ &\leq \|W^t - W^*\|_F^2 - 2\gamma_t\alpha (f(W^t) - f(W^*)) + \gamma_t^2\alpha \|\partial f(W^t)\|_F^2 \\ &\stackrel{(91)}{=} \|W^t - W^*\|_F^2 - \frac{\alpha (f(W^t) - f(W^*))^2}{\|\partial f(W^t)\|_F^2} \\ &\stackrel{(84)}{\leq} \|W^t - W^*\|_F^2 - \frac{\alpha (f(W^t) - f(W^*))^2}{L_0^2}. \end{aligned} \quad (92)$$

Taking expectation of both sides of (92) again and applying the tower property

$$\mathbb{E} \left[ \|W^{t+1} - W^*\|_F^2 \right] \leq \mathbb{E} \left[ \|W^t - W^*\|_F^2 \right] - \frac{\alpha \mathbb{E} \left[ (f(W^t) - f(W^*))^2 \right]}{L_0^2} \quad (93)$$

Since  $f$  is convex, by Jensen's inequality (16) and the Cauchy-Bunyakovsky-Schwarz inequality (14) with  $X := f(W^t) - f(W^*)$  and  $Y := 1$ , we have

$$\begin{aligned} \mathbb{E} \left[ f_i(\overline{W}^T) - f(W^*) \right] &\stackrel{(16)}{\leq} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} f(W^t) - f(W^*) \right] \\ &\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [f(W^t) - f(W^*)] \\ &\stackrel{(14)}{\leq} \frac{1}{T} \sum_{t=0}^{T-1} \sqrt{\mathbb{E} [(f(W^t) - f(W^*))^2]} \\ &\leq \sqrt{\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [(f(W^t) - f(W^*))^2]} \\ &\stackrel{(93)}{\leq} \frac{R^0 L_0}{\sqrt{\alpha T}}, \end{aligned} \quad (94)$$

which matches the optimal rate achieved by the constant stepsize strategy with optimal tuning.  $\square$

## C.2 Analysis of Bernoulli-LoRA-SGD

---

### Algorithm 4 Bernoulli-LoRA-SGD

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling factor  $\alpha > 0$ , chain
   length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top g(W^t)$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta g(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13: end for

```

---

Earlier findings were derived utilizing full gradient computations. Nonetheless, this method proves impractical in deep learning applications, where obtaining full gradients is rarely feasible. Our focus moves to a framework that employs **Stochastic Gradient Descent (SGD)** while incorporating a more flexible and generalized data sampling strategy, enabling greater adaptability in the selection and utilization of data throughout the training process. General sampling techniques for strongly convex functions have been thoroughly examined in [Gower et al., 2019]. For broader convex optimization problems, Khaled et al. [2023] provide a comprehensive study of how **SGD** performs under different sampling strategies. In non-convex scenarios, the works of Khaled and Richtárik [2023] and [Demidovich et al., 2023b] investigate the effects of generalized sampling methods on **SGD**'s convergence and efficiency, offering valuable insights into its adaptability for diverse machine learning applications. In this section we focus on **Bernoulli-LoRA-SGD**, a method, designed in the scope of **Bernoulli-LoRA** framework, based on the classical **SGD** algorithm.

For convergence analysis, we notice the gradient step in Algorithm 4 is equivalent to the following update

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}, \quad (95)$$

where  $G^t = g(W^t)$  is an unbiased stochastic gradient, which satisfies Assumption 4.

### C.2.1 Convergence for Smooth Non-Convex Functions

**Theorem 11.** *Let Assumptions 2, 3, and 4 hold, and stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{1}{\sqrt{LA_1 \lambda_{\max}^p T}}, \frac{1}{LB_1} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}.$$

*Then iterates generated by Bernoulli-LoRA-SGD (Algorithm 4) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\tilde{W}^T) \right\|_F^2 \right] \leq \frac{6(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T} + \gamma LC_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$



where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\widetilde{W}^T$  is chosen at random from  $\{W^0, W^1, \dots, W^{T-1}\}$  with probabilities  $\{\frac{w_t}{\mathcal{W}_{T-1}}\}_{t=0}^{T-1}$ , where  $w_t = \frac{w_{t-1}}{(1+\gamma^2 L A_1 \lambda_{\max}^p)}$ ,  $\mathcal{W}_{T-1} = \sum_{t=0}^{T-1} w_t$ , and  $w^{-1} > 0$ .

*Proof.* We start with smoothness of function  $f$ :

$$\begin{aligned} f(W^{t+1}) &\leq f(W^t) + \langle \nabla f(W^t), W^{t+1} - W^t \rangle + \frac{L}{2} \|W^{t+1} - W^t\|_F^2 \\ &\stackrel{(95)}{=} f(W^t) - \gamma \langle \nabla f(W^t), \hat{G}^t \rangle + \frac{\gamma^2 L}{2} \|\hat{G}^t\|_F^2. \end{aligned} \quad (96)$$

Taking a conditional expectation by  $W^t$ , we bound the second and the third terms from inequality (96):

$$\begin{aligned} \mathbb{E} [\langle \nabla f(W^t), \hat{G}^t \rangle | W^t] &= \langle \nabla f(W^t), \mathbb{E} [\hat{G}^t | W^t] \rangle \\ &\stackrel{(95)}{=} p \langle \nabla f(W^t), \mathbb{E} [H_B^t G^t | W^t] \rangle + (1-p) \langle \nabla f(W^t), \mathbb{E} [G^t H_A^t | W^t] \rangle \\ &\stackrel{(*)}{=} p \langle \nabla f(W^t), \mathbb{E} [H_B^t | W^t] \mathbb{E} [G^t | W^t] \rangle + (1-p) \langle \nabla f(W^t), \mathbb{E} [G^t | W^t] \mathbb{E} [H_A^t | W^t] \rangle \\ &= p \langle \nabla f(W^t), \mathbb{E} [H_B^t | W^t] \nabla f(W^t) \rangle + (1-p) \langle \nabla f(W^t), \nabla f(W^t) \mathbb{E} [H_A^t | W^t] \rangle \\ &\geq \underbrace{(p\lambda_{\min}(\mathbb{E} [H_B^t]) + (1-p)\lambda_{\min}(\mathbb{E} [H_A^t]))}_{:=\lambda_{\min}^p} \|\nabla f(W^t)\|_F^2 \\ &= \lambda_{\min}^p \|\nabla f(W^t)\|_F^2, \end{aligned} \quad (97)$$

where in  $(*)$  we used that  $H_B^t$ ,  $H_A^t$  and  $G^t$  are independent. Now we bound the third term:

$$\begin{aligned} \mathbb{E} [\|\hat{G}^t\|_F^2 | W^t] &\stackrel{(95)}{=} p \mathbb{E} [\|H_B^t G^t\|_F^2 | W^t] + (1-p) \mathbb{E} [\|G^t H_A^t\|_F^2 | W^t] \\ &= p \mathbb{E} [\langle H_B^t G^t, H_B^t G^t \rangle | W^t] + (1-p) \mathbb{E} [\langle G^t H_A^t, G^t H_A^t \rangle | W^t] \\ &\stackrel{(**)}{=} p \mathbb{E} [\langle G^t, H_B^t G^t \rangle | W^t] + (1-p) \mathbb{E} [\langle G^t, G^t H_A^t \rangle | W^t], \end{aligned}$$

where in  $(**)$  we used property of projection matrices  $H_B^t, H_A^t$ . By the independence of  $H_B^t, H_A^t, G^t$ , we obtain

$$\begin{aligned} \mathbb{E} [\|\hat{G}^t\|_F^2 | W^t] &= p \mathbb{E} [\langle G^t, \mathbb{E} [H_B^t | W^t] G^t \rangle | W^t] + (1-p) \mathbb{E} [\langle G^t, G^t \mathbb{E} [H_A^t | W^t] \rangle | W^t] \\ &\leq p \lambda_{\max}(\mathbb{E} [H_B^t | W^t]) \mathbb{E} [\|G^t\|_F^2 | W^t] + (1-p) \lambda_{\max}(\mathbb{E} [H_A^t | W^t]) \mathbb{E} [\|G^t\|_F^2 | W^t] \\ &= \underbrace{(p\lambda_{\max}(\mathbb{E} [H_B^t | W^t]) + (1-p)\lambda_{\max}(\mathbb{E} [H_A^t | W^t]))}_{:=\lambda_{\max}^p} \mathbb{E} [\|G^t\|_F^2 | W^t] \\ &= \lambda_{\max}^p \mathbb{E} [\|G^t\|_F^2 | W^t]. \end{aligned} \quad (98)$$

Plugging (97) and (98) into (96), we obtain

$$\begin{aligned} \mathbb{E} [f(W^{t+1}) | W^t] &\leq f(W^t) - \gamma \mathbb{E} [\langle \nabla f(W^t), \hat{G}^t \rangle | W^t] + \frac{\gamma^2 L}{2} \mathbb{E} [\|\hat{G}^t\|_F^2 | W^t] \\ &\leq f(W^t) - \gamma \lambda_{\min}^p \|\nabla f(W^t)\|_F^2 + \frac{\gamma^2 \lambda_{\max}^p L}{2} \mathbb{E} [\|G^t\|_F^2 | W^t]. \end{aligned}$$

By Assumption 4,

$$\begin{aligned}
\mathbb{E} [f(W^{t+1}) - f^* | W^t] &\leq f(W^t) - \gamma \mathbb{E} [\langle \nabla f(W^t), \hat{G}^t \rangle | W^t] + \frac{\gamma^2 L}{2} \mathbb{E} [\|\hat{G}^t\|_F^2 | W^t] \\
&\leq f(W^t) - f^* - \gamma \lambda_{\min}^p \|\nabla f(W^t)\|_F^2 \\
&\quad + \frac{\gamma^2 \lambda_{\max}^p L}{2} \left( 2A_1(f(W^t) - f^*) + B_1 \|\nabla f(W^t)\|_F^2 + C_1 \right) \\
&\leq (1 + \gamma^2 \lambda_{\max}^p L A_1) (f(W^t) - f^*) - \gamma \lambda_{\min}^p \left( 1 - \frac{\gamma L B_1 \lambda_{\max}^p}{2 \lambda_{\min}^p} \right) \|\nabla f(W^t)\|_F^2 \\
&\quad + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}.
\end{aligned}$$

Taking mathematical expectation and selecting a stepsize as  $0 < \gamma \leq \frac{1}{LB_1} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1}$ , we get

$$\begin{aligned}
\mathbb{E} [f(W^{t+1}) - f^*] &\leq (1 + \gamma^2 \lambda_{\max}^p L A_1) \mathbb{E} [f(W^t) - f^*] \\
&\quad - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}.
\end{aligned} \tag{99}$$

Defining  $\delta^t := \mathbb{E} [f(W^t) - f^*]$ ,  $r^t := \mathbb{E} [\|\nabla f(W^t)\|_F^2]$  for every  $t \geq 0$ , we have

$$\delta^{t+1} \leq (1 + \gamma^2 \lambda_{\max}^p L A_1) \delta^t - \frac{\gamma \lambda_{\min}^p}{2} r^t + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}.$$

Fixing  $w^{-1} > 0$  and defining  $w_t = \frac{w_{t-1}}{1 + \gamma^2 L A_1 \lambda_{\max}^p}$  for all  $t \geq 0$ , we have

$$\begin{aligned}
\frac{1}{2} \lambda_{\min}^p w_t r^t &\leq \frac{w_t}{\gamma} (1 + \gamma^2 \lambda_{\max}^p L A_1) \delta^t - \frac{w_t}{\gamma} \delta^{t+1} + \frac{1}{2} \gamma L C_1 \lambda_{\max}^p w_t \\
&= \frac{w_{t-1} \delta^t}{\gamma} - \frac{w_t \delta^{t+1}}{\gamma} + \frac{1}{2} \gamma L C_1 \lambda_{\max}^p w_t.
\end{aligned}$$

Summing over  $t$  from 0 to  $T - 1$ , we have

$$\sum_{t=0}^{T-1} w_t r^t \leq \frac{2w_{-1} \delta^0}{\gamma \lambda_{\min}^p} - \frac{2w_{T-1} \delta^T}{\gamma \lambda_{\min}^p} + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \sum_{t=0}^{T-1} w_t.$$

Defining  $\mathcal{W}_{T-1} = \sum_{t=0}^{T-1} w_t$ , we acquire

$$\sum_{t=0}^{T-1} \frac{w_t}{\mathcal{W}_{T-1}} r^t \leq \frac{2w_{-1} \delta^0}{\gamma \lambda_{\min}^p \mathcal{W}_{T-1}} + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p}.$$

Using the next chain of inequalities

$$W_{T-1} = \sum_{t=0}^{T-1} w_t \geq T \min_{0 \leq t \leq T-1} w_t = T w_{T-1} = \frac{T w_{-1}}{(1 + \gamma^2 \lambda_{\max}^p L A_1)^T},$$

we have

$$\sum_{t=0}^{T-1} \frac{w_t}{\mathcal{W}_{T-1}} r^t \leq \frac{2(1 + \gamma^2 \lambda_{\max}^p L A_1)^T}{\gamma T \lambda_{\min}^p} (f(W^0) - f^*) + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p}.$$

Selecting  $0 < \gamma \leq \frac{1}{\sqrt{LA_1 \lambda_{\max}^p T}}$ , and using  $(1 + \gamma^2 \lambda_{\max}^p LA_1)^T \leq \exp(\gamma^2 \lambda_{\max}^p LA_1 T) \leq \exp(1) \leq 3$ , we obtain

$$\sum_{t=0}^{T-1} \frac{w_t}{\mathcal{W}^{T-1}} r^t \leq \frac{6\delta^0}{\gamma T \lambda_{\min}^p} + \gamma LC_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p}.$$

□

Next we show convergence of **Bernoulli-LoRA-SGD** under additional Assumption 6.

### C.2.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 12.** *Let Assumptions 2, 3, 4, and 6 hold, and stepsize satisfy*

$0 < \gamma \leq \min \left\{ \frac{\mu \lambda_{\min}^p}{2LA_1 \lambda_{\max}^p}, \frac{2}{\mu \lambda_{\min}^p}, \frac{1}{LB_1} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}$ . Then iterates generated by **Bernoulli-LoRA-SGD** (Algorithm 4) satisfy

$$\mathbb{E} [f(W^T) - f^*] \leq \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^T (f(W^0) - f^*) + \frac{\gamma LC_1}{\mu} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ .

*Proof.* We start our proof with inequality 99. Using PL-inequality (see Assumption 6), we have

$$\begin{aligned} \mathbb{E} [f(W^{t+1}) - f^*] &\leq (1 + \gamma^2 \lambda_{\max}^p LA_1) \mathbb{E} [f(W^t) - f^*] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{2} \\ &\leq (1 - \gamma \mu \lambda_{\min}^p + \gamma^2 \lambda_{\max}^p LA_1) \mathbb{E} [f(W^t) - f^*] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{2}. \end{aligned}$$

Taking the stepsize as  $0 < \gamma \leq \min \left\{ \frac{\mu \lambda_{\min}^p}{2LA_1 \lambda_{\max}^p}, \frac{2}{\mu \lambda_{\min}^p} \right\}$ , we obtain

$$\begin{aligned} \mathbb{E} [f(W^{t+1}) - f^*] &\leq \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right) \mathbb{E} [f(W^t) - f^*] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{2} \\ &\leq \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^{t+1} \mathbb{E} [f(W^0) - f^*] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{2} \sum_{\tau=0}^t \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^{t-\tau} \\ &\leq \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^{t+1} \mathbb{E} [f(W^0) - f^*] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{2} \sum_{\tau=0}^{\infty} \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^{\tau} \\ &= \left( 1 - \frac{1}{2} \gamma \mu \lambda_{\min}^p \right)^{t+1} \mathbb{E} [f(W^0) - f^*] + \frac{\gamma^2 \lambda_{\max}^p LC_1}{\gamma \mu \lambda_{\min}^p}, \end{aligned}$$

where in the last equation we use the formula of the sum of geometric progression. □

### C.3 Analysis of Bernoulli-LoRA-MVR

---

**Algorithm 5** Bernoulli-LoRA-MVR

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ ,  $G^0 \in \mathbb{R}^{m \times n}$  rank  $r \ll \min\{m, n\}$ , scaling factor
    $\alpha > 0$ , chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ , momentum parameter
    $b \in [0, 1]$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta G^t (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13:   Sample  $\xi^{t+1} \sim \mathcal{D}$ 
14:    $G^{t+1} = \nabla f_{\xi^{t+1}}(W^{t+1}) + (1 - b) (G^t - \nabla f_{\xi^{t+1}}(W^t))$ 
15: end for

