BERNOULLI-LORA: A THEORETICAL FRAMEWORK FOR RANDOMIZED LOW-RANK ADAPTATION

Anonymous authors

000

001

002003004

010

011

012

013

014

015

016

018

019

021

023

024

025

026

027

028

029

031

032

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

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.

1 Introduction

Fine-tuning adapts pre-trained models to new datasets and is central to modern deep learning, especially in NLP (Peters et al., 2018; Devlin et al., 2019). While effective, full fine-tuning is computationally intensive for large models. Parameter-Efficient Fine-Tuning (PEFT) methods (He et al., 2021) address this by updating only a fraction of parameters (Richtárik & Takáč, 2016; Demidovich et al., 2023a), matching full fine-tuning performance with significantly lower training cost (Radford et al., 2019; Brown et al., 2020; Han et al., 2024).

Pre-trained models have a low intrinsic dimensionality (Li et al., 2018; Aghajanyan et al., 2020), so fine-tuning is effective in a reduced-dimensional subspace. Low-Rank Adaptation (LoRA) (Hu et al., 2022) is a prominent reparameterization technique that exploits this. It avoids updating large weight matrices by optimizing two small low-rank matrices whose product forms the update:

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

where $W^0 \in \mathbb{R}^{m \times n}$ is fixed, and $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$ are trainable. Typically, A is initialized from a Gaussian distribution and B as zero, though other strategies exist (Zhu et al., 2024; Hayou et al., 2024; Meng et al., 2024; Wang et al., 2024). Choosing a rank $r \ll \min\{m,n\}$ drastically reduces trainable parameters. While not always matching full fine-tuning, LoRA better mitigates catastrophic forgetting and enhances output diversity (Biderman et al., 2024). It is easy to implement and performs well across many tasks (Hu et al., 2022). Research has also improved LoRA's computational efficiency (Cherniuk et al., 2023; Mao et al., 2025).

To bridge the gap with full fine-tuning, Xia et al. (2024) introduced Chain of LoRA (COLA), an iterative framework that builds higher-rank updates from a sequence of low-rank components at

no extra computational cost. It iteratively trains a LoRA module, merges its updates into the fixed parameters, and repeats. This yields successive updates:

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

Unlike standard LoRA, COLA uses sequential low-rank decompositions to approximate higher-rank updates, improving adaptation efficiency.

2 PROBLEM STATEMENT

Supervised learning is an optimization problem that minimizes a loss function. We focus on this challenge in fine-tuning, using a general, model-agnostic formulation:

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

Here, W^0 represents the pre-trained parameters, ΔW is the trainable adaptation, and f is the empirical loss. Since $m \times n$ is very large, ΔW requires a simple, trainable structure.

For our stochastic methods, we consider these objective structures:

 Finite-Sum Setting: The objective is an average over N data samples, used in methods like Bernoulli-LoRA-PAGE:

$$f(W^{0} + \Delta W) = \frac{1}{N} \sum_{i=1}^{N} f_{i}(W^{0} + \Delta W).$$
 (2)

 Expectation Setting: The objective is an expectation over a data distribution D, for methods like Bernoulli-LoRA-MVR:

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

We also address the **distributed optimization setting** for our proposed Federated Learning (FL) algorithms (e.g., Fed-Bernoulli-LoRA-QGD). Here, the goal is to minimize a global objective averaged over M clients:

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

where f_l is the local loss for client l. The goal is to find ΔW that minimizes this global objective.

3 MOTIVATION

Despite their success, Low-Rank Adaptation (LoRA) and its variants like Chain of LoRA (COLA) lack a solid theoretical foundation. Key gaps exist. First, LoRA's re-parameterization makes smooth loss functions non-smooth, creating significant theoretical hurdles (Sun et al., 2024). Second, existing COLA analysis ignores its core low-rank updates by focusing on full-rank optimization, thus failing to explain its efficiency (Xia et al., 2024). Consequently, most LoRA-based methods are heuristics without convergence guarantees, making them sensitive to hyperparameters (Khodak et al., 2021; Kuang et al., 2024). Malinovsky et al. (2024) even showed COLA can diverge and introduced RAC-LoRA, the first framework to establish convergence rates for LoRA-style updates. However, the RAC-LoRA framework is limited. It lacks optimal variance-reduced techniques for non-convex problems and fails to address advanced Federated Learning (FL) scenarios incorporating communication compression and error feedback (Alistarh et al., 2018; Wen et al., 2017; Horvóth et al., 2022; Panferov et al., 2024). The need for distributed optimization like FL is driven by the challenge of training massive models (Brown et al., 2020; Kolesnikov et al., 2020; Goyal et al., 2017; You et al., 2019; Le Scao et al., 2023). Our work aims to bridge this gap by extending a theoretically sound LoRA framework to these vital, practical optimization settings.

4 CONTRIBUTIONS

LoRA-based methods are sensitive to hyperparameters (Khodak et al., 2021; Kuang et al., 2024) and require a stronger theoretical foundation. While Malinovsky et al. (2024) provided an initial framework with RAC-LoRA, we aim to further advance the theory and versatility of low-rank adaptation.

Low-rank PEFT updates two matrices, A and B, either individually or alternating deterministically (Malinovsky et al., 2024; Xia et al., 2024; Zhu et al., 2024). Our main contribution, Bernoulli-LoRA, is a generic framework with a probabilistic update: at each step, a Bernoulli trial selects either A or B for optimization while the other is fixed. This randomized selection unifies and generalizes existing update strategies. Similar to COLA (Xia et al., 2024), our framework applies a sequence of these probabilistic low-rank updates.

Our analysis uses standard non-convex optimization assumptions, like L-smoothness. We instantiate Bernoulli-LoRA with several algorithms, from foundational gradient methods to advanced stochastic, variance-reduced, and federated learning variants. We establish rigorous convergence guarantees for each method. Our key contributions include:

- Foundational Algorithmic Variants: We establish the framework's properties with two fundamental schemes to understand the interplay between randomized selection and standard descent.
 - Bernoulli-LoRA-GD (Algorithm 2) uses full gradients, providing a foundational analysis of convergence in an idealized setting.
 - Bernoulli-LoRA-SGD (Algorithm 4) uses practical stochastic gradients, offering insights into the interplay of stochasticity and randomized adaptation for large-scale tasks.
- ◆ Advanced Variance Reduction for Non-Convex Optimization: To counter variance from stochastic gradients, we develop VR-enhanced variants, providing the first theoretical analysis of LoRA-type methods with advanced VR schemes in L-smooth non-convex settings.
 - Bernoulli-LoRA-PAGE (Algorithm 6) adapts the optimal and simple PAGE (Li et al., 2021) for the finite-sum setting (2).
 - Bernoulli-LoRA-MVR (Algorithm 5) uses Momentum Variance Reduction inspired by STORM (Cutkosky & Orabona, 2019) for the expectation setting, proving its effectiveness in our framework.
- ◆ Communication-Efficient Federated Learning Extensions: We extend Bernoulli-LoRA to FL, addressing communication overhead. We provide the first comprehensive analysis of LoRA-type methods integrated with established communication-saving techniques like quantization, gradient difference compression, and error feedback.
 - Fed-Bernoulli-LoRA-QGD (Algorithm 7) incorporates QSGD-style quantization (Alistarh et al., 2017; Wen et al., 2017; Horvóth et al., 2022; Panferov et al., 2024) to compress gradients and reduce communication bandwidth.
 - Fed-Bernoulli-LoRA-MARINA (Algorithm 8) adapts the MARINA strategy (Gorbunov et al., 2021) to efficiently compress gradient differences.
 - Fed-Bernoulli-LoRA-EF21 (Algorithm 9) integrates the EF21 error feedback mechanism (Richtárik et al., 2021) to stabilize training with contractive compressors.
- ◆ Analysis for Non-Smooth Convex Functions: We broaden our framework's applicability by providing the first theoretical analysis of LoRA-type methods for non-smooth convex optimization. We present a version of Bernoulli-LoRA-GD (Algorithm 3) and establish its convergence rates with different stepsize policies.

5 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.

Setting	Method	NC convergence rate	PŁ convergence rate
(1)	Bernoulli-LoRA-GD (Alg. 2)	$rac{\Delta^0}{\gamma \lambda_{ ext{min}} T}$	$(1 - \gamma \mu \lambda_{\min})^T \Delta^0$
(1)	Bernoulli-LoRA-SGD (Alg. 4)	$\frac{\Delta^0}{\gamma \lambda_{\min} T} + \frac{\gamma L C_1 \lambda_{\max}}{\lambda_{\min}}$	$(1 - \gamma \mu \lambda_{\min})^T \Delta^0 + \frac{\gamma L C_1 \lambda_{\max}}{\mu \lambda_{\min}}$
(1)+(3)	Bernoulli-LoRA-MVR (Alg. 5)	$\frac{\Phi_1}{\gamma \lambda_{\min} T} + \frac{b\sigma^2 \lambda_{\max}}{(2-b)\lambda_{\min}} $ (1)	$(1 - \gamma \mu \lambda_{\min})^T \Phi_1 + \frac{b\sigma^2 \lambda_{\max}}{(2-b)\mu \lambda_{\min}}$
(1)+(2)	Bernoulli-LoRA-PAGE (Alg. 6)	$\frac{\Phi_2}{\gamma \lambda_{\min} T}$ (2)	$(1 - \gamma \mu \lambda_{\min})^T \Phi_2^{(2)}$
(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)	$\frac{\Phi_2}{\gamma \lambda_{\min} T}$ (2)	$(1 - \gamma \mu \lambda_{\min})^T \Phi_2^{(2)}$
(1)+(4)	Fed-Bernoulli-LoRA-EF21 (Alg. 9)	$\frac{\Phi_3}{\gamma \lambda_{\min} T} \ ^{(3)}$	$(1 - \gamma \mu \lambda_{\min})^T \Phi_3^{(3)}$

 $[\]begin{split} & \overset{(1)}{\Phi_1} & = \Delta^0 + \frac{\gamma}{b(2-b)} \mathcal{G}^0; \\ & \overset{(2)}{\Phi_2} & = \Delta^0 + \frac{\gamma}{q} \mathcal{G}^0; \\ & \overset{(3)}{\Phi_3} & = \Delta^0 + \frac{\gamma}{1-\sqrt{1-\beta}} \hat{\mathcal{G}}^0 \end{split}$

173 174 175

176

177

184

185

186

187

188 189

190 191

192

193

194 195

196 197

198

199 200 201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

13: **end for**

Table 1: Summary of the convergence rates for the proposed methods, presented for smooth non-convex functions ("NC") and for functions satisfying the PŁ-condition ("PŁ"). 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; ω 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; β is the contraction parameter; b is the momentum parameter; b is the probability of a full gradient computation; b is the contraction parameter; b is the momentum parameter; b is the momentum parameter; b is the probability of a full gradient computation; b is the contraction b is the momentum parameter; b is the m