```

---

Recently, there has been a significant surge of interest in variance-reduced methods for addressing finite-sum problems [J Reddi et al., 2015, Shang et al., 2018, Malinovsky et al., 2022, Richtárik et al., 2024]. It has gained prominence as a formidable alternative to stochastic gradient descent (SGD) in tackling non-convex optimization problems. Notably, it has been pivotal in introducing the first algorithms capable of surpassing SGD’s convergence rate for locating first-order critical points. Despite these advancements, variance reduction methods often come with challenges, including the necessity for meticulously tuned learning rates and the reliance on overly large batch sizes to realize their benefits. To address some of these limitations, Momentum Variance Reduction (MVR) was proposed specifically for server-only stochastic non-convex optimization [Cutkosky and Orabona, 2019]. This approach leverages a modified form of momentum to achieve variance reduction while eliminating the dependence on large batch sizes. A proof on MVR technique with better dependence on momentum parameter was obtained by Tyurin and Richtárik [2023]. In the context of Federated Learning, Karagulyan et al. [2024] proposed the SPAM method. On the server side, MVR is utilized to enhance optimization efficiency, while the client side incorporates the Stochastic Proximal Point Method updates. This section is devoted to Bernoulli-LoRA-MVR, a method, designed in the scope of Bernoulli-LoRA framework, based on the MVR technique.

To show convergence guarantees for Bernoulli-LoRA-MVR, the iterates of the method can be rewritten in following way

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases} \quad (100)$$

$$G^{t+1} = \nabla f_{\xi^{t+1}}(W^{t+1}) + (1 - b) (G^t - \nabla f_{\xi^{t+1}}(W^t)). \quad (101)$$

First of all, we reprove descent lemma from the paper of Li et al. [2021] for generic gradient step (100).

**Lemma 5.** Let Assumptions 1, 3 hold. Then, iterates defined as (100) satisfy

$$\begin{aligned} \mathbb{E} [f(W^{t+1}) - f^* \mid W^t] &\leq f(W^t) - f^* - \frac{\gamma \lambda_{\min}^p}{2} \|\nabla f(W^t)\|_F^2 \\ &\quad + \frac{\gamma \lambda_{\max}^p}{2} \|G^t - \nabla f(W^t)\|_F^2 - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbb{E} [\|W^{t+1} - W^t\|_F^2 \mid W^t]. \end{aligned}$$

*Proof.* By Assumption 3, we have

$$\begin{aligned} f(W^{t+1}) &\leq f(W^t) + \langle \nabla f(W^t), W^{t+1} - W^t \rangle_F + \frac{L}{2} \|W^{t+1} - W^t\|_F^2 \\ &= f(W^t) - \gamma \langle \nabla f(W^t), \hat{G}^t \rangle_F + \frac{L}{2} \|W^{t+1} - W^t\|_F^2. \end{aligned} \quad (102)$$

To continue our proof, we need to bound the second term from (102). Taking conditional expectation by  $H^t, W^t$ , we obtain

$$\begin{aligned} \mathbb{E} [\langle \nabla f(W^t), \hat{G}^t \rangle_F \mid H^t, W^t] &\stackrel{(100)}{=} p \langle \nabla f(W^t), H_B^t G^t \rangle_F + (1-p) \langle \nabla f(W^t), G^t H_A^t \rangle_F \\ &= p \langle H_B^t \nabla f(W^t), H_B^t G^t \rangle_F + (1-p) \langle \nabla f(W^t) H_A^t, G^t H_A^t \rangle_F \\ &= \frac{p}{2} \left( \|H_B^t \nabla f(W^t)\|_F^2 + \|H_B^t G^t\|_F^2 - \|H_B^t G^t - H_B^t \nabla f(W^t)\|_F^2 \right) \\ &\quad + \frac{1-p}{2} \left( \|\nabla f(W^t) H_A^t\|_F^2 + \|G^t H_A^t\|_F^2 - \|G^t H_A^t - \nabla f(W^t) H_A^t\|_F^2 \right) \\ &\geq \frac{1}{2} \left( p \|H_B^t \nabla f(W^t)\|_F^2 + (1-p) \|\nabla f(W^t) H_A^t\|_F^2 \right) + \frac{1}{2} \mathbb{E} [\|\hat{G}^t\|_F^2 \mid H^t, W^t] \\ &\quad - \frac{1}{2} \left( p \|H_B^t G^t - H_B^t \nabla f(W^t)\|_F^2 + (1-p) \|G^t H_A^t - \nabla f(W^t) H_A^t\|_F^2 \right). \end{aligned}$$

Taking conditional expectation by  $W^t$ , we have

$$\begin{aligned} \mathbb{E} [\langle \nabla f(W^t), \hat{G}^t \rangle_F \mid W^t] &\geq \frac{1}{2} \left( p \mathbb{E} [\|H_B^t \nabla f(W^t)\|_F^2 \mid W^t] + (1-p) \mathbb{E} [\|\nabla f(W^t) H_A^t\|_F^2 \mid W^t] \right) + \frac{1}{2} \mathbb{E} [\|\hat{G}^t\|_F^2 \mid W^t] \\ &\quad - \frac{1}{2} \left( p \mathbb{E} [\|H_B^t G^t - H_B^t \nabla f(W^t)\|_F^2 \mid W^t] + (1-p) \mathbb{E} [\|G^t H_A^t - \nabla f(W^t) H_A^t\|_F^2 \mid W^t] \right) \\ &\stackrel{(*)}{\geq} \frac{1}{2} \underbrace{(p \lambda_{\min}(\mathbb{E}[H_B^t]) + (1-p) \lambda_{\min}(\mathbb{E}[H_A^t]))}_{:=\lambda_{\min}^p} \|\nabla f(W^t)\|_F^2 + \frac{1}{2} \mathbb{E} [\|\hat{G}^t\|_F^2 \mid W^t] \\ &\quad - \frac{1}{2} \underbrace{(p \lambda_{\max}(\mathbb{E}[H_B^t]) + (1-p) \lambda_{\max}(\mathbb{E}[H_A^t]))}_{:=\lambda_{\max}^p} \|G^t - \nabla f(W^t)\|_F^2 \\ &\stackrel{(100)}{=} \frac{\lambda_{\min}^p}{2} \|\nabla f(W^t)\|_F^2 + \frac{1}{2\gamma^2} \mathbb{E} [\|W^{t+1} - W^t\|_F^2 \mid W^t] - \frac{\lambda_{\max}^p}{2} \|G^t - \nabla f(W^t)\|_F^2, \end{aligned} \quad (103)$$

where in  $(*)$  we used the following inequalities for any matrix  $V \in \mathbb{R}^{m \times n}$

$$\begin{aligned} \mathbb{E} [\|H_B^t V\|_F^2] &= \mathbb{E} [\langle H_B^t V, H_B^t V \rangle_F] = \langle \mathbb{E}[H_B^t] V, V \rangle_F \geq \lambda_{\min}(\mathbb{E}[H_B^t]) \|V\|_F^2, \\ \mathbb{E} [\|H_B^t V\|_F^2] &\leq \lambda_{\max}(\mathbb{E}[H_B^t]) \|V\|_F^2, \\ \mathbb{E} [\|V H_A^t\|_F^2] &= \mathbb{E} [\langle V H_A^t, V H_A^t \rangle_F] = \langle V \mathbb{E}[H_A^t], V \rangle_F \geq \lambda_{\min}(\mathbb{E}[H_A^t]) \|V\|_F^2, \\ \mathbb{E} [\|V H_A^t\|_F^2] &\leq \lambda_{\max}(\mathbb{E}[H_A^t]) \|V\|_F^2. \end{aligned}$$

Plugging in (103) into (102), we get

$$\begin{aligned}\mathbb{E} [f(W^{t+1}) | W^t] &\leq f(W^t) - \frac{\gamma\lambda_{\min}^p}{2} \|\nabla f(W^t)\|_F^2 - \frac{1}{2\gamma} \mathbb{E} [\|W^{t+1} - W^t\|_F^2 | W^t] \\ &\quad + \frac{\gamma\lambda_{\max}^p}{2} \|G^t - \nabla f(W^t)\|_F^2 + \frac{L}{2} \mathbb{E} [\|W^{t+1} - W^t\|_F^2 | W^t].\end{aligned}$$

□

**Lemma 6.** *Let Assumptions 3, 5 hold. Then, iterates generated by Bernoulli-LoRA-MVR (Algorithm 5) satisfy*

$$\mathbb{E} [\|G^{t+1} - \nabla f(W^{t+1})\|_F^2] \leq (1-b)^2 \mathbb{E} [\|G^t - \nabla f(W^t)\|_F^2] + 2(1-b)^2 L^2 \mathbb{E} [\|W^{t+1} - W^t\|_F^2] + 2b^2 \sigma^2 \quad (104)$$

*Proof.* Taking conditional expectation by  $\mathcal{F}^{t+1} = \{W^{t+1}, G^t\}$ , we obtain

$$\begin{aligned}\mathbb{E} [\|G^{t+1} - \nabla f(W^{t+1})\|_F^2 | \mathcal{F}^{t+1}] &\stackrel{(101)}{=} \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1}) + (1-b)(G^t - \nabla f_{\xi^{t+1}}(W^t))\|_F^2 | \mathcal{F}^{t+1}] \\ &\stackrel{(15)}{=} (1-b)^2 \|G^t - \nabla f(W^t)\|_F^2 \\ &\quad + \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1}) + (1-b)(\nabla f(W^t) - \nabla f_{\xi^{t+1}}(W^t))\|_F^2 | \mathcal{F}^{t+1}] \\ &\leq (1-b)^2 \|G^t - \nabla f(W^t)\|_F^2 + 2b^2 \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1})\|_F^2 | \mathcal{F}^{t+1}] \\ &\quad + 2(1-b)^2 \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f_{\xi^{t+1}}(W^t) - \nabla f(W^{t+1}) + \nabla f(W^t)\|_F^2 | \mathcal{F}^{t+1}] \\ &\leq (1-b)^2 \|G^t - \nabla f(W^t)\|_F^2 + 2b^2 \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1})\|_F^2 | \mathcal{F}^{t+1}] \\ &\quad + 2(1-b)^2 \mathbb{E} [\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f_{\xi^{t+1}}(W^t)\|_F^2 | \mathcal{F}^{t+1}] \\ &\leq (1-b)^2 \|G^t - \nabla f(W^t)\|_F^2 + 2(1-b)^2 L^2 \|W^{t+1} - W^t\|_F^2 + 2b^2 \sigma^2,\end{aligned}$$

where in the last inequality we used smoothness of  $f_{\xi}$  and bounded variance assumption. Taking math expectation, we conclude the proof. □

### C.3.1 Convergence for Smooth Non-Convex Functions

**Theorem 13.** *Let Assumptions 1, 2, 3, and 5 hold, and let the stepsize satisfy  $0 < \gamma \leq \frac{1}{L(1 + \sqrt{\frac{2\lambda_{\max}^p(1-b)^2}{b}})}$ .*

*Then the iterates of Bernoulli-LoRA-MVR (Algorithm 5) satisfy*

$$\mathbb{E} [\|\nabla f(\widetilde{W}^T)\|_F^2] \leq \frac{2(f(W^0) - f^*)}{\lambda_{\min}^p \gamma T} + \frac{\|G^0 - \nabla f(W^0)\|_F^2}{b(2-b)T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p} + \frac{2b\sigma^2}{2-b} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p}, \quad (105)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ ,  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma\lambda_{\max}^p}{2b(2-b)} \|G^t - \nabla f(W^t)\|_F^2. \quad (106)$$