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

$$\Delta W = -\frac{\alpha}{r} B_S \hat{A},\tag{5}$$

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

$$\Delta W = -\frac{\alpha}{r}\hat{B}A_S,\tag{6}$$

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.
    for t = 0, 1, ..., T - 1 do
         Sample c^t \sim \text{Be}(p)
                                                                                                      Bernoulli random variable
        if c^t = 1 then
 4:
 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).
            W^{t+1} = W^t + \frac{\alpha}{\pi} B_S^t \hat{A}^t.
 7:
 8:
        else
            Sample A_S^t \sim \mathcal{D}_S^A
 9:
                                                                                                                     (Right sketch)
10:
            Using a chosen optimizer, approximately solve \hat{B}^t \approx \arg\min_B f(W^t + \frac{\alpha}{r}BA_S^t).
            W^{t+1} = W^t + \frac{\alpha}{\pi} \hat{B}^t A_S^t.
11:
12:
         end if
```

_	_	_	
	1		
	1		
	1		
	1		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	2		
	3		
	3		
	3		
	3		
	3		
	3		
	3		
	3		
	3		
	3		
	4		
	4		
	4		
	4		
	4		
	4		
	4		
	4		
	4		
	4		
	5		
	5		
	5		
	5		
	5		
	5		
2	5	6	
2	5	7	
	5		
2	5	9	
2	6	0	
2	6	1	
2	6	2	
2	6	3	
2	6	4	
2	6	5	
2	6	6	
2	6	7	
2	6	8	

Setting	Method	Base Gradient Estimator G^t	Thms. #
(1)	Bernoulli-LoRA-GD (Algs. 2 & 3)	$G^t = \nabla f(W^t)$	1 & 9 & 10
(1)	Bernoulli-LoRA-SGD (Alg. 4)	$G^t = g(W^t)$	11 & 12
(1)+(3)	Bernoulli-LoRA-MVR (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)	Bernoulli-LoRA-PAGE (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)	Fed-Bernoulli-LoRA-QGD (Alg. 7)	$G^{t} = \frac{1}{M} \sum_{l=1}^{M} \mathcal{Q}_{l}^{t}(\nabla f_{l}(W^{t}))$	17 & 18
(1)+(4)	Fed-Bernoulli-LoRA-MARINA (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)	Fed-Bernoulli-LoRA-EF21 (Alg. 9)	$ \begin{aligned} G^t &= \frac{1}{M} \sum_{l=1}^{M} G_l^t \\ \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 \end{aligned} $	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.

5.1 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. More discussion can be found in Appendix E.

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:

$$\hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}$$
 (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 Convergence Results

The convergence properties of our framework hinge on the spectral properties of the expected projection matrix, which is introduced in Section 5.1. 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 D (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 ε -suboptimal solution: a random matrix $\hat{W} \in \mathbb{R}^{m \times n}$ that satisfies

$$\mathbb{E}\left[f(\hat{W}) - f(W^*)\right] \le \varepsilon,\tag{9}$$

where $\mathbb{E}\left[\cdot\right]$ 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 ε -stationary point: a random matrix $\hat{W} \in \mathbb{R}^{m \times n}$ satisfying

$$\mathbb{E}\left[\left\|\nabla f(\hat{W})\right\|_{\mathcal{F}}^{2}\right] \leq \varepsilon^{2}.\tag{10}$$

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 & 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

$$\left\|\nabla f(W) - \nabla f(V)\right\|_{F} \le L \left\|W - V\right\|_{F},$$

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 H.1.

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

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

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 & Richtárik, 2023)). The stochastic gradient estimator g(W) satisfies

$$\mathbb{E}\left[\left\|g(W)\right\|_{\mathrm{F}}^{2}\right] \leq 2A_{1}\left(f(W) - f^{*}\right) + B_{1} \cdot \left\|\nabla f(W)\right\|_{\mathrm{F}}^{2} + C_{1},$$

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 H.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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \leq \frac{6\Delta^0}{\gamma \lambda_{\min}^p T} + \gamma L C_1 \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

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

To analyze our variance-reduced methods, we consider a specific bounded variance assumption.

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

$$\mathbb{E}\left[\nabla f_{\xi}(W)\right] = \nabla f(W), \qquad \mathbb{E}\left[\left\|\nabla f_{\xi}(W) - \nabla f(W)\right\|_{F}^{2}\right] \leq \sigma^{2}.$$

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

Theorem 3. Let Assumptions 1, 2, 3, and 5 hold, and let the stepsize satisfy $0 < \gamma \le \frac{1}{L\left(1+\sqrt{\frac{2\lambda_{\max}^{p}(1-b)^{2}}{b}}\right)}$. Then the iterates of Bernoulli-LoRA-MVR (Algorithm 5) satisfy

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

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

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

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

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

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

where
$$\Delta^0 := f(W^0) - f^*$$
 and $\mathcal{G}^0 := \left\| G^0 - \nabla f(W^0) \right\|_{\mathrm{F}}^2$.

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 I.1.

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

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

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

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 I.2.

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

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

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

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

Theorem 7. Let Assumptions 1, 2, and 3 hold, and let the stepsize satisfy $0 < \gamma \le \frac{1}{L\left(1 + \frac{\sqrt{\lambda_{\max}^{p}(1-\beta)}}{1 - \sqrt{1-\beta}}\right)}$

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

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

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

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; Lojasiewicz, 1963)). There exists $\mu > 0$ such that

$$\frac{1}{2} \left\| \nabla f(W) \right\|_{\mathcal{F}}^2 \ge \mu \left(f(W) - f^* \right).$$

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

Theorem 8. Let Assumptions 2, 3, 4, and 6 hold, and let the stepsize satisfy
$$0 < \gamma \le \min\left\{\frac{\mu\lambda_{\min}^p}{2LA_1\lambda_{\max}^p}, \frac{2}{\mu\lambda_{\min}^p}, \frac{1}{LB_1}\left(\frac{\lambda_{\max}^p}{\lambda_{\min}^p}\right)^{-1}\right\}$$
. Then the iterates of Bernoulli-LoRA-SGD (Algorithm 4) satisfy

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

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

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

EXPERIMENTS

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

7.1 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 (\cdot) for pre-training quantities and (\cdot) for fine-tuning.

During the **pre-training phase**, we solve
$$\min_{x \in \mathbb{R}^n} \left\{ \widetilde{f}(x) := \frac{1}{2\widetilde{m}} \left\| \widetilde{D}x - \widetilde{b} \right\|_2^2 + \widetilde{\lambda} \sum_{j=1}^d \frac{x_j^2}{1 + x_j^2} \right\}$$

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 $\|\nabla f(\widetilde{x}^*)\|^2 \leq 10^{-8}$ to obtain \widetilde{x}^* . For the **fine-tuning phase**, we

use
$$\widetilde{x}^*$$
 as the initialization and then solve $\min_{x \in \mathbb{R}^n} \left\{ \widehat{f}(x) := \frac{1}{2\widehat{m}} \left\| \widehat{D}x - \widehat{b} \right\|_2^2 + \widehat{\lambda} \sum_{j=1}^d \frac{x_j^2}{1 + x_j^2} \right\}$,

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 $\hat{\lambda} = \|\hat{D}\|_2 \approx 4101.7$,

yielding $\ddot{L} \approx 12305.3$. 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 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.

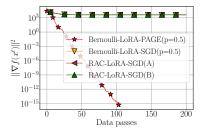


Figure 1: Comparison of RAC-LoRA-SGD, Bernoulli-LoRA-SGD and Bernoulli-LoRA-PAGE on linear regression fine-tuning. Curves with $p=0.01,0.2,\ldots$ indicate Bernoulli-LoRA sampling parameters. RAC-LoRA-SGD(A) trains B after resampling A, while RAC-LoRA-SGD(B) does the reverse. All methods use $\gamma=c/\hat{L}$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

7.2 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).

Methodology. We first pre-train a three-layer MLP on MNIST digits –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 2×10^{-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 ± 1.60	1K
LoRA	Zero	Gaussian	89.82 ± 0.90	1K
COLA	Gaussian	Zero	93.32 ± 0.50	1K
COLA	Zero	Gaussian	96.55 ± 0.20	1K
AsymmLoRA	Gaussian	Zero	64.04 ± 6.90	133
AsymmLoRA	Zero	Gaussian	74.52 ± 7.20	912
RAC-LoRA	Gaussian	Zero	93.02 ± 0.50	133
RAC-LoRA	Zero	Gaussian	96.49 ± 0.20	912
Bernoulli-LoRA	² Zero ¹	Gaussian	96.46 ± 0.17	≈ 904

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

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 J.

² 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.