By Lemma 5 and Lemma 6, we have

$$\begin{aligned}
\mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2}\mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\
&\quad + \frac{\gamma\lambda_{\max}^p}{2}\mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma(1-b)^2\lambda_{\max}^p}{2b(2-b)}\mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\
&\quad + \frac{\gamma(1-b)^2L^2\lambda_{\max}^p}{2b(2-b)}\mathbb{E}[\|W^{t+1} - W^t\|_F^2] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \\
&\leq \mathbb{E}[\Phi_t] - \frac{\gamma\lambda_{\min}^p}{2}\mathbb{E}[\|\nabla f(W^t)\|_F^2] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \\
&\quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-b)^2L^2\lambda_{\max}^p}{2b(2-b)}\right)\mathbb{E}[\|W^{t+1} - W^t\|_F^2].
\end{aligned}$$

Selecting  $0 < \gamma \leq \frac{1}{L\left(1 + \sqrt{\frac{(1-b)^2}{b(2-b)}\lambda_{\max}^p}\right)}$ , we obtain

$$\mathbb{E}[\Phi_{t+1}] \leq \mathbb{E}[\Phi_t] - \frac{\gamma\lambda_{\min}^p}{2}\mathbb{E}[\|\nabla f(W^t)\|_F^2] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b}.$$

Summing over  $t$  from 0 to  $T-1$ , we get

$$\frac{\gamma\lambda_{\min}^p}{2} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W^t)\|_F^2] \leq \mathbb{E}[\Phi_0] - \mathbb{E}[\Phi_T] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b}T.$$

Finally, dividing both sides by  $\frac{\gamma\lambda_{\min}^p}{2}$  yields

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_F^2\right] \leq \frac{2\Phi_0}{\lambda_{\min}^p\gamma T} + \frac{2b\sigma^2}{2-b} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .  $\square$

Next we show convergence guarantee for **Bernoulli-LoRA-MVR**, supposing additionally Assumption 6 holds.

### C.3.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 14.** *Let Assumptions 1, 2, 3, 5, and 6 hold, and let the stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1-b)^2}{b(2-b)}\lambda_{\max}^p} \right)}, \frac{b}{2\mu\lambda_{\min}^p} \right\}.$$

*Then the iterates of Bernoulli-LoRA-MVR (Algorithm 5) satisfy*

$$\mathbb{E}[f(W^T) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0 + \frac{b\sigma^2}{(2-b)\mu} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p}, \quad (107)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\Phi_0 = f(W^0) - f^* + \frac{\gamma\lambda_{\max}^p}{b(2-b)}\|G^0 - \nabla f(W^0)\|_F^2$ .



*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma\lambda_{\max}^p}{b(2-b)} \|G^t - \nabla f(W^t)\|_F^2. \quad (108)$$

By Lemma 5 and Lemma 6, we have

$$\begin{aligned} \mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\ &\quad + \frac{\gamma\lambda_{\max}^p}{2} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma(1-b)^2\lambda_{\max}^p}{b(2-b)} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad + \frac{\gamma(1-b)^2L^2\lambda_{\max}^p}{b(2-b)} \mathbb{E}[\|W^{t+1} - W^t\|_F^2] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \\ &\leq \max\left\{1 - \gamma\mu\lambda_{\min}^p, 1 - \frac{b}{2}\right\} \mathbb{E}[\Phi_t] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \\ &\quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-b)^2L^2\lambda_{\max}^p}{b(2-b)}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2], \end{aligned}$$

where in the last inequality we used Assumption 6. Selecting positive stepsize  $\gamma$  satisfying the upper bound assumed in the theorem statement, we obtain

$$\begin{aligned} \mathbb{E}[\Phi_{t+1}] &\leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[\Phi_t] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \\ &\leq (1 - \gamma\mu\lambda_{\min}^p)^{t+1} \mathbb{E}[\Phi_0] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \sum_{\tau=0}^t (1 - \gamma\mu\lambda_{\min}^p)^{t-\tau} \\ &\leq (1 - \gamma\mu\lambda_{\min}^p)^{t+1} \mathbb{E}[\Phi_0] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{2-b} \sum_{\tau=0}^{\infty} (1 - \gamma\mu\lambda_{\min}^p)^{\tau} \\ &= (1 - \gamma\mu\lambda_{\min}^p)^{t+1} \mathbb{E}[\Phi_0] + \frac{\gamma\lambda_{\max}^pb\sigma^2}{(2-b)\gamma\mu\lambda_{\min}^p}, \end{aligned}$$

where, in the last equation, we used the formula for the sum of a geometric progression. □

## C.4 Analysis of Bernoulli-LoRA-PAGE

---

### Algorithm 6 Bernoulli-LoRA-PAGE

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , a vector  $G^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling
   factor  $\alpha > 0$ , chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ , probability  $q$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta g(W^t) (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger A_S^t$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13:   Sample  $i_{t+1}$  uniformly at random from  $[n]$ 
14:    $G^{t+1} = \begin{cases} \nabla f(W^{t+1}), & \text{with probability } q \\ G^t + (\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t)), & \text{with probability } 1 - q \end{cases}$ 
15: end for

```

---

There exist several optimal methods for solving a general non-convex optimization problem, e.g. SPIDER [Fang et al., 2018] and SARAH [Pham et al., 2020]. However, the known lower bound used to establish their optimality works only in the small data regime. Probabilistic Gradient Estimator (PAGE) [Li et al., 2021] is a very simple and easy to implement algorithm, known for achieving optimal convergence results in non-convex optimization. PAGE uses the full gradient update with probability  $q_t$ , or reuses the previous gradient with a small adjustment (at a low computational cost) with probability  $1 - q_t$ . A general version of PAGE on Riemannian manifolds is considered in [Demidovich et al., 2024a]. In this section we present Bernoulli-LoRA-PAGE, a new method within Bernoulli-LoRA framework, based on PAGE algorithm.

Notice, that the iterates of Bernoulli-LoRA-PAGE (Algorithm 6) can be rewritten in the following simple way

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases} \quad (109)$$

$$G^{t+1} = \begin{cases} \nabla f(W^{t+1}), & \text{with probability } q \\ G^t + (\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t)), & \text{with probability } 1 - q \end{cases} \quad (110)$$

**Lemma 7.** *Let Assumption 3 hold. Then, iterates generated by Bernoulli-LoRA-PAGE*

$$\mathbb{E} \left[ \|G^{t+1} - \nabla f(W^{t+1})\|_F^2 \right] \leq (1-q) \mathbb{E} \left[ \|G^t - \nabla f(W^t)\|_F^2 \right] + (1-q) L^2 \mathbb{E} \left[ \|W^{t+1} - W^t\|_F^2 \right]. \quad (111)$$

*Proof.* Taking the full mathematical expectation, we obtain

$$\begin{aligned}
\mathbb{E} \left[ \|G^{t+1} - \nabla f(W^{t+1})\|_F^2 \right] &\stackrel{(110)}{=} (1-q)\mathbb{E} \left[ \|G^t - \nabla f(W^{t+1}) + (\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t))\|_F^2 \right] \\
&\stackrel{(15)}{=} (1-q)\mathbb{E} \left[ \|G^t - \nabla f(W^t)\|_F^2 \right] \\
&\quad + (1-q)\mathbb{E} \left[ \|(\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t)) - (\nabla f(W^{t+1}) - \nabla f(W^t))\|_F^2 \right] \\
&\leq (1-q)\mathbb{E} \left[ \|G^t - \nabla f(W^t)\|_F^2 \right] \\
&\quad + (1-q)\mathbb{E} \left[ \|\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t)\|_F^2 \right] \\
&\leq (1-q)\mathbb{E} \left[ \|G^t - \nabla f(W^t)\|_F^2 \right] + (1-q)L^2\mathbb{E} \left[ \|W^{t+1} - W^t\|_F^2 \right],
\end{aligned}$$

where in the last inequality we used smoothness of each  $f_i$ .  $\square$

#### C.4.1 Convergence for Smooth Non-Convex Functions

**Theorem 15.** Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy

$$0 < \gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{1-q}{q} \lambda_{\max}^p} \right)}.$$

Then the iterates of **PAGE-Bernoulli-LoRA** (Algorithm 6) satisfy

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2(f(W^0) - f^*)}{\lambda_{\min}^p \gamma T} + q \frac{\|G^0 - \nabla f(W^0)\|_F^2}{T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p}, \quad (112)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ ,  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\max}^p}{2q} \|G^t - \nabla f(W^t)\|_F^2. \quad (113)$$

By Lemma 5 and Lemma 7, we have

$$\begin{aligned}
\mathbb{E} [\Phi_{t+1}] &\leq \mathbb{E} [f(W^t)] - f^* - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] \\
&\quad - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbb{E} [\|W^{t+1} - W^t\|_F^2] + \frac{\gamma \lambda_{\max}^p}{2} \mathbb{E} [\|G^t - \nabla f(W^t)\|_F^2] \\
&\quad + \frac{\gamma \lambda_{\max}^p (1-q)}{2q} \mathbb{E} [\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma \lambda_{\max}^p (1-q) L^2}{2q} \mathbb{E} [\|W^{t+1} - W^t\|_F^2] \\
&\leq \mathbb{E} [\Phi_t] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma (1-q) L^2 \lambda_{\max}^p}{2q} \right) \mathbb{E} [\|W^{t+1} - W^t\|_F^2].
\end{aligned}$$

Selecting  $0 < \gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{1-q}{q} \lambda_{\max}^p} \right)}$ , we obtain

$$\mathbb{E} [\Phi_{t+1}] \leq \mathbb{E} [\Phi_t] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2].$$

Summing over  $t$  from 0 to  $T - 1$ , we get

$$\frac{\gamma\lambda_{\min}^p}{2} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla f(W^t)\|_F^2 \right] \leq \mathbb{E}[\Phi_0] - \mathbb{E}[\Phi_T].$$

Finally, dividing both sides by  $\frac{\gamma\lambda_{\min}^p}{2}$  yields

$$\mathbb{E} \left[ \|\nabla f(\widetilde{W}^T)\|_F^2 \right] \leq \frac{2\Phi_0}{\gamma\lambda_{\min}^p T}.$$

where  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .  $\square$

#### C.4.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 16.** *Let Assumptions 1, 2, 3, and 6 hold, and let the stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + 2\sqrt{\frac{1-q}{q}} \lambda_{\max}^p \right)}, \frac{q}{2\mu\lambda_{\min}^p} \right\}.$$

*Then the iterates of Bernoulli-LoRA-PAGE (Algorithm 6) satisfy*

$$\mathbb{E} [f(W^T) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0, \quad (114)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ , and  $\Phi_0 = f(W^0) - f^* + \frac{\gamma\lambda_{\max}^p}{q} \|G^0 - \nabla f(W^0)\|_F^2$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma\lambda_{\max}^p}{q} \|G^t - \nabla f(W^t)\|_F^2. \quad (115)$$

By Lemma 5 and Lemma 7, we have

$$\begin{aligned} \mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\ &\quad + \frac{\gamma\lambda_{\max}^p}{2} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma(1-q)\lambda_{\max}^p}{q} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad + \frac{\gamma(1-q)L^2\lambda_{\max}^p}{q} \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\ &\leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[f(W^t) - f^*] + \left( 1 - \frac{q}{2} \right) \frac{\gamma\lambda_{\max}^p}{q} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad - \left( \frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-q)L^2\lambda_{\max}^p}{q} \right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2], \end{aligned}$$

where in the last inequality we used Assumption 6. Selecting  $0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + 2\sqrt{\frac{1-q}{q}} \lambda_{\max}^p \right)}, \frac{q}{2\mu\lambda_{\min}^p} \right\}$ , we obtain

$$\mathbb{E}[\Phi_{t+1}] \leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[\Phi_t].$$

Unrolling the recursion, we obtain

$$\mathbb{E}[\Phi_T] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0.$$

$\square$

## D Proofs for Federated Learning Extensions

In recent years, distributed optimization problems and algorithms have become a focal point in the Machine Learning (ML) community. This surge in interest is driven by the need to train modern deep neural networks, which often involve billions of parameters and massive datasets [Brown et al., 2020, Kolesnikov et al., 2020]. To achieve practical training times [Li, 2020], parallelizing computations, such as stochastic gradient evaluations, has emerged as a natural solution, leading to the widespread adoption of distributed training algorithms [Goyal et al., 2017, You et al., 2019, Le Scao et al., 2023]. Additionally, distributed methods are crucial when data is inherently distributed across multiple devices or clients, often accompanied by privacy constraints—a common scenario in Federated Learning (FL) [Konečný et al., 2016, McMahan et al., 2016, Kairouz et al., 2019, Demidovich et al., 2024b, Sadiiev et al., 2024, Yi et al., 2024].

We develop several FL methods within the **Bernoulli-LoRA** framework and provide a convergence analysis for them.

### D.1 Analysis of Fed-Bernoulli-LoRA-QGD

---

#### Algorithm 7 Fed-Bernoulli-LoRA-QGD

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ , rank  $r \ll \min\{m, n\}$ , scaling factor  $\alpha > 0$ , chain
   length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probabilities  $p$  and  $q$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   for any client  $l \in [M]$  in parallel do
4:     Compute gradient  $\nabla f_l(W^{t+1})$  and send compressed version  $G_l^t = \mathcal{Q}_l^t(\nabla f_l(W^{t+1}))$  to the
       server
5:   end for
6:    $G^t = \frac{1}{M} \sum_{l=1}^M G_l^t$ 
7:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
8:   if  $c^t = 1$  then
9:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
10:     $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
12:   else
13:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
14:     $\hat{B}^t = -\eta G^t (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger$ 
15:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
16:   end if
17:   Broadcast  $W^{t+1}$  to each client  $l \in [M]$ 
18: end for

```

---

Parallel implementations of SGD have become a prominent area of study due to their impressive scalability. However, one of the primary challenges in parallelizing SGD lies in the substantial communication overhead required to exchange gradient updates across nodes. To address this, numerous lossy compression techniques have been developed, enabling nodes to transmit quantized gradients instead of full gradients. While these methods often work well in practice, they are not universally reliable and may fail to ensure convergence.

To overcome these limitations, Quantized SGD (QSGD) by [Alistarh et al. \[2017\]](#) introduces a family of compression techniques that provide both theoretical convergence guarantees and strong empirical performance. QSGD offers a flexible mechanism for balancing communication bandwidth and convergence speed. By adjusting the number of bits transmitted per iteration, nodes can reduce bandwidth usage, albeit at the potential cost of increased variance in the gradient estimates. Different variants of QSGD were considered by [Horvóth et al. \[2022\]](#), [Wen et al. \[2017\]](#), [Panferov et al. \[2024\]](#).

We consider the following distributed optimization problem:

$$\min_{W \in \mathbb{R}^{m \times n}} \frac{1}{M} \sum_{l=1}^M f_l(W),$$

where  $M$  represents the number of clients. In Federated Learning, a primary bottleneck is the communication overhead between clients and the central server. A common approach to mitigate this issue is communication compression.

**Definition 3.** A randomized operator  $\mathcal{Q} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$  is called an unbiased compression operator (or compressor) if there exists a constant  $\omega > 0$  such that, for any matrix  $W \in \mathbb{R}^{m \times n}$ , the following conditions hold:

$$\mathbb{E}[\mathcal{Q}(W)] = W, \quad \text{and} \quad \mathbb{E}[\|\mathcal{Q}(W) - W\|_F^2] \leq \omega \|W\|_F^2. \quad (116)$$

To analyze the optimization process, we introduce the following assumption regarding function dissimilarity:

**Assumption 11.** Let  $f^* := \inf_W f(W)$  and  $f_l^* := \inf_W f_l$  for each  $l \in [M]$ . In the non-convex case, the difference at the optimum is defined as:

$$\Delta^* := f^* - \frac{1}{M} \sum_{l=1}^M f_l^* \geq 0. \quad (117)$$

This assumption quantifies the discrepancy between the global optimal function value and the average of the local optimal function values between the clients.

To start convergence analysis, we rewrite the updates for  $W^t$  and  $G^t$  generated by [Fed-Bernoulli-LoRA-QGD](#) (Algorithm 7) as follows

$$G^t = \frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t(\nabla f_l(W^t)); \quad (118)$$

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}. \quad (119)$$

To establish the convergence guarantee for [Fed-Bernoulli-LoRA-QGD](#) (Algorithm 7), we first demonstrate that the gradient estimator  $G^t$  satisfies Assumption 4. Once this is verified, the convergence rate follows directly using the same reasoning as in the proof of Theorem 11.

**Lemma 8.** Let Assumptions 2, 3, and 11 hold. Then,  $G^t$  defined in Algorithm 7 (see (118)) satisfies Assumption 4 with the following constants:

$$A_1 = \frac{L\omega}{M}, \quad B_1 = 1, \quad C_1 = 2 \frac{L\omega\Delta^*}{M}.$$

*Proof.* First, we show  $G^t$  is an unbiased estimator of  $\nabla f(W^t)$ :

$$\mathbb{E}[G^t|W^t] = \frac{1}{M} \sum_{l=1}^M \mathbb{E}[\mathcal{Q}_l^t(\nabla f_l(W^t))|W^t] \stackrel{(116)}{=} \frac{1}{M} \sum_{l=1}^M \nabla f_l(W^t) = \nabla f(W^t).$$

Now we establish that  $G^t$  satisfies Assumption 4. Taking the conditional expectation with respect to  $W^t$ , we have

$$\begin{aligned} \mathbb{E}[\|G^t\|_F^2|W^t] &= \mathbb{E}\left[\left\|\frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t(\nabla f_l(W^t)) - \nabla f(W^t) + \nabla f(W^t)\right\|_F^2|W^t\right] \\ &\stackrel{(15)}{=} \mathbb{E}\left[\left\|\frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t(\nabla f_l(W^t)) - \nabla f(W^t)\right\|_F^2|W^t\right] + \|\nabla f(W^t)\|_F^2 \\ &= \frac{1}{M^2} \sum_{l=1}^M \mathbb{E}[\|\mathcal{Q}_l^t(\nabla f_l(W^t)) - \nabla f_l(W^t)\|_F^2|W^t] + \|\nabla f(W^t)\|_F^2 \\ &\stackrel{(116)}{\leq} \frac{\omega}{M^2} \sum_{l=1}^M \|\nabla f_l(W^t)\|_F^2 + \|\nabla f(W^t)\|_F^2 \\ &\stackrel{(*)}{\leq} \frac{2L\omega}{M^2} \sum_{l=1}^M (f_l(W^t) - f_l^*) + \|\nabla f(W^t)\|_F^2 \\ &= 2\frac{L\omega}{M} (f(W^t) - f^*) + \|\nabla f(W^t)\|_F^2 + 2\frac{L\omega}{M} \underbrace{\left(f^* - \frac{1}{M} \sum_{l=1}^M f_l^*\right)}_{:=\Delta^*}, \end{aligned}$$

where in  $(*)$  we used smoothness of each  $f_l$ . Thus, we have shown that  $G^t$  satisfies Assumption 4 with following constants

$$A_1 = \frac{L\omega}{M}, \quad B_1 = 1, \quad C_1 = 2\frac{L\omega\Delta^*}{M}.$$

□

### D.1.1 Convergence for Smooth Non-Convex Functions

**Theorem 17.** Let Assumptions 1, 2, and 3 hold, and stepsize satisfy

$$0 < \gamma \leq \min \left\{ \frac{1}{L\sqrt{\frac{\omega}{M}\lambda_{\max}^p T}}, \frac{1}{L} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}.$$

Then iterates generated by Fed-Bernoulli-LoRA-QGD (Algorithm 7) satisfy

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{6(f(W^0) - f^*)}{\gamma\lambda_{\min}^p T} + \frac{2\gamma L\omega\Delta^*}{M} \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\widetilde{W}^T$  is chosen at random from  $\{W^0, W^1, \dots, W^{T-1}\}$  with probabilities  $\{\frac{w_t}{\mathcal{W}_{T-1}}\}_{t=0}^{T-1}$ , where  $w_t = \frac{w_{t-1}}{(1+\gamma^2 L^2 \lambda_{\max}^p \omega/M)}$ ,  $\mathcal{W}_{T-1} = \sum_{t=0}^{T-1} w_t$ , and  $w^{-1} > 0$ .