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 (eds.), *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)*, pp. 4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-[]1423. URL 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*, pp. 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*, pp. 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*, pp. 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 Horvóth, Chen-Yu Ho, Ludovit Horvath, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. In *Mathematical and Scientific Machine Learning*, pp. 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 Poczos, 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, Badih 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, pp. 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.
 - 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*, pp. 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 Farkhoor, 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*, pp. 6286–6295. PMLR, 2021.
 - Stanisław 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-[]40663-[]9. URL 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 (eds.), *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)*, pp. 2227–2237. Association for Computational Linguistics, 2018. doi: 10.18653/V1/N18-[]1202. URL https://doi.org/10.18653/v1/n18-[]1202.
 - 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*, pp. 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*, pp. 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, I and Licheng Jiao. Vr-sgd: A simple stochastic variance reduction method for mach <i>IEEE Transactions on Knowledge and Data Engineering</i> , 32(1):188–202, 2018.	
Igor Sokolov and Peter Richtárik. MARINA-P: Superior performance in non-smoo optimization with adaptive stepsizes. <i>arXiv preprint arXiv:2412.17082</i> , 2024.	th federated
Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD was Advances in Neural Information Processing Systems, 31, 2018.	ith memory.
Ruo-Yu Sun. Optimization for deep learning: An overview. <i>Journal of the Operatio Society of China</i> , 8(2):249–294, 2020.	ns Research
Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. Improving LoRA in privacy-preserving learning. <i>arXiv preprint arXiv:2403.12313</i> , 2024.	ng federated
Alexander Tyurin and Peter Richtárik. DASHA: Distributed nonconvex optimization munication compression and optimal oracle complexity. In <i>The Eleventh Internation ence on Learning Representations</i> , 2023. URL https://openreview.net/fvAlYpcNr7ul.	onal Confer-
Roman Vershynin. High-dimensional probability, 2009.	
Evgeniya Vorontsova, Roland Hildebrand, Alexander Gasnikov, and Fedor Stonyak optimization. <i>arXiv preprint arXiv:2106.01946</i> , 2021.	in. Convex
Hanqing Wang, Yixia Li, Shuo Wang, Guanhua Chen, and Yun Chen. MiLoRA: Harne singular components for parameter-efficient LLM finetuning. <i>arXiv preprint arXiv</i> : 2024.	
Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai L Ternary gradients to reduce communication in distributed deep learning. <i>Advance Information Processing Systems</i> , 30, 2017.	
Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of LoRA: efficient fine-tuning models via residual learning. <i>arXiv preprint arXiv:2401.04151</i> , 2024.	of language
Kai Yi, Timur Kharisov, Igor Sokolov, and Peter Richtárik. Cohort squeeze: Beyond a si nication round per cohort in cross-device federated learning. <i>arXiv preprint arXiv</i> : 2024.	
Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapa Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimizati learning: Training BERT in 76 minutes. <i>arXiv preprint arXiv:1904.00962</i> , 2019.	
Jiacheng Zhu, Kristjan Greenewald, Kimia Nadjahi, Haitz Saez De Ocariz Borde, Ri Gabrielsson, Leshem Choshen, Marzyeh Ghassemi, Mikhail Yurochkin, and Just Asymmetry in low-rank adapters of foundation models. <i>arXiv preprint arXiv:2402.1</i>	in Solomon.
A APPENDIX	
CONTENTS	
1 Introduction	1
2 Problem Statement	2
3 Motivation	2
4 Contributions	3

5	Bernoulli-LoRA Framework			
	5.1 Reformulation as a Projected Gradient Step			
6	Convergence Results			
7	Empirements			
7	Experiments			
	7.1 Linear Regression with Non-convex Regularization.7.2 MLP on MNIST.			
	7.2 WEI ON WHOIST			
A	Appendix			
В	Basic Facts and Useful Inequalities			
C	Notation 1			
D	Discussion on Positive Expected Projection (Assumption 1)			
E	Reformulation as a Projected Gradient Step			
F	Core Algorithmic Variants			
G	Extensions for Federated Learning			
Н	Proofs for Core Algorithmic Variants			
	H.1 Analysis of Bernoulli-LoRA-GD			
	H.1.1 Convergence for Smooth Non-Convex Functions			
	H.1.2 Convergence under Polyak-Łojasiewicz Condition			
	H.1.3 Convergence for Non-Smooth Convex Functions			
	H.2 Analysis of Bernoulli-LoRA-SGD			
	H.2.1 Convergence for Smooth Non-Convex Functions			
	H.2.2 Convergence under Polyak-Łojasiewicz Condition			
	H.3 Analysis of Bernoulli-LoRA-MVR			
	H.3.1 Convergence for Smooth Non-Convex Functions			
	H.3.2 Convergence under Polyak-Łojasiewicz Condition			
	H.4 Analysis of Bernoulli-LoRA-PAGE			
	H.4.1 Convergence for Smooth Non-Convex Functions			
	H.4.2 Convergence under Polyak-Łojasiewicz Condition			
Ι	Proofs for Federated Learning Extensions			
	I.1 Analysis of Fed-Bernoulli-LoRA-QGD			
	I.1.1 Convergence for Smooth Non-Convex Functions			
	I.1.2 Convergence under Polyak-Łojasiewicz Condition			
	I.2 Analysis of Fed-Bernoulli-LoRA-MARINA			

		I.2.1	Convergence for Smooth Non-Convex Functions	50	
		I.2.2	Convergence under Polyak-Łojasiewicz Condition	51	
	I.3	Analy	sis of Fed-Bernoulli-LoRA-EF21	52	
		I.3.1	Convergence for Smooth Non-Convex Functions	53	
		I.3.2	Convergence under Polyak-Łojasiewicz Condition	54	
J	J Experiments: Missing Details				
	J.1 Linear Regression with Non-convex Regularization			55	

B BASIC FACTS AND USEFUL INEQUALITIES

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

$$\mathbb{E}\left[\mathbb{E}\left[X\mid Y\right]\right] = \mathbb{E}\left[X\right].\tag{11}$$

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

$$|\mathbb{E}[XY]| \le \sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]}.$$
(12)

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

$$\mathbb{E}\left[\|X - c\|_{2}^{2}\right] = \mathbb{E}\left[\|X - \mathbb{E}\left[X\right]\|_{2}^{2}\right] + \|\mathbb{E}\left[X\right] - c\|_{2}^{2}.$$
(13)

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

$$g(\mathbb{E}[X]) \le \mathbb{E}[g(X)].$$
 (14)

C 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 $\|\cdot\|_{\mathrm{F}}$, defined as $\|W\|_{\mathrm{F}} = \sqrt{\mathrm{Tr}(W^{\top}W)}$, where $\mathrm{Tr}(\cdot)$ denotes the matrix trace. The inner product between two matrices A and B is denoted by $\langle A, B \rangle = \mathrm{Tr}(A^{\top}B)$. 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 $\Delta^0 := f(W^0) - f^*$, $\mathcal{G}^0 := \|G^0 - \nabla f(W^0)\|_{\mathrm{F}}^2$ and $\hat{\mathcal{G}}^0 := \frac{1}{M} \sum_{l=1}^M \|G_l^0 - \nabla f_l(W^0)\|_{\mathrm{F}}^2$. 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.

D 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

$$M = QMQ^{\top}$$
 for all orthonormal matrices $Q \in \mathbb{R}^{n \times n}$. (15)

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 λ . 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 λ . If every unit vector is an eigenvector with eigenvalue λ , 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., $M = \lambda I_n$.

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:

$$H_{QB} = (QB) ((QB)^{\top} QB)^{\dagger} (QB)^{\top}$$

$$= QB (B^{\top} Q^{\top} QB)^{\dagger} B^{\top} Q^{\top}$$

$$= QB (B^{\top} B)^{\dagger} B^{\top} Q^{\top}$$

$$= Q (B(B^{\top} B)^{\dagger} B^{\top}) Q^{\top} = QH_B Q^{\top}.$$

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 $r \le n$). Thus, the rank of H_B is r.

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

By linearity of expectation and trace, we also have:

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

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.

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).

E REFORMULATION AS A PROJECTED GRADIENT STEP

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 :

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

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

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

where we define the effective stepsize $\gamma := \frac{\alpha \eta}{r}$ and the projection matrix $H_B^t := B_S^t \left(\left(B_S^t \right)^\top B_S^t \right)^\dagger \left(B_S^t \right)^\top$. 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 $H_A^t := (A_S^t)^\top \left(A_S^t \left(A_S^t \right)^\top \right)^\dagger A_S^t$. 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.

F 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 & 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 .

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 & 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 & 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.

G 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_l^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.

H PROOFS FOR CORE ALGORITHMIC VARIANTS

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

7:
$$W^{t+1} = \overset{\searrow}{W}^t + \frac{\alpha}{r} B_S^t \overset{?}{A}^t$$

8: else

9: Sample
$$A_S^t \sim \mathcal{D}_S^A$$

Right sketch

10:
$$\hat{B}^t = -\eta \nabla f(W^t) \left(A_S^t \right)^\top \left(A_S^t \left(A_S^t \right)^\top \right)^\dagger$$