*Proof.* By Lemma 8, and Theorem 11, we directly obtain the statement of the theorem. □



### D.1.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 18.** *Let Assumptions 1, 2, 3, and 6 hold, and stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{\mu}{2L^2\omega/M} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1}, \frac{2}{\mu\lambda_{\min}^p}, \frac{1}{L} \left( \frac{\lambda_{\max}^p}{\lambda_{\min}^p} \right)^{-1} \right\}.$$

*Then iterates generated by Fed-Bernoulli-LoRA-QGD (Algorithm 7) satisfy*

$$\mathbb{E} [f(W^T) - f^*] \leq \left( 1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p \right)^T (f(W^0) - f^*) + \frac{2\gamma L^2}{\mu} \cdot \frac{\omega}{M} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ .

*Proof.* By Lemma 8, and Theorem 12, we directly obtain the statement of the theorem. □

## D.2 Analysis of Fed-Bernoulli-LoRA-MARINA

---

### Algorithm 8 Fed-Bernoulli-LoRA-MARINA

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ ,  $\{G_l^0\}_{l \in [M]} \in \mathbb{R}^{m \times n}$  rank  $r \ll \min\{m, n\}$ , scaling
   factor  $\alpha > 0$ , chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probabilities  $p$  and  $q$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta G^t (A_S^t)^\top \left( A_S^t (A_S^t)^\top \right)^\dagger$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:  end if
13:  Broadcast  $W^{t+1}$  to each client  $l \in [M]$ 
14:  Sample  $s^t \sim \text{Be}(q)$ 
15:  for any client  $l \in [M]$  in parallel do
16:    Compute gradient  $\nabla f_l(W^{t+1})$ 
17:     $G_l^{t+1} = \begin{cases} \nabla f_l(W^{t+1}), & \text{with probability } q \\ G_l^t + \mathcal{Q}_l^t (\nabla f_l(W^{t+1}) - \nabla f_l(W^t)), & \text{with probability } 1 - q \end{cases}$ 
18:    Send  $G_l^{t+1}$  to the server
19:  end for
20:   $G^{t+1} = \frac{1}{M} \sum_{l=1}^M G_l^{t+1}$ 
21: end for

```

---

**MARINA** [Gorbunov et al., 2021] is an advanced method that significantly enhances communication efficiency in non-convex distributed learning across heterogeneous datasets. Its core innovation lies in a communication reduction mechanism that compresses the differences between gradients. The communication complexity bounds for **MARINA** are known to be better than those of all previous first-order methods. Non-smooth convex analysis of **MARINA** with different stepsize strategies can be found in [Sokolov and Richtárik, 2024]. This section is devoted to **Fed-Bernoulli-LoRA-MARINA** (Algorithm 8), a method within the **Bernoulli-LoRA** framework, based on **MARINA** algorithm.

In order to start convergence analysis, we rewrite the updates  $W^t, G^t$  generated by **Fed-Bernoulli-LoRA-MARINA** (Algorithm 8):

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases} \quad (120)$$

$$G_l^{t+1} = \begin{cases} \nabla f_l(W^{t+1}), & \text{with probability } q \\ G_l^t + \mathcal{Q}_l^t (\nabla f_l(W^{t+1}) - \nabla f_l(W^t)), & \text{with probability } 1 - q \end{cases} \quad (121)$$

$$G^{t+1} = \frac{1}{M} \sum_{l=1}^M G_l^{t+1}. \quad (122)$$

**Lemma 9.** Let Assumption 3 hold. Then iterates generated by Fed-Bernoulli-LoRA-MARINA satisfy

$$\mathbb{E} \left[ \|G^{t+1} - \nabla f(W^{t+1})\|_F^2 \right] \leq (1-q)\mathbb{E} \left[ \|G^t - \nabla f(W^t)\|_F^2 \right] + (1-q)\frac{\omega L^2}{M}\mathbb{E} \left[ \|W^{t+1} - W^t\|_F^2 \right]. \quad (123)$$

*Proof.* Taking the conditional expectation with respect to  $W^{t+1}$  and defining  $D_l^{t+1} := \nabla f_l(W^{t+1}) - \nabla f_l(W^t)$ ,  $D^{t+1} = \frac{1}{M} \sum_{l=1}^M D_l^{t+1}$ , we obtain

$$\begin{aligned} \mathbb{E} \left[ \|G^{t+1} - \nabla f(W^{t+1})\|_F^2 | W^{t+1} \right] &= (1-q)\mathbb{E} \left[ \left\| G^t - \nabla f(W^t) + \frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t (\nabla f_l(W^{t+1}) - \nabla f_l(W^t)) \right\|_F^2 | W^{t+1} \right] \\ &\stackrel{(15)}{=} (1-q) \|G^t - \nabla f(W^t)\|_F^2 + (1-q)\mathbb{E} \left[ \left\| \frac{1}{M} \sum_{l=1}^M \mathcal{Q}_l^t (D_l^{t+1}) - D^{t+1} \right\|_F^2 | W^{t+1} \right] \\ &= (1-q) \|G^t - \nabla f(W^t)\|_F^2 + \frac{1-q}{M^2} \sum_{m=1}^M \mathbb{E} \left[ \left\| \mathcal{Q}_l^t (D_l^{t+1}) - D_l^{t+1} \right\|_F^2 | W^{t+1} \right] \\ &\stackrel{(16)}{\leq} (1-q) \|G^t - \nabla f(W^t)\|_F^2 + \frac{(1-q)\omega}{M^2} \sum_{l=1}^M \mathbb{E} \left[ \|\nabla f_l(W^{t+1}) - \nabla f_l(W^t)\|_F^2 \right] \\ &\leq (1-q) \|G^t - \nabla f(W^t)\|_F^2 + \frac{(1-q)\omega L^2}{M} \|W^{t+1} - W^t\|_F^2, \end{aligned}$$

where in the last inequality we used that the gradient of each  $f_l$  is Lipschitz continuous.  $\square$

## D.2.1 Convergence for Smooth Non-Convex Functions

**Theorem 19.** Let Assumptions 1, 2, 3, and hold, and let the stepsize satisfy

$$0 < \gamma \leq \frac{1}{L \left( 1 + \sqrt{\lambda_{\max}^p \frac{1-q}{q} \cdot \frac{\omega}{M}} \right)}.$$

Then the iterates of Fed-Bernoulli-LoRA-MARINA (Algorithm 8) satisfy

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T} + \frac{\|G^0 - \nabla f(W^0)\|_F^2}{qT} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p}, \quad (124)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\max}^p}{2q} \|G^t - \nabla f(W^t)\|_F^2. \quad (125)$$

By Lemma 5 and Lemma 9, we have

$$\begin{aligned} \mathbb{E} [\Phi_{t+1}] &\leq \mathbb{E} [f(W^t)] - f^* - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbb{E} [\|W^{t+1} - W^t\|_F^2] \\ &\quad + \frac{\gamma \lambda_{\max}^p}{2} \mathbb{E} [\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma(1-q)\lambda_{\max}^p}{2q} \mathbb{E} [\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad + \frac{\gamma(1-q)L^2\omega\lambda_{\max}^p}{2qM} \mathbb{E} [\|W^{t+1} - W^t\|_F^2] \\ &\leq \mathbb{E} [\Phi_t] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E} [\|\nabla f(W^t)\|_F^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-q)L^2\omega\lambda_{\max}^p}{2qM} \right) \mathbb{E} [\|W^{t+1} - W^t\|_F^2]. \end{aligned}$$

Selecting  $0 < \gamma \leq \frac{1}{L(1 + \sqrt{\lambda_{\max}^p \frac{1-q}{q} \cdot \frac{\omega}{M}})}$ , we obtain

$$\mathbb{E}[\Phi_{t+1}] \leq \mathbb{E}[\Phi_t] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2].$$

Summing over, we get

$$\frac{\gamma \lambda_{\min}^p}{2} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W^t)\|_F^2] \leq \mathbb{E}[\Phi_0] - \mathbb{E}[\Phi_T].$$

Finally, we derive

$$\mathbb{E}[\|\nabla f(\widetilde{W}^T)\|_F^2] \leq \frac{2\Phi_0}{\lambda_{\min}^p \gamma T}.$$

where  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ . □

### D.2.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 20.** *Let Assumptions 1, 2, 3, and 6 hold, and let the stepsize satisfy*

$$0 < \gamma \leq \min \left\{ \frac{1}{L(1 + \sqrt{2\lambda_{\max}^p \frac{1-q}{q} \cdot \frac{\omega}{M}})}, \frac{q}{2\mu\lambda_{\min}^p} \right\}.$$

*Then the iterates of Fed-Bernoulli-LoRA-MARINA (Algorithm 8) satisfy*

$$\mathbb{E}[f(W^T) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0, \quad (126)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\Phi_0 = f(W^0) - f^* + \frac{\gamma\lambda_{\max}^p}{q} \|G^0 - \nabla f(W^0)\|_F^2$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma\lambda_{\max}^p}{q} \|G^t - \nabla f(W^t)\|_F^2. \quad (127)$$

By Lemma 5 and Lemma 7, we have

$$\begin{aligned} \mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\ &\quad + \frac{\gamma\lambda_{\max}^p}{2} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma(1-q)\lambda_{\max}^p}{q} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad + \frac{\gamma(1-q)L^2\lambda_{\max}^p}{q} \cdot \frac{\omega}{M} \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\ &\leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[f(W^t) - f^*] + \left(1 - \frac{q}{2}\right) \frac{\gamma\lambda_{\max}^p}{q} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] \\ &\quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-q)L^2\lambda_{\max}^p}{q} \cdot \frac{\omega}{M}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2], \end{aligned}$$

where in the last inequality we used Assumption 6. Selecting  $0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1-q)\omega}{qM}} \lambda_{\max}^p \right)}, \frac{q}{2\mu\lambda_{\min}^p} \right\}$ , we obtain

$$\mathbb{E} [\Phi_{t+1}] \leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E} [\Phi_t].$$

Taking recursion, we have

$$\mathbb{E} [\Phi_T] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0.$$

□

### D.3 Analysis of Fed-Bernoulli-LoRA-EF21

---

**Algorithm 9** Fed-Bernoulli-LoRA-EF21

---