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:

$$\hat{A}^{t} := \arg\min_{A} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \nabla f(W^{t}), B_{S}^{t} A \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma r^{2}} \left\| B_{S}^{t} A \right\|_{F}^{2} \right\},$$

$$\hat{B}^{t} := \arg\min_{B} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \nabla f(W^{t}), BA_{S}^{t} \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma r^{2}} \left\| BA_{S}^{t} \right\|_{F}^{2} \right\}. \tag{16}$$

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

$$W^{t+1} = W^t - \gamma G^t, \tag{17}$$

where G^t is defined by

$$G^{t} = \begin{cases} H_{B}^{t} \nabla f(W^{t}), & \text{with probability } p \\ \nabla f(W^{t}) H_{A}^{t}, & \text{with probability } 1 - p \end{cases}$$
(18)

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

$$H_A^t := \left(A_S^t\right)^\top \left(A_S^t \left(A_S^t\right)^\top\right)^\dagger A_S^t \quad \text{and} \quad H_B^t := B_S^t \left(\left(B_S^t\right)^\top B_S^t\right)^\dagger \left(B_S^t\right)^\top, \tag{19}$$

where † 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:

$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t. \tag{20}$$

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

$$\frac{\alpha}{r} \left(B_S^t \right)^\top \nabla f(W^t) + \frac{\alpha^2}{\gamma r^2} \left(B_S^t \right)^\top B_S^t \hat{A}^t = 0;$$

$$\left(B_S^t \right)^\top B_S^t \hat{A}^t = -\frac{\gamma r}{\alpha} \left(B_S^t \right)^\top \nabla f(W^t);$$

$$\hat{A}^t = -\frac{\gamma r}{\alpha} \left(\left(B_S^t \right)^\top B_S^t \right)^\dagger \left(B_S^t \right)^\top \nabla f(W^t). \quad (21)$$

This leads to the Left sketch update:

$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$$

$$= W^t - \gamma B_S^t \left(\left(B_S^t \right)^\top B_S^t \right)^\dagger \left(B_S^t \right)^\top \nabla f(W^t)$$

$$= W^t - \gamma H_B^t \nabla f(W^t), \tag{22}$$

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

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

$$W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t. \tag{23}$$

First, observe that:

$$\left\| \hat{B}^t A_S^t \right\|_{\mathrm{F}}^2 = \left\langle \hat{B}^t A_S^t, \hat{B}^t A_S^t \right\rangle_{\mathrm{F}} = \left\langle A_S^t, \left(\hat{B}^t \right)^\top \hat{B}^t A_S^t \right\rangle_{\mathrm{F}}.$$
 (24)

For the linear term from (16):

$$\frac{\alpha}{r} \left\langle \nabla f(W^t), \hat{B}^t A_S^t \right\rangle_{\mathcal{F}} = \frac{\alpha}{r} \operatorname{Tr} \left(\left(\nabla f(W^t) \right)^{\top} \hat{B}^t A_S^t \right), \tag{25}$$

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

$$\frac{\alpha^2}{\gamma r^2} \hat{B}^t A_S^t \left(A_S^t \right)^\top. \tag{26}$$

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

$$\hat{B}^{t} = -\frac{\gamma r}{\alpha} \nabla f(W^{t}) \left(A_{S}^{t} \right)^{\top} \left(A_{S}^{t} \left(A_{S}^{t} \right)^{\top} \right)^{\dagger}, \tag{27}$$

which yields the Right sketch update:

$$W^{t+1} = W^{t} + \frac{\alpha}{r} \hat{B}^{t} A_{S}^{t}$$

$$= W^{t} - \gamma \nabla f(W^{t}) \left(A_{S}^{t} \right)^{\top} \left(A_{S}^{t} \left(A_{S}^{t} \right)^{\top} \right)^{\dagger} A_{S}^{t}$$

$$= W^{t} - \gamma \nabla f(W^{t}) H_{A}^{t}, \tag{28}$$

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

Combined Update Rule. Combining equations (22) and (28), we obtain the unified update:

$$W^{t+1} = W^t - \gamma G^t, \tag{29}$$

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

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

H.1.1 CONVERGENCE FOR SMOOTH NON-CONVEX FUNCTIONS

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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \le \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T},\tag{30}$$

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

$$W^{t+1} = W^t - \gamma G^t, \tag{31}$$

where G^t takes the form

$$G^{t} = \begin{cases} H_{B}^{t} \nabla f(W^{t}), & \text{with probability } p \\ \nabla f(W^{t}) H_{A}^{t}, & \text{with probability } 1 - p \end{cases}$$
(32)

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 :

$$\mathbb{E}\left[G^{t} \mid W^{t}, H^{t}\right] = pH_{B}^{t} \nabla f(W^{t}) + (1-p)\nabla f(W^{t})H_{A}^{t},$$

$$\mathbb{E}\left[\left\|G^{t}\right\|_{F}^{2} \mid W^{t}, H^{t}\right] = p\left\|H_{B}^{t} \nabla f(W^{t})\right\|_{F}^{2} + (1-p)\left\|\nabla f(W^{t})H_{A}^{t}\right\|_{F}^{2},$$
(33)

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

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

$$-\gamma p \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F} + \frac{L\gamma^{2}}{2} p \left\| H_{B}^{t} \nabla f(W^{t}) \right\|_{F}^{2}$$

$$= -\gamma p \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F} + \frac{L\gamma^{2}}{2} p \left\langle H_{B}^{t} \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F}$$

$$= -\gamma p \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F} + \frac{L\gamma^{2}}{2} p \left\langle \nabla f(W^{t}), \left(H_{B}^{t}\right)^{\top} H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F}$$

$$= p \left(-\gamma \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F} + \frac{L\gamma^{2}}{2} \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F} \right)$$

$$\stackrel{\gamma \leq 1/L}{\leq} -\frac{\gamma}{2} p \left\langle \nabla f(W^{t}), H_{B}^{t} \nabla f(W^{t}) \right\rangle_{F}. \tag{34}$$

For any projection matrix H_A^t , we have:

$$\langle \nabla f(W^{t}) H_{A}^{t}, \nabla f(W^{t}) H_{A}^{t} \rangle_{F} = \operatorname{Tr} \left(\left(H_{A}^{t} \right)^{\top} \left(\nabla f(W^{t}) \right)^{\top} \nabla f(W^{t}) H_{A}^{t} \right)$$

$$= \operatorname{Tr} \left(\left(\nabla f(W^{t}) \right)^{\top} \nabla f(W^{t}) H_{A}^{t} \left(H_{A}^{t} \right)^{\top} \right)$$

$$= \operatorname{Tr} \left(\left(\nabla f(W^{t}) \right)^{\top} \nabla f(W^{t}) H_{A}^{t} \right)$$

$$= \langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \rangle_{F}. \tag{35}$$

Therefore:

$$-\gamma(1-p)\left\langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F} + \frac{L\gamma^{2}}{2} (1-p) \left\| \nabla f(W^{t}) H_{A}^{t} \right\|_{F}^{2}$$

$$= -\gamma(1-p) \left\langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F} + \frac{L\gamma^{2}}{2} (1-p) \left\langle \nabla f(W^{t}) H_{A}^{t}, \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F}$$

$$= -\gamma(1-p) \left\langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F} + \frac{L\gamma^{2}}{2} (1-p) \left\langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F}$$

$$\stackrel{\gamma \leq 1/L}{\leq} -\frac{\gamma}{2} (1-p) \left\langle \nabla f(W^{t}), \nabla f(W^{t}) H_{A}^{t} \right\rangle_{F}. \tag{36}$$

Using the Lipschitz gradient condition and the above bounds:

$$\mathbb{E}\left[f(W^{t+1}) \mid W^{t}, H^{t}\right] \leq f(W^{t}) + \mathbb{E}\left[\left\langle\nabla f(W^{t}), W^{t+1} - W^{t}\right\rangle_{F} \mid W^{t}, H^{t}\right] \\ + \frac{L}{2}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2} \mid W^{t}, H^{t}\right] \\ = f(W^{t}) - \gamma\left\langle\nabla f(W^{t}), \mathbb{E}\left[G^{t} \mid W^{t}, H^{t}\right]\right\rangle_{F} + \frac{L\gamma^{2}}{2}\mathbb{E}\left[\left\|G^{t}\right\|_{F}^{2} \mid W^{t}, H^{t}\right] \\ = f(W^{t}) - \gamma p\left\langle\nabla f(W^{t}), H_{B}^{t}\nabla f(W^{t})\right\rangle_{F} - \gamma(1-p)\left\langle\nabla f(W^{t}), \nabla f(W^{t})H_{A}^{t}\right\rangle_{F} \\ + \frac{L\gamma^{2}}{2}p\left\|H_{B}^{t}\nabla f(W^{t})\right\|_{F}^{2} + \frac{L\gamma^{2}}{2}(1-p)\left\|\nabla f(W^{t})H_{A}^{t}\right\|_{F}^{2} \\ + \frac{(34).(36)}{2} \int f(W^{t}) - \frac{\gamma}{2}\left(p\left\langle\nabla f(W^{t}), H_{B}^{t}\nabla f(W^{t})\right\rangle_{F} + (1-p)\left\langle\nabla f(W^{t}), \nabla f(W^{t})H_{A}^{t}\right\rangle_{F}\right).$$

$$(37)$$

For the first term:

$$-\left\langle \nabla f(W^{t}), \mathbb{E}\left[H_{B}^{t}\right] \nabla f(W^{t})\right\rangle_{F} = -\operatorname{Tr}\left(\left(\nabla f(W^{t})\right)^{\top} \mathbb{E}\left[H_{B}^{t}\right] \nabla f(W^{t})\right)$$

$$\leq -\lambda_{\min}\left(\mathbb{E}\left[H_{B}^{t}\right]\right) \operatorname{Tr}\left(\left(\nabla f(W^{t})\right)^{\top} \nabla f(W^{t})\right)$$

$$= -\lambda_{\min}^{H_{B}} \left\|\nabla f(W^{t})\right\|_{F}^{2}. \tag{38}$$

Similarly, for the second term:

$$-\left\langle \nabla f(W^{t}), \nabla f(W^{t}) \mathbb{E}\left[H_{A}^{t}\right]\right\rangle_{F} = -\operatorname{Tr}\left(\left(\nabla f(W^{t})\right)^{\top} \nabla f(W^{t}) \mathbb{E}\left[H_{A}^{t}\right]\right)$$

$$= -\operatorname{Tr}\left(\mathbb{E}\left[H_{A}^{t}\right] \left(\nabla f(W^{t})\right)^{\top} \nabla f(W^{t})\right)$$

$$\leq -\lambda_{\min}^{H_{A}} \left\|\nabla f(W^{t})\right\|_{F}^{2}. \tag{39}$$

Therefore:

$$\mathbb{E}\left[f(W^{t+1}) \mid W^{t}\right] = \mathbb{E}\left[\mathbb{E}\left[f(W^{t+1}) \mid W^{t}, H^{t}\right] \mid W^{t}\right] \\
\leq f(W^{t}) - \frac{\gamma}{2} \left(p \left\langle \nabla f(W^{t}), \mathbb{E}\left[H_{B}^{t}\right] \nabla f(W^{t})\right\rangle_{F} + (1-p) \left\langle \nabla f(W^{t}), \nabla f(W^{t}) \mathbb{E}\left[H_{A}^{t}\right]\right\rangle_{F}\right) \\
\leq f(W^{t}) - \frac{\gamma}{2} \left(p \lambda_{\min}^{H_{B}} + (1-p) \lambda_{\min}^{H_{A}}\right) \left\|\nabla f(W^{t})\right\|_{F}^{2} \\