```

1: Parameters: pre-trained model  $W^0 \in \mathbb{R}^{m \times n}$ ,  $\{G_l^0\}_{l \in [M]} \in \mathbb{R}^{m \times n}$  rank  $r \ll \min\{m, n\}$ , scaling
   factor  $\alpha > 0$ , chain length  $T$ , sketch distribution  $\mathcal{D}_S^B$  or  $\mathcal{D}_S^A$ , Bernoulli probability  $p$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $c^t \sim \text{Be}(p)$  Bernoulli random variable
4:   if  $c^t = 1$  then
5:     Sample  $B_S^t \sim \mathcal{D}_S^B$  Left sketch
6:      $\hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t$ 
7:      $W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$ 
8:   else
9:     Sample  $A_S^t \sim \mathcal{D}_S^A$  Right sketch
10:     $\hat{B}^t = -\eta G^t (A_S^t)^\top (A_S^t (A_S^t)^\top)^\dagger$ 
11:     $W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$ 
12:   end if
13:   Broadcast  $W^{t+1}$  to each client  $l \in [M]$ 
14:   for any client  $l \in [M]$  in parallel do
15:     Compute gradient  $\nabla f_l(W^{t+1})$ 
16:      $G_l^{t+1} = G_l^t + \mathcal{C}_l^t (\nabla f_l(W^{t+1}) - G_l^t)$ 
17:     Send  $G_l^{t+1}$  to the server
18:   end for
19:    $G^{t+1} = \frac{1}{M} \sum_{l=1}^M G_l^{t+1}$ 
20: end for

```

---

Error Feedback (EF) [Seide et al., 2014, Stich et al., 2018, Alistarh et al., 2018, Richtárik et al., 2021, Fatkhullin et al., 2021, Richtárik et al., 2022, Khirirat et al., 2024], often referred to as error compensation, is an exceptionally influential mechanism for stabilizing convergence in distributed training of supervised machine learning models, particularly when contractive communication compression techniques are employed. We design Fed-Bernoulli-LoRA-EF21 within the Bernoulli-LoRA framework, based on EF-21 method. Our theoretical analysis, built on standard assumptions, applies to distributed training in heterogeneous data settings and achieves the best known convergence rates.

Compared to Fed-Bernoulli-LoRA-MARINA, in this section we work with the wider class of compression operators called contractive.

**Definition 4.** A randomized operator  $\mathcal{C} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$  is called a contractive compression operator (compressor) if it satisfies the following condition: there exists a constant  $0 < \beta \leq 1$  such that

$$\mathbb{E} \left[ \|\mathcal{C}(W) - W\|_F^2 \right] \leq (1 - \beta) \|W\|_F^2, \quad \forall W \in \mathbb{R}^{m \times n}. \quad (128)$$

The iterates of Fed-Bernoulli-LoRA-EF21 can be rewritten as follows

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases} \quad (129)$$

$$G_l^{t+1} = G_l^t + \mathcal{C}_l^t (\nabla f_l(W^{t+1}) - G_l^t), \quad \forall l \in [M] \quad (130)$$

$$G^{t+1} = \frac{1}{M} \sum_{l=1}^M G_l^{t+1}. \quad (131)$$

**Lemma 10.** *Let Assumption 3 hold. Then for the iterates generated by Fed-Bernoulli-LoRA-EF21 (Algorithm 9) satisfy*

$$\mathbb{E} \left[ \|G_l^{t+1} - \nabla f_l(W^{t+1})\|_F^2 \right] \leq \sqrt{1-\beta} \mathbb{E} \left[ \|G_l^t - \nabla f_l(W^t)\|_F^2 \right] + \frac{(1-\beta)L^2}{1-\sqrt{1-\beta}} \mathbb{E} \left[ \|W^{t+1} - W^t\|_F^2 \right]$$

*Proof.* For each  $l \in [M]$  we have

$$\begin{aligned} \mathbb{E} \left[ \|G_l^{t+1} - \nabla f_l(W^{t+1})\|_F^2 \right] &\stackrel{(130),(131)}{=} \mathbb{E} \left[ \mathbb{E} \left[ \|C_l^t (\nabla f_l(W^{t+1}) - G_l^t) - (\nabla f_l(W^{t+1}) - G_l^t)\|_F^2 \mid G_l^{t+1}, W^{t+1} \right] \right] \\ &\stackrel{(128)}{\leq} (1-\beta) \mathbb{E} \left[ \|G_l^t - \nabla f_l(W^{t+1})\|_F^2 \right] \\ &\leq (1-\beta) (1+\theta) \mathbb{E} \left[ \|G_l^t - \nabla f_l(W^t)\|_F^2 \right] \\ &\quad + (1-\beta) \left( 1 + \frac{1}{\theta} \right) \mathbb{E} \left[ \|\nabla f_l(W^{t+1}) - \nabla f_l(W^t)\|_F^2 \right], \end{aligned}$$

where in the last inequality we used  $\|U + V\|_F^2 \leq (1+\theta) \|U\|_F^2 + (1+\frac{1}{\theta}) \|V\|_F^2$  for any constant  $\theta > 0$ , and matrices  $U, V \in \mathbb{R}^{m \times n}$ . Taking  $\theta = \frac{1}{\sqrt{1-\beta}} - 1$ , we acquire

$$\begin{aligned} \mathbb{E} \left[ \|G_l^{t+1} - \nabla f_l(W^{t+1})\|_F^2 \right] &\leq \sqrt{1-\beta} \mathbb{E} \left[ \|G_l^t - \nabla f_l(W^t)\|_F^2 \right] + \frac{1-\beta}{1-\sqrt{1-\beta}} \mathbb{E} \left[ \|\nabla f_l(W^{t+1}) - \nabla f_l(W^t)\|_F^2 \right] \\ &\leq \sqrt{1-\beta} \mathbb{E} \left[ \|G_l^t - \nabla f_l(W^t)\|_F^2 \right] + \frac{(1-\beta)L^2}{1-\sqrt{1-\beta}} \mathbb{E} \left[ \|W^{t+1} - W^t\|_F^2 \right], \end{aligned}$$

where in the last inequality we used that the gradient of each  $f_l$  is Lipschitz continuous. Summing over  $l$  from 1 to  $M$ , we finish the proof.  $\square$

### D.3.1 Convergence for Smooth Non-Convex Functions

**Theorem 21.** *Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy*

$$0 < \gamma \leq \frac{1}{L \left( 1 + \frac{\sqrt{\lambda_{\max}^p (1-\beta)}}{1-\sqrt{1-\beta}} \right)}.$$

*Then the iterates of Fed-Bernoulli-LoRA-EF21 (Algorithm 9) satisfy*

$$\mathbb{E} \left[ \left\| \nabla f(\widetilde{W}^T) \right\|_F^2 \right] \leq \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T} + \frac{\mathcal{G}^0}{(1-\sqrt{1-\beta})T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p}, \quad (132)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ , and  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ ,  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ , and  $\mathcal{G}^0 := \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_F^2$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\max}^p}{2(1-\sqrt{1-\beta})} \cdot \frac{1}{M} \sum_{l=1}^M \|G_l^t - \nabla f_l(W^t)\|_F^2. \quad (133)$$



By Lemma 5 and Lemma 10, we have

$$\begin{aligned}
\mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\
&\quad + \frac{\gamma\lambda_{\max}^p}{2} \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma\lambda_{\max}^p\sqrt{1-\beta}}{2(1-\sqrt{1-\beta})} \cdot \frac{1}{M} \sum_{l=1}^M \mathbb{E}[\|G_l^t - \nabla f_l(W^t)\|_F^2] \\
&\quad + \frac{\gamma\lambda_{\max}^p L^2(1-\beta)}{2(1-\sqrt{1-\beta})^2} \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\
&\leq \mathbb{E}[\Phi_t] - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma\lambda_{\max}^p L^2(1-\beta)}{2(1-\sqrt{1-\beta})^2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2].
\end{aligned}$$

Selecting  $0 < \gamma \leq \frac{1}{L\left(1 + \frac{\sqrt{\lambda_{\max}^p(1-\beta)}}{1-\sqrt{1-\beta}}\right)}$ , we obtain

$$\mathbb{E}[\Phi_{t+1}] \leq \mathbb{E}[\Phi_t] - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2].$$

Summing over  $t$  from 0 to  $T-1$ , we get

$$\frac{\gamma\lambda_{\min}^p}{2} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W^t)\|_F^2] \leq \mathbb{E}[\Phi_0] - \mathbb{E}[\Phi_T].$$

Finally, dividing both sides by  $\frac{\gamma\lambda_{\min}^p}{2}$  yields

$$\mathbb{E}[\|\nabla f(\widetilde{W}^T)\|_F^2] \leq \frac{2\Phi_0}{\gamma\lambda_{\min}^p T}.$$

where  $\widetilde{W}^T$  is drawn uniformly at random from the iterate sequence  $\{W^0, W^1, \dots, W^{T-1}\}$ .  $\square$

### D.3.2 Convergence under Polyak-Łojasiewicz Condition

**Theorem 22.** Let Assumptions 1, 2, 3, and 6 hold, and let the stepsize satisfy

$$0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + \frac{\sqrt{2\lambda_{\max}^p(1-\beta)}}{1-\sqrt{1-\beta}} \right)}, \frac{1 + \sqrt{1-\beta}}{2\mu\lambda_{\min}^p} \right\}$$

. Then the iterates of Fed-Bernoulli-LoRA-EF21 (Algorithm 9) satisfy

$$\mathbb{E}[f(W^T) - f^*] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0, \quad (134)$$

where  $\lambda_{\min}^p := p\lambda_{\min}^{H_B} + (1-p)\lambda_{\min}^{H_A}$ ,  $\lambda_{\max}^p := p\lambda_{\max}^{H_B} + (1-p)\lambda_{\max}^{H_A}$ , and  $\Phi_0 = f(W^0) - f^* + \frac{\gamma\lambda_{\max}^p}{1-\sqrt{1-\beta}} \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_F^2$ .

*Proof.* Denote Lyapunov function  $\Phi_t$  as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma\lambda_{\max}^p}{1-\sqrt{1-\beta}} \cdot \frac{1}{M} \sum_{l=1}^M \|G_l^t - \nabla f_l(W^t)\|_F^2. \quad (135)$$

By Lemma 5 and Lemma 10, we have

$$\begin{aligned}
\mathbb{E}[\Phi_{t+1}] &\leq \mathbb{E}[f(W^t)] - f^* - \frac{\gamma\lambda_{\min}^p}{2} \mathbb{E}[\|\nabla f(W^t)\|_F^2] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\
&\quad + \frac{\gamma\lambda_{\max}^p}{2} \cdot \mathbb{E}[\|G^t - \nabla f(W^t)\|_F^2] + \frac{\gamma\lambda_{\max}^p\sqrt{1-\beta}}{1-\sqrt{1-\beta}} \cdot \frac{1}{M} \sum_{l=1}^M \mathbb{E}[\|G_l^t - \nabla f_l(W^t)\|_F^2] \\
&\quad + \frac{\gamma\lambda_{\max}^p(1-\beta)L^2}{(1-\sqrt{1-\beta})^2} \mathbb{E}[\|W^{t+1} - W^t\|_F^2] \\
&\leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[f(W^t) - f^*] + \frac{\gamma\lambda_{\max}^p(1+\sqrt{1-\beta})}{2(1-\sqrt{1-\beta})} \cdot \frac{1}{M} \sum_{l=1}^M \mathbb{E}[\|G_l^t - \nabla f_l(W^t)\|_F^2] \\
&\quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma\lambda_{\max}^p(1-\beta)L^2}{(1-\sqrt{1-\beta})^2}\right) \mathbb{E}[\|W^{t+1} - W^t\|_F^2],
\end{aligned}$$

where in the last inequality we used Assumption 6. Selecting  $0 < \gamma \leq \min \left\{ \frac{1}{L \left( 1 + \frac{\sqrt{2\lambda_{\max}^p(1-\beta)}}{1-\sqrt{1-\beta}} \right)}, \frac{1+\sqrt{1-\beta}}{2\mu\lambda_{\min}^p} \right\}$ ,

we obtain

$$\mathbb{E}[\Phi_{t+1}] \leq (1 - \gamma\mu\lambda_{\min}^p) \mathbb{E}[\Phi_t].$$

Taking the recursion, we have

$$\mathbb{E}[\Phi_T] \leq (1 - \gamma\mu\lambda_{\min}^p)^T \Phi_0.$$

□

complete it was that from new reps

## E Experiments: Missing Details

In this section, we provide additional details regarding the experimental setting from Section 8.

### E.1 Linear Regression with Non-convex Regularization

**Full gradient setting.** We begin by evaluating these methods in a standard optimization setting where full gradients are computed at each iteration. In this regime, we compare **Bernoulli-LoRA-GD** and **RAC-LoRA-GD**.

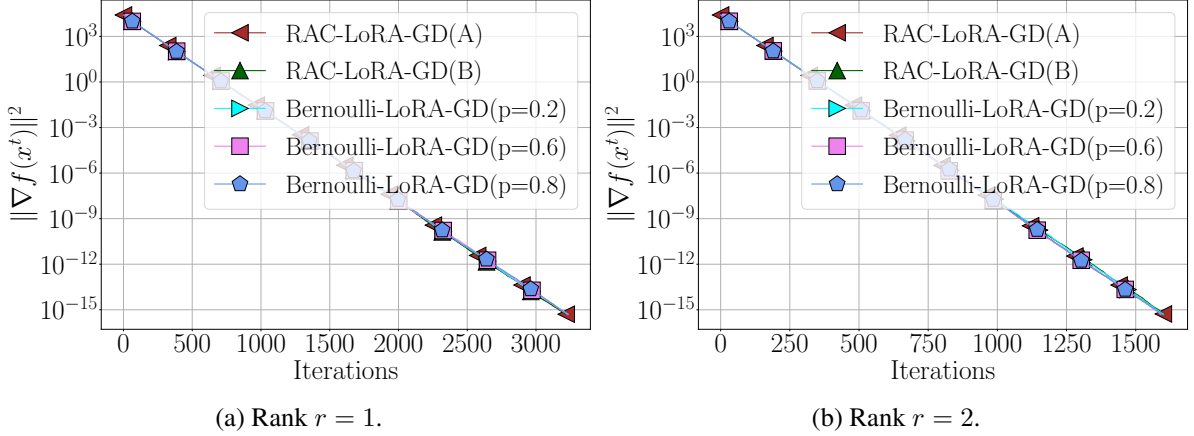


Figure 2: Comparison of **RAC-LoRA-GD** and **Bernoulli-LoRA-GD** on linear regression fine-tuning. Curves with  $p = 0.01, 0.2, \dots$  indicate **Bernoulli-LoRA-GD** sampling parameters. **RAC-LoRA-GD(A)** trains  $B$  after resampling  $A$ , while **RAC-LoRA-GD(B)** does the reverse. All methods use  $\gamma = c/\hat{L}$  with  $c \in \{1, 2\}$  tuned individually.

Figure 2 shows that, across all tested probabilities, **Bernoulli-LoRA-GD** and both variants of **RAC-LoRA-GD** exhibit similar convergence on the linear regression task. This numerical stability suggests that the ratio of updates between  $A$  and  $B$  has little effect on the performance for this problem. We also observe that higher ranks  $r$  produce faster convergence, which aligns with the theoretical  $r/n$  factor in our analysis.

**Hardware and Software.** All algorithms were implemented in Python 3.10 and executed on three different CPU cluster node types:

1. AMD EPYC 7702 64-Core,
2. Intel(R) Xeon(R) Gold 6148 CPU @ 2.40GHz,
3. Intel(R) Xeon(R) Gold 6248 CPU @ 2.50GHz.

**Implementation Details.** For each method, we set the stepsize to  $\gamma = c/\hat{L}$ , where  $c$  is a constant multiplier tuned individually for every algorithm. Convergence was monitored by computing the squared norm of the full gradient at each iteration. The algorithms terminated when either a maximum iteration limit was reached or the criterion  $\|\nabla f(x^t)\|_2^2 \leq 5 \times 10^{-16}$  was satisfied. To ensure reliability, each method was run 20 times using different random seeds, and all figures show the median performance over these trials.

**Datasets.** The synthetic pre-training dataset  $(\tilde{D}, \tilde{b})$  was generated using

```
sklearn.datasets.make_regression
```

with moderate noise and a controlled rank structure:

```
1 wt_D, wt_b = make_regression(n_samples=90000, n_features=4096,  
2                             n_informative=4096, noise=20.0,  
3                             bias=0.0, tail_strength=0.8,  
4                             effective_rank=64, random_state=42)
```

followed by standard scaling. The fine-tuning dataset  $(\hat{D}, \hat{b})$  was produced similarly:

```
1 h_D, h_b = make_regression(n_samples=10000, n_features=4096,  
2                             n_informative=4096//2, noise=50.0,  
3                             bias=10.0, tail_strength=0.9,  
4                             effective_rank=32, random_state=84)
```

and subsequently adjusted with a biased scaling (mean 1, standard deviation 2).