= f(W^{t}) - \frac{\gamma}{2} \lambda_{\min}^{p} \left\|\nabla f(W^{t})\right\|_{F}^{2}, \tag{40}$$

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

$$\mathbb{E}\left[\mathbb{E}\left[f(W^{t+1})\mid W^{t}, H^{t}\right]\mid W^{t}\right] - f^{\star} \leq f(W^{t}) - f^{\star} - \frac{\gamma}{2}\lambda_{\min}^{p}\left\|\nabla f(W^{t})\right\|_{F}^{2}.$$
(41)

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

$$\mathbb{E}\left[f(W^T) - f^{\star}\right] \leq f(W^0) - f^{\star} - \frac{\gamma}{2}\lambda_{\min}^p \sum_{t=0}^{T-1} \mathbb{E}\left[\left\|\nabla f(W^t)\right\|_{\mathrm{F}}^2\right]. \tag{42}$$

By rearranging terms, we get:

$$\frac{\gamma}{2} \lambda_{\min}^p \sum_{t=0}^{T-1} \mathbb{E}\left[\left\| \nabla f(W^t) \right\|_{\mathcal{F}}^2 \right] \le f(W^0) - f^*. \tag{43}$$

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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \le \frac{2(f(W^0) - f^\star)}{\gamma \lambda_{\min}^p T},\tag{44}$$

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

H.1.2 Convergence under Polyak-Łojasiewicz Condition

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

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

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:

$$\mathbb{E}\left[f(W^{t+1}) \mid W^t\right] \le f(W^t) - \frac{\gamma}{2} \lambda_{\min}^p \left\|\nabla f(W^t)\right\|_{\mathrm{F}}^2. \tag{45}$$

By invoking the Polyak-Łojasiewicz condition (Assumption 6), which states that $\frac{1}{2} \|\nabla f(W)\|_{\mathrm{F}}^2 \ge \mu(f(W) - f^*)$, we can further bound the right-hand side of the inequality (45):

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

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

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

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

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

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

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

This completes the proof.

H.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: $\hat{A}^{t} = \arg\min_{A} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \partial f(W^{t}), B_{S}^{t} A \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma_{t}r^{2}} \left\| B_{S}^{t} A \right\|_{F}^{2} \right\}$

7: $W^{t+1} = W^t + \frac{\alpha}{\pi} B_S^t \hat{A}^t$

8: else

9: Sample $A_S^t \sim \mathcal{D}_S^A$

Right sketch

10: $\hat{B}^{t} = \arg\min_{B} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \partial f(W^{t}), BA_{S}^{t} \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma_{t}r^{2}} \left\| BA_{S}^{t} \right\|_{F}^{2} \right\}$

11: $W^{t+1} = W^t + \frac{\alpha}{2} \hat{B}^t A_S^t$

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

$$|f(W) - f(V)| \le L_0 \|W - V\|_{F}, \quad \forall W, V \in \mathbb{R}^{m \times n}.$$
 (46)

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:

$$\|\partial f(W)\|_{\mathcal{F}} \le L_0, \quad \forall W \in \mathbb{R}^{m \times n}.$$
 (47)

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:

$$\hat{A}^{t} := \arg\min_{A} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \partial f\left(W^{t}\right), B_{S}^{t} A \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma_{t}r^{2}} \left\| B_{S}^{t} A \right\|_{F}^{2} \right\},
\hat{B}^{t} := \arg\min_{B} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \partial f\left(W^{t}\right), BA_{S}^{t} \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma_{t}r^{2}} \left\| BA_{S}^{t} \right\|_{F}^{2} \right\}.$$
(48)

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

$$W^{t+1} = W^t - \gamma_t G^t, \tag{49}$$

where G^t is defined by

$$G^{t} = \begin{cases} H_{B}^{t} \partial f\left(W^{t}\right), & \text{with probability } p \\ \partial f\left(W^{t}\right) H_{A}^{t}, & \text{with probability } 1 - p \end{cases}$$

$$(50)$$

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

$$H_A^t := \left(A_S^t\right)^\top \left(A_S^t \left(A_S^t\right)^\top\right)^\dagger A_S^t \quad \textit{and} \quad H_B^t := B_S^t \left(\left(B_S^t\right)^\top B_S^t\right)^\dagger \left(B_S^t\right)^\top, \tag{51}$$

where † 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:

$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t. \tag{52}$$

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

$$\hat{A}^{t} := \underset{A}{\operatorname{arg\,min}} \left\{ f(W^{t}) + \frac{\alpha}{r} \left\langle \partial f\left(W^{t}\right), B_{S}^{t} A \right\rangle_{F} + \frac{\alpha^{2}}{2\gamma_{t} r^{2}} \left\| B_{S}^{t} A \right\|_{F}^{2} \right\}. \tag{53}$$

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

$$\frac{\alpha}{r} \left(B_S^t \right)^\top \partial f \left(W^t \right) + \frac{\alpha^2}{\gamma_t r^2} \left(B_S^t \right)^\top B_S^t \hat{A}^t = 0;$$

$$\hat{A}^t = -\frac{\gamma_t r}{\alpha} \left(\left(B_S^t \right)^\top B_S^t \right)^\dagger \left(B_S^t \right)^\top \partial f \left(W^t \right). \tag{54}$$

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

$$W^{t+1} = W^{t} + \frac{\alpha}{r} B_{S}^{t} \hat{A}^{t}$$

$$= W^{t} - \gamma_{t} B_{S}^{t} \left(\left(B_{S}^{t} \right)^{\top} B_{S}^{t} \right)^{\dagger} \left(B_{S}^{t} \right)^{\top} \partial f \left(W^{t} \right)$$

$$= W^{t} - \gamma_{t} H_{B}^{t} \partial f \left(W^{t} \right). \tag{55}$$

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

$$W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t. \tag{56}$$

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

$$\hat{B}^{t} = -\frac{\gamma_{t}r}{\alpha}\partial f\left(W^{t}\right)\left(A_{S}^{t}\right)^{\top}\left(A_{S}^{t}\left(A_{S}^{t}\right)^{\top}\right)^{\dagger},\tag{57}$$

which leads to the Right sketch update:

$$W^{t+1} = W^{t} + \frac{\alpha}{r} \hat{B}^{t} A_{S}^{t}$$

$$= W^{t} - \gamma_{t} \partial f \left(W^{t} \right) \left(A_{S}^{t} \right)^{\top} \left(A_{S}^{t} \left(A_{S}^{t} \right)^{\top} \right)^{\dagger} A_{S}^{t}$$

$$= W^{t} - \gamma_{t} \partial f \left(W^{t} \right) H_{A}^{t}. \tag{58}$$

Combined Update Rule. By combining equations (55) and (58), we arrive at the unified update rule:

$$W^{t+1} = W^t - \gamma_t G^t, \tag{59}$$

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

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

$$\mathbb{E}[H] = \alpha I,\tag{60}$$

where a constant $\alpha > 0$.

Theorem 10. Let Assumptions 1, 7, 8, 9, and 10 hold. Let us define the following quantities: $\overline{W}^T := \frac{1}{T} \sum_{t=0}^{T-1} W^t$ 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

$$\mathbb{E}\left[f(\overline{W}^T) - f(W^*)\right] \le \frac{R_0^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2}.$$
(61)

Moreover, with the optimal stepsize $\gamma_* = \sqrt{\frac{(R^0)^2}{T\alpha L_0^2}}$, we obtain:

$$\mathbb{E}\left[f(\overline{W}^T) - f(W^*)\right] \le \frac{R^0 L_0}{\sqrt{\rho T}}.$$
(62)

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

$$\gamma_t = \frac{(f(W^t) - f(W^*))}{\|\partial f(W^t)\|_{\mathcal{P}}^2},\tag{63}$$

then

$$\mathbb{E}\left[f(\overline{W}^T) - f(W^*)\right] \le \frac{R^0 L_0}{\sqrt{\rho T}}.$$
(64)

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

$$W^{t+1} = W^t - \gamma_t G^t, \tag{65}$$

where G^t takes the form

$$G^{t} = \begin{cases} H_{B}^{t} \partial f(W^{t}), & \text{with probability } p \\ \partial f(W^{t}) H_{A}^{t}, & \text{with probability } 1 - p \end{cases}$$

$$(66)$$

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 :

$$\mathbb{E}\left[G^t \mid W^t, H^t\right] = pH_B^t \partial f(W^t) + (1-p)\partial f(W^t)H_A^t, \tag{67}$$

$$\mathbb{E}\left[\left\|G^{t}\right\|_{F}^{2} \mid W^{t}, H^{t}\right] = p \left\|H_{B}^{t} \partial f(W^{t})\right\|_{F}^{2} + (1-p) \left\|\partial f(W^{t}) H_{A}^{t}\right\|_{F}^{2}, \tag{68}$$

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

By the definition of subgradient, we have:

$$f(W^*) \geq f(W^t) + \langle \partial f(W^t), W^* - W^t \rangle_{\mathcal{F}}, \tag{69}$$

which implies:

$$\left\langle \partial f(W^t), W^t - W^* \right\rangle_{\mathcal{F}} \ge f(W^t) - f(W^*).$$
 (70)

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

$$-2\gamma_{t}\mathbb{E}\left[\left\langle G^{t}, W^{t} - W^{*}\right\rangle_{F} \mid W^{t}, H^{t}\right] \stackrel{(67)}{=} -2\gamma_{t}p\left\langle H_{B}^{t}\partial f(W^{t}), W^{t} - W^{*}\right\rangle_{F} -2\gamma_{t}(1-p)\left\langle \partial f(W^{t})H_{A}^{t}, W^{t} - W^{*}\right\rangle_{F}. \quad (71)$$

For projection matrices, we have the following properties:

$$\begin{aligned} \left\| \partial f(W^{t}) H_{A}^{t} \right\|_{F}^{2} &= \left\langle \partial f(W^{t}) H_{A}^{t}, \partial f(W^{t}) H_{A}^{t} \right\rangle_{F} \\ &= \operatorname{Tr} \left(\left(H_{A}^{t} \right)^{\top} \left(\partial f(W^{t}) \right)^{\top} \partial f(W^{t}) H_{A}^{t} \right) \\ &= \operatorname{Tr} \left(\left(\nabla f(W^{t}) \right)^{\top} \nabla f(W^{t}) H_{A}^{t} \left(H_{A}^{t} \right)^{\top} \right) \\ &= \operatorname{Tr} \left(\left(\partial f(W^{t}) \right)^{\top} \partial f(W^{t}) H_{A}^{t} \right) \\ &= \left\langle \partial f(W^{t}), \partial f(W^{t}) H_{A}^{t} \right\rangle_{F}, \end{aligned}$$
(72)

and similarly, one can show that

$$\left\| H_B^t \partial f(W^t) \right\|_{\mathcal{F}}^2 = \left\langle \partial f(W^t), H_B^t \partial f(W^t) \right\rangle_{\mathcal{F}}. \tag{73}$$

This allows us to express the second moment term as:

$$\gamma_{t}^{2}\mathbb{E}\left[\left\|G^{t}\right\|_{\mathrm{F}}^{2}\mid W^{t}, H^{t}\right] \stackrel{(68)}{=} \gamma_{t}^{2}p\left\|H_{B}^{t}\partial f(W^{t})\right\|_{\mathrm{F}}^{2} + \gamma_{t}^{2}(1-p)\left\|\partial f(W^{t})H_{A}^{t}\right\|_{\mathrm{F}}^{2}$$

$$\stackrel{(72),(73)}{=} \gamma_{t}^{2}p\left\langle\partial f(W^{t}), H_{B}^{t}\partial f(W^{t})\right\rangle_{\mathrm{F}} + \gamma_{t}^{2}(1-p)\left\langle\partial f(W^{t}), \partial f(W^{t})H_{A}^{t}\right\rangle_{\mathrm{F}}.$$

$$(74)$$

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

$$\mathbb{E}\left[\|W^{t+1} - W^*\|_{F}^{2} \mid W^{t}, H^{t}\right] = \mathbb{E}\left[\|W^{t} - \gamma_{t}G^{t} - W^*\|_{F}^{2} \mid W^{t}, H^{t}\right] \\
= \|W^{t} - W^*\|_{F}^{2} - 2\gamma_{t}\mathbb{E}\left[\left\langle G^{t}, W^{t} - W^*\right\rangle_{F} \mid W^{t}, H^{t}\right] \\
+ \gamma_{t}^{2}\mathbb{E}\left[\|G^{t}\|_{F}^{2} \mid W^{t}, H^{t}\right] \\
\stackrel{(71)}{=} \|W^{t} - W^*\|_{F}^{2} - 2\gamma_{t}p\left\langle H_{B}^{t}\partial f(W^{t}), W^{t} - W^*\right\rangle_{F} \\
- 2\gamma_{t}(1 - p)\left\langle \partial f(W^{t})H_{A}^{t}, W^{t} - W^*\right\rangle_{F} + \gamma_{t}^{2}p\left\langle \partial f(W^{t}), H_{B}^{t}\partial f(W^{t})\right\rangle_{F} \\
+ \gamma_{t}^{2}(1 - p)\left\langle \partial f(W^{t}), \partial f(W^{t})H_{A}^{t}\right\rangle_{F}. \tag{75}$$

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

$$\langle \partial f(W^{t}), \mathbb{E} \left[H_{B}^{t} \right] \partial f(W^{t}) \rangle_{F} = \operatorname{Tr} \left(\left(\partial f(W^{t}) \right)^{\top} \mathbb{E} \left[H_{B}^{t} \right] \partial f(W^{t}) \right)$$

$$= \alpha \operatorname{Tr} \left(\left(\partial f(W^{t}) \right)^{\top} \partial f(W^{t}) \right)$$

$$= \alpha \left\| \partial f(W^{t}) \right\|_{F}^{2}, \tag{76}$$

and similarly,

$$\left\langle \partial f(W^t), \partial f(W^t) \mathbb{E}\left[H_A^t\right] \right\rangle_{\mathcal{F}} = \alpha \left\| \partial f(W^t) \right\|_{\mathcal{F}}^2. \tag{77}$$

Taking expectation of both sides of (75) again, we get

$$\mathbb{E}\left[\|W^{t+1} - W^*\|_{F}^{2} \mid W^{t}\right] = \mathbb{E}\left[\mathbb{E}\left[\|W^{t+1} - W^*\|_{F}^{2} \mid W^{t}, H^{t}\right] \mid W^{t}\right]$$
(78)
$$= \|W^{t} - W^*\|_{F}^{2} - 2\gamma_{t}p \left\langle \mathbb{E}\left[H_{B}^{t}\right] \partial f\left(W^{t}\right), W^{t} - W^*\right\rangle_{F}$$
(79)
$$-2\gamma_{t}(1 - p) \left\langle \partial f\left(W^{t}\right) \mathbb{E}\left[H_{A}^{t}\right], W^{t} - W^*\right\rangle_{F}$$
(79)
$$+\gamma_{t}^{2}p \left\langle \partial f\left(W^{t}\right), \mathbb{E}\left[H_{B}^{t}\right] \partial f\left(W^{t}\right)\right\rangle_{F} + \gamma_{t}^{2}(1 - p) \left\langle \partial f\left(W^{t}\right), \partial f\left(W^{t}\right) \mathbb{E}\left[H_{A}^{t}\right]\right\rangle_{F}$$
(80)
$$-2\gamma_{t}(1 - p)\alpha \left\langle \partial f\left(W^{t}\right), W^{t} - W^*\right\rangle_{F} + \gamma_{t}^{2}\alpha \left\|\partial f\left(W^{t}\right)\right\|_{F}^{2}$$

$$= \|W^{t} - W^{*}\|_{F}^{2} - 2\gamma_{t}\alpha \left\langle \partial f\left(W^{t}\right), W^{t} - W^{*}\right\rangle_{F} + \gamma_{t}^{2}\alpha \left\|\partial f\left(W^{t}\right)\right\|_{F}^{2}$$

$$= \|W^{t} - W^{*}\|_{F}^{2} - 2\gamma_{t}\alpha \left\langle \partial f\left(W^{t}\right), W^{t} - W^{*}\right\rangle_{F} + \gamma_{t}^{2}\alpha \left\|\partial f\left(W^{t}\right)\right\|_{F}^{2}$$
(81)

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

$$\|\partial f(W)\|_{\mathcal{F}} \le L_0 \quad \forall W \in \mathbb{R}^{m \times n}.$$
 (82)

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 (78) again, applying tower property (11) and using the bound (82), we obtain:

$$\mathbb{E}\left[\left\|W^{t+1} - W^*\right\|_{\mathrm{F}}^2\right] \le \mathbb{E}\left[\left\|W^t - W^*\right\|_{\mathrm{F}}^2\right] - 2\gamma\alpha\mathbb{E}\left[f(W^t) - f(W^*)\right] + \gamma^2\alpha L_0^2. \tag{83}$$

Rearranging terms in (83):

$$2\gamma \alpha \mathbb{E}\left[f(W^{t}) - f(W^{*})\right] \le \mathbb{E}\left[\left\|W^{t} - W^{*}\right\|_{F}^{2}\right] - \mathbb{E}\left[\left\|W^{t+1} - W^{*}\right\|_{F}^{2}\right] + \gamma^{2} \alpha L_{0}^{2}.$$
 (84)

Summing inequality (84) for t = 0, ..., T - 1:

$$2\gamma\alpha \sum_{t=0}^{T-1} \mathbb{E}\left[f(W^{t}) - f(W^{*})\right] \leq \sum_{t=0}^{T-1} \left(\mathbb{E}\left[\left\|W^{t} - W^{*}\right\|_{F}^{2}\right] - \mathbb{E}\left[\left\|W^{t+1} - W^{*}\right\|_{F}^{2}\right]\right) + T\gamma^{2}\alpha L_{0}^{2}$$

$$= \mathbb{E}\left[\left\|W^{0} - W^{*}\right\|_{F}^{2}\right] - \mathbb{E}\left[\left\|W^{T} - W^{*}\right\|_{F}^{2}\right] + T\gamma^{2}\alpha L_{0}^{2}$$

$$\leq \left\|W^{0} - W^{*}\right\|_{F}^{2} + T\gamma^{2}\alpha L_{0}^{2}, \tag{85}$$

where the last inequality follows from the non-negativity of $\left\|W^T - W^*\right\|_{\mathrm{F}}^2$.

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

$$\mathbb{E}\left[f(\overline{W}^{T}) - f(W^{*})\right] \leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[f(W^{t}) - f(W^{*})\right]$$

$$\leq \frac{\|W^{0} - W^{*}\|_{F}^{2}}{2\gamma\alpha T} + \frac{\gamma L_{0}^{2}}{2}$$

$$= \frac{(R^{0})^{2}}{2\gamma\alpha T} + \frac{\gamma L_{0}^{2}}{2},$$
(86)

1728 where we denoted $(R^0)^2 := \|W^0 - W^*\|_{\mathrm{F}}^2$

 To optimize this bound, we minimize it with respect to γ . The optimal stepsize γ_* solves:

$$\gamma_* = \underset{\gamma>0}{\operatorname{arg\,min}} \left(\frac{(R^0)^2}{2\gamma\alpha T} + \frac{\gamma L_0^2}{2} \right)$$
$$= \sqrt{\frac{(R^0)^2}{T\alpha L_0^2}}.$$
 (87)

Substituting γ_* back into (86), we obtain the optimal convergence rate:

$$\mathbb{E}\left[f(\overline{W}^T) - f(W^*)\right] \le \frac{R^0 L_0}{\sqrt{\rho T}}.$$
(88)

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

$$\gamma_{t} = \underset{\gamma>0}{\operatorname{arg\,min}} \left\{ \left\| W^{t} - W^{*} \right\|_{F}^{2} - 2\gamma\alpha \left(f(W^{t}) - f(W^{*}) \right) + \gamma^{2}\alpha \left\| \partial f \left(W^{t} \right) \right\|_{F}^{2} \right\} \\
= \frac{\left(f(W^{t}) - f(W^{*}) \right)}{\left\| \partial f(W^{t}) \right\|_{F}^{2}}.$$
(89)

Substituting this stepsize into inequality (78):

$$\mathbb{E}\left[\left\|W^{t+1} - W^{*}\right\|_{F}^{2} \mid W^{t}\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|W^{t+1} - W^{*}\right\|_{F}^{2} \mid W^{t}, H^{t}\right] \mid W^{t}\right] \\ \leq \left\|W^{t} - W^{*}\right\|_{F}^{2} - 2\gamma_{t}\alpha\left(f(W^{t}) - f(W^{*})\right) + \gamma_{t}^{2}\alpha\left\|\partial f\left(W^{t}\right)\right\|_{F}^{2} \\ \stackrel{(89)}{=} \left\|W^{t} - W^{*}\right\|_{F}^{2} - \frac{\alpha\left(f(W^{t}) - f(W^{*})\right)^{2}}{\left\|\partial f(W^{t})\right\|_{F}^{2}} \\ \stackrel{(82)}{\leq} \left\|W^{t} - W^{*}\right\|_{F}^{2} - \frac{\alpha\left(f(W^{t}) - f(W^{*})\right)^{2}}{L_{2}^{2}}. \tag{90}$$

Taking expectation of both sides of (90) again and applying the tower property

$$\mathbb{E}\left[\left\|W^{t+1} - W^*\right\|_{\mathrm{F}}^2\right] \le \mathbb{E}\left[\left\|W^t - W^*\right\|_{\mathrm{F}}^2\right] - \frac{\alpha \mathbb{E}\left[\left(f(W^t) - f(W^*)\right)^2\right]}{L_0^2}$$
(91)

Since f is convex, by Jensen's inequality (14) and the Cauchy-Bunyakovsky-Schwarz inequality (12) with $X := f(W^t) - f(W^*)$ and Y := 1, we have

$$\mathbb{E}\left[f_{i}(\overline{W}^{T}) - f(W^{*})\right] \stackrel{(14)}{\leq} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1}f(W^{t}) - f(W^{*})\right]$$

$$\leq \frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\left[f(W^{t}) - f(W^{*})\right]$$

$$\stackrel{(12)}{\leq} \frac{1}{T}\sum_{t=0}^{T-1}\sqrt{\mathbb{E}\left[\left(f(W^{t}) - f(W^{*})\right)^{2}\right]}$$

$$\leq \sqrt{\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\left[\left(f(W^{t}) - f(W^{*})\right)^{2}\right]}$$

$$\stackrel{(91)}{\leq} \frac{R^{0}L_{0}}{\sqrt{\sigma T}}, \qquad (92)$$

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

H.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
1786
                    length T, sketch distribution \mathcal{D}_S^B or \mathcal{D}_S^A, Bernoulli probability p
1787
                    for t = 0, 1, \dots, T - 1 do
1788
                         Sample c^t \sim \text{Be}(p)
                                                                                                                                   Bernoulli random variable
1789
                         if c^t = 1 then
               4:
1790
                            Sample B_S^t \sim \mathcal{D}_S^B
               5:
                                                                                                                                                            Left sketch
1791
                            \hat{A}^t = -\eta \left( \left( B_S^t \right)^\top B_S^t \right)^\dagger \left( B_S^t \right)^\top g(W^t)
1792
1793
                            W^{t+1} = \overset{\searrow}{W}^t + \frac{\alpha}{\pi} B_S^t \overset{?}{A}^t
1794
               8:
                            Sample A_S^t \sim \mathcal{D}_S^A
               9:
                                                                                                                                                          Right sketch
1796
                            \hat{B}^{t} = -\eta g(W^{t}) \left( A_{S}^{t} \right)^{\top} \left( A_{S}^{t} \left( A_{S}^{t} \right)^{\top} \right)^{\dagger}
              10:
1797
                             W^{t+1} = W^t + \frac{\alpha}{2} \hat{B}^t A_S^t
1798
1799
                        end if
              12:
              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 & 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

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}, \tag{93}$$

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

H.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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \leq \frac{6(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T} + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p},$$

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 $\{\frac{w_t}{W_{T-1}}\}_{t=0}^{T-1}$, where $w_t = \frac{w_{t-1}}{(1+\gamma^2 L A_1 \lambda_{\max}^p)}$, $W_{T-1} = \sum_{t=0}^{T-1} w_t$, and $w^{-1} > 0$.

Proof. We start with smoothness of function f:

$$f(W^{t+1}) \leq f(W^t) + \langle \nabla f(W^t), W^{t+1} - W^t \rangle + \frac{L}{2} \|W^{t+1} - W^t\|_{\mathrm{F}}^2$$

$$\stackrel{(93)}{=} f(W^t) - \gamma \langle \nabla f(W^t), \hat{G}^t \rangle + \frac{\gamma^2 L}{2} \|\hat{G}^t\|_{\mathrm{F}}^2. \tag{94}$$

Taking a conditional expectation by W^t , we bound the second and the third terms from inequality (94):

$$\mathbb{E}\left[\langle \nabla f(W^{t}), \hat{G}^{t} \rangle | W^{t}\right] = \langle \nabla f(W^{t}), \mathbb{E}\left[\hat{G}^{t} | W^{t}\right] \rangle$$

$$\stackrel{(93)}{=} p\langle \nabla f(W^{t}), \mathbb{E}\left[H_{B}^{t} G^{t} | W^{t}\right] \rangle + (1-p)\langle \nabla f(W^{t}), \mathbb{E}\left[G^{t} H_{A}^{t} | W^{t}\right] \rangle$$

$$\stackrel{(*)}{=} p\langle \nabla f(W^{t}), \mathbb{E}\left[H_{B}^{t} | W^{t}\right] \mathbb{E}\left[G^{t} | W^{t}\right] \rangle + (1-p)\langle \nabla f(W^{t}), \mathbb{E}\left[G^{t} | W^{t}\right] \mathbb{E}\left[H_{A}^{t} | W^{t}\right] \rangle$$

$$= p\langle \nabla f(W^{t}), \mathbb{E}\left[H_{B}^{t} | W^{t}\right] \nabla f(W^{t}) \rangle + (1-p)\langle \nabla f(W^{t}), \nabla f(W^{t}) \mathbb{E}\left[H_{A}^{t} | W^{t}\right] \rangle$$

$$\geq \underbrace{\left(p\lambda_{\min}(\mathbb{E}\left[H_{B}^{t}\right]) + (1-p)\lambda_{\min}(\mathbb{E}\left[H_{A}^{t}\right])\right)}_{:=\lambda_{\min}^{p}} \|\nabla f(W^{t})\|_{F}^{2}$$

$$= \lambda_{\min}^{p} \|\nabla f(W^{t})\|_{F}^{2}, \tag{95}$$

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

$$\mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] \stackrel{(93)}{=} p\mathbb{E}\left[\left\|H_{B}^{t}G^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] + (1-p)\mathbb{E}\left[\left\|G^{t}H_{A}^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right]$$

$$= p\mathbb{E}\left[\left\langle H_{B}^{t}G^{t}, H_{B}^{t}G^{t}\right\rangle|W^{t}\right] + (1-p)\mathbb{E}\left[\left\langle G^{t}H_{A}^{t}, G^{t}H_{A}^{t}\right\rangle|W^{t}\right]$$

$$\stackrel{(**)}{=} p\mathbb{E}\left[\left\langle G^{t}, H_{B}^{t}G^{t}\right\rangle|W^{t}\right] + (1-p)\mathbb{E}\left[\left\langle G^{t}, G^{t}H_{A}^{t}\right\rangle|W^{t}\right],$$

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

$$\mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] = p\mathbb{E}\left[\left\langle G^{t}, \mathbb{E}\left[H_{B}^{t}|W^{t}\right]G^{t}\right\rangle|W^{t}\right] + (1-p)\mathbb{E}\left[\left\langle G^{t}, G^{t}\mathbb{E}\left[H_{A}^{t}|W^{t}\right]\right\rangle|W^{t}\right] \\
\leq p\lambda_{\max}(\mathbb{E}\left[H_{B}^{t}|W^{t}\right])\mathbb{E}\left[\left\|G^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] + (1-p)\lambda_{\max}(\mathbb{E}\left[H_{A}^{t}|W^{t}\right])\mathbb{E}\left[\left\|G^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] \\
= \underbrace{\left(p\lambda_{\max}(\mathbb{E}\left[H_{B}^{t}|W^{t}\right]) + (1-p)\lambda_{\max}(\mathbb{E}\left[H_{A}^{t}|W^{t}\right])\right)}_{:=\lambda_{\max}^{p}}\mathbb{E}\left[\left\|G^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right] \\
= \lambda_{\max}^{p}\mathbb{E}\left[\left\|G^{t}\right\|_{\mathrm{F}}^{2}|W^{t}\right]. \tag{96}$$

Plugging (95) and (96) into (94), we obtain

$$\mathbb{E}\left[f(W^{t+1})|W^{t}\right] \leq f(W^{t}) - \gamma \mathbb{E}\left[\langle \nabla f(W^{t}), \hat{G}^{t} \rangle |W^{t}\right] + \frac{\gamma^{2}L}{2} \mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{F}^{2} |W^{t}\right]$$

$$\leq f(W^{t}) - \gamma \lambda_{\min}^{p} \left\|\nabla f(W^{t})\right\|_{F}^{2} + \frac{\gamma^{2}\lambda_{\max}^{p}L}{2} \mathbb{E}\left[\left\|G^{t}\right\|_{F}^{2} |W^{t}\right].$$

By Assumption 4,

$$\mathbb{E}\left[f(W^{t+1}) - f^{*}|W^{t}\right] \leq f(W^{t}) - \gamma \mathbb{E}\left[\langle \nabla f(W^{t}), \hat{G}^{t} \rangle |W^{t}\right] + \frac{\gamma^{2}L}{2} \mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{F}^{2} |W^{t}\right] \\
\leq f(W^{t}) - f^{*} - \gamma \lambda_{\min}^{p} \left\|\nabla f(W^{t})\right\|_{F}^{2} \\
+ \frac{\gamma^{2}\lambda_{\max}^{p}L}{2} \left(2A_{1}(f(W^{t}) - f^{*}) + B_{1} \left\|\nabla f(W^{t})\right\|_{F}^{2} + C_{1}\right) \\
\leq \left(1 + \gamma^{2}\lambda_{\max}^{p}LA_{1}\right) \left(f(W^{t}) - f^{*}\right) - \gamma \lambda_{\min}^{p} \left(1 - \frac{\gamma LB_{1}\lambda_{\max}^{p}}{2\lambda_{\min}^{p}}\right) \left\|\nabla f(W^{t})\right\|_{F}^{2} \\
+ \frac{\gamma^{2}\lambda_{\max}^{p}LC_{1}}{2}.$$

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

1892
1893
$$\mathbb{E}\left[f(W^{t+1}) - f^*\right] \leq \left(1 + \gamma^2 \lambda_{\max}^p L A_1\right) \mathbb{E}\left[f(W^t) - f^*\right] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E}\left[\left\|\nabla f(W^t)\right\|_{\mathrm{F}}^2\right] + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}. \tag{97}$$

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

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

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

$$\frac{1}{2}\lambda_{\min}^{p}w_{t}r^{t} \leq \frac{w_{t}}{\gamma}\left(1+\gamma^{2}\lambda_{\max}^{p}LA_{1}\right)\delta^{t}-\frac{w_{t}}{\gamma}\delta^{t+1}+\frac{1}{2}\gamma LC_{1}\lambda_{\max}^{p}w_{t}$$

$$=\frac{w_{t-1}\delta^{t}}{\gamma}-\frac{w_{t}\delta^{t+1}}{\gamma}+\frac{1}{2}\gamma LC_{1}\lambda_{\max}^{p}w_{t}.$$

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 $W_{T-1} = \sum_{t=0}^{T-1} w_t$, 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

$$W_{T-1} = \sum_{t=0}^{T-1} w_t \ge T \min_{0 \le t \le T-1} w_t = Tw_{T-1} = \frac{Tw_{-1}}{(1 + \gamma^2 \lambda_{\max}^p LA_1)^T},$$

we have

$$\sum_{t=0}^{T-1} \frac{w_t}{\mathcal{W}^{T-1}} r^t \leq \frac{2(1+\gamma^2 \lambda_{\max}^p LA_1)^T}{\gamma T \lambda_{\min}^p} (f(W^0) - f^*) + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p}.$$

Selecting $0 < \gamma \le \frac{1}{\sqrt{LA_1\lambda_{\max}^p T}}$, and using $(1+\gamma^2\lambda_{\max}^p LA_1)^T \le \exp\left(\gamma^2\lambda_{\max}^p LA_1T\right) \le \exp\left(1\right) \le 3$, we obtain

$$\sum_{t=0}^{T-1} \frac{w_t}{\mathcal{W}^{T-1}} r^t \leq \frac{6\delta^0}{\gamma T \lambda_{\min}^p} + \gamma L C_1 \frac{\lambda_{\max}^p}{\lambda_{\min}^p}$$

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

H.2.2 CONVERGENCE UNDER POLYAK-ŁOJASIEWICZ CONDITION

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

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

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

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 97. Using PL-inequality (see Assumption 6), we have

$$\mathbb{E}\left[f(W^{t+1}) - f^*\right] \leq \left(1 + \gamma^2 \lambda_{\max}^p L A_1\right) \mathbb{E}\left[f(W^t) - f^*\right] - \frac{\gamma \lambda_{\min}^p}{2} \mathbb{E}\left[\left\|\nabla f(W^t)\right\|_F^2\right] + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}$$

$$\leq \left(1 - \gamma \mu \lambda_{\min}^p + \gamma^2 \lambda_{\max}^p L A_1\right) \mathbb{E}\left[f(W^t) - f^*\right] + \frac{\gamma^2 \lambda_{\max}^p L C_1}{2}.$$

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

$$\begin{split} \mathbb{E}\left[f(W^{t+1}) - f^*\right] & \leq \left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)\mathbb{E}\left[f(W^t) - f^*\right] + \frac{\gamma^2\lambda_{\max}^pLC_1}{2} \\ & \leq \left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)^{t+1}\mathbb{E}\left[f(W^0) - f^*\right] + \frac{\gamma^2\lambda_{\max}^pLC_1}{2}\sum_{\tau=0}^t\left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)^{t-\tau} \\ & \leq \left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)^{t+1}\mathbb{E}\left[f(W^0) - f^*\right] + \frac{\gamma^2\lambda_{\max}^pLC_1}{2}\sum_{\tau=0}^{\infty}\left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)^{\tau} \\ & = \left(1 - \frac{1}{2}\gamma\mu\lambda_{\min}^p\right)^{t+1}\mathbb{E}\left[f(W^0) - f^*\right] + \frac{\gamma^2\lambda_{\max}^pLC_1}{\gamma\mu\lambda_{\min}^p}, \end{split}$$

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

H.3 ANALYSIS OF BERNOULLI-LORA-MVR

Algorithm 5 Bernoulli-LoRA-MVR

1998

1999 2000

20012002

2003

2004

2005

2006

2007

2008

2009 2010

2011

2012

2013

2014 2015

2016

201720182019

202020212022

2023

2024

2025

2026

2027

2028

2029

2030

2031

2032

2033

2035

20362037

2038

2039

2040

20412042

2043 2044

2045

2046

2047 2048

20492050

2051

```
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
           Sample c^t \sim \text{Be}(p)
                                                                                                                          Bernoulli random variable
           if c^t = 1 then
 4:
              Sample B_S^t \sim \mathcal{D}_S^B
                                                                                                                                                   Left sketch
              \hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger (B_S^t)^\top G^t
W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t
 7:
 8:
              Sample A_S^t \sim \mathcal{D}_S^A
                                                                                                                                                 Right sketch
              \hat{B}^{t} = -\eta G^{t} \left( A_{S}^{t} \right)^{\top} \left( A_{S}^{t} \left( A_{S}^{t} \right)^{\top} \right)^{\dagger}
10:
              W^{t+1} = W^t + \frac{\alpha}{\pi} \hat{B}^t A_S^t
11:
12:
           end if
          Sample \xi^{t+1} \sim \mathcal{D}

G^{t+1} = \nabla f_{\xi^{t+1}}(W^{t+1}) + (1-b) \left( G^t - \nabla f_{\xi^{t+1}}(W^t) \right)
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 & 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 & 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

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}$$
 (98)

$$G^{t+1} = \nabla f_{\xi^{t+1}}(W^{t+1}) + (1-b)\left(G^t - \nabla f_{\xi^{t+1}}(W^t)\right). \tag{99}$$

First of all, we reprove descent lemma from the paper of Li et al. (2021) for generic gradient step (98).

Lemma 5. Let Assumptions 1, 3 hold. Then, iterates defined as (98) satisfy

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

Proof. By Assumption 3, we have

$$f(W^{t+1}) \leq f(W^{t}) + \langle \nabla f(W^{t}), W^{t+1} - W^{t} \rangle_{F} + \frac{L}{2} \|W^{t+1} - W^{t}\|_{F}^{2}$$

$$= f(W^{t}) - \gamma \langle \nabla f(W^{t}), \hat{G}^{t} \rangle_{F} + \frac{L}{2} \|W^{t+1} - W^{t}\|_{F}^{2}.$$
(100)

To continue our proof, we need to bound the second term from (100). Taking conditional expectation by H^t, W^t , we obtain

$$\begin{split} \mathbb{E}\left[\langle \nabla f(W^t), \hat{G}^t \rangle_F \mid H^t, W^t \right] & \stackrel{(98)}{=} \quad p \langle \nabla f(W^t), H_B^t G^t \rangle_F + (1-p) \langle \nabla f(W^t), G^t H_A^t \rangle_F \\ & = \quad p \langle H_B^t \nabla f(W^t), H_B^t G^t \rangle_F + (1-p) \langle \nabla f(W^t) H_A^t, G^t H_A^t \rangle_F \\ & = \quad \frac{p}{2} \left(\left\| H_B^t \nabla f(W^t) \right\|_F^2 + \left\| H_B^t G^t \right\|_F^2 - \left\| H_B^t G^t - H_B^t \nabla f(W^t) \right\|_F^2 \right) \\ & \quad + \frac{1-p}{2} \left(\left\| \nabla f(W^t) H_A^t \right\|_F^2 + \left\| G^t H_A^t \right\|_F^2 - \left\| G^t H_A^t - \nabla f(W^t) H_A^t \right\|_F^2 \right) \\ & \geq \quad \frac{1}{2} \left(p \left\| H_B^t \nabla f(W^t) \right\|_F^2 + (1-p) \left\| \nabla f(W^t) H_A^t \right\|_F^2 \right) + \frac{1}{2} \mathbb{E} \left[\left\| \hat{G}^t \right\|_F^2 \mid H^t, W^t \right] \\ & \quad - \frac{1}{2} \left(p \left\| H_B^t G^t - H_B^t \nabla f(W^t) \right\|_F^2 + (1-p) \left\| G^t H_A^t - \nabla f(W^t) H_A^t \right\|_F^2 \right). \end{split}$$

Taking conditional expectation by W^t , we have

$$\mathbb{E}\left[\langle\nabla f(W^{t}),\hat{G}^{t}\rangle_{F}|W^{t}\right] \geq \frac{1}{2}\left(p\mathbb{E}\left[\left\|H_{B}^{t}\nabla f(W^{t})\right\|_{F}^{2}|W^{t}\right] + (1-p)\mathbb{E}\left[\left\|\nabla f(W^{t})H_{A}^{t}\right\|_{F}^{2}|W^{t}\right]\right) + \frac{1}{2}\mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{F}^{2}|W^{t}\right] \\
-\frac{1}{2}\left(p\mathbb{E}\left[\left\|H_{B}^{t}G^{t} - H_{B}^{t}\nabla f(W^{t})\right\|_{F}^{2}|W^{t}\right] + (1-p)\mathbb{E}\left[\left\|G^{t}H_{A}^{t} - \nabla f(W^{t})H_{A}^{t}\right\|_{F}^{2}|W^{t}\right]\right) \\
\stackrel{(*)}{\geq} \frac{1}{2}\underbrace{\left(p\lambda_{\min}(\mathbb{E}\left[H_{B}^{t}\right]) + (1-p)\lambda_{\min}(\mathbb{E}\left[H_{A}^{t}\right])\right)}_{:=\lambda_{\min}^{p}}\left\|\nabla f(W^{t})\right\|_{F}^{2} + \frac{1}{2}\mathbb{E}\left[\left\|\hat{G}^{t}\right\|_{F}^{2}|W^{t}\right] \\
-\frac{1}{2}\underbrace{\left(p\lambda_{\max}(\mathbb{E}\left[H_{B}^{t}\right]) + (1-p)\lambda_{\max}(\mathbb{E}\left[H_{A}^{t}\right])\right)}_{:=\lambda_{\max}^{p}}\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2} \\
\stackrel{(98)}{=} \frac{\lambda_{\min}^{p}}{2}\left\|\nabla f(W^{t})\right\|_{F}^{2} + \frac{1}{2\gamma^{2}}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}|W^{t}\right] - \frac{\lambda_{\max}^{p}}{2}\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}, \quad (101)$$

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

$$\begin{split} & \mathbb{E}\left[\left\|H_{B}^{t}V\right\|_{\mathrm{F}}^{2}\right] &= \mathbb{E}\left[\left\langle H_{B}^{t}V, H_{B}^{t}V\right\rangle_{F}\right] = \left\langle \mathbb{E}\left[H_{B}^{t}\right]V, V\right\rangle_{F} \geq \lambda_{\min}\left(\mathbb{E}\left[H_{B}^{t}\right]\right)\left\|V\right\|_{\mathrm{F}}^{2}, \\ & \mathbb{E}\left[\left\|H_{B}^{t}V\right\|_{\mathrm{F}}^{2}\right] &\leq \lambda_{\max}\left(\mathbb{E}\left[H_{B}^{t}\right]\right)\left\|V\right\|_{\mathrm{F}}^{2}, \\ & \mathbb{E}\left[\left\|VH_{A}^{t}\right\|_{\mathrm{F}}^{2}\right] &= \mathbb{E}\left[\left\langle VH_{A}^{t}, VH_{A}^{t}\right\rangle_{F}\right] = \left\langle V\mathbb{E}\left[H_{A}^{t}\right], V\right\rangle_{F} \geq \lambda_{\min}\left(\mathbb{E}\left[H_{A}^{t}\right]\right)\left\|V\right\|_{\mathrm{F}}^{2}, \\ & \mathbb{E}\left[\left\|VH_{A}^{t}\right\|_{\mathrm{F}}^{2}\right] &\leq \lambda_{\max}\left(\mathbb{E}\left[H_{A}^{t}\right]\right)\left\|V\right\|_{\mathrm{F}}^{2}. \end{split}$$

Plugging in (101) into (100), we get

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

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

$$\mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2}\right] \leq (1-b)^{2} \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2}\right] + 2(1-b)^{2} L^{2} \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \sigma^{2} \left(102\right) + 2b^{2} \mathcal{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \mathcal{E}\left[\left\|W^{t}\right\|_{\mathrm{F}}^{2}\right] + 2b^{2} \mathcal{E}\left[\left\|W^{t}\right\|_{\mathrm{F}}^{2$$

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

$$\mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \stackrel{(99)}{=} \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1}) + (1-b)\left(G^{t} - \nabla f_{\xi^{t+1}}(W^{t})\right)\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ \stackrel{(13)}{=} (1-b)^{2} \left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2} \\ + \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1}) + (1-b)\left(\nabla f(W^{t}) - \nabla f_{\xi^{t+1}}(W^{t})\right)\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ \leq (1-b)^{2} \left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2} + 2b^{2} \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ + 2(1-b)^{2} \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f_{\xi^{t+1}}(W^{t}) - \nabla f(W^{t+1}) + \nabla f(W^{t})\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ \leq (1-b)^{2} \left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2} + 2b^{2} \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ + 2(1-b)^{2} \mathbb{E}\left[\left\|\nabla f_{\xi^{t+1}}(W^{t+1}) - \nabla f_{\xi^{t+1}}(W^{t})\right\|_{\mathrm{F}}^{2} |\mathcal{F}^{t+1}\right] \\ \leq (1-b)^{2} \left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2} + 2(1-b)^{2} L^{2} \left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2} + 2b^{2}\sigma^{2},$$

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

H.3.1 Convergence for Smooth Non-Convex Functions

Theorem 13. Let Assumptions 1, 2, 3, and 5 hold, and let the stepsize satisfy $0 < \gamma \le \frac{1}{L\left(1+\sqrt{\frac{2\lambda_{\max}^p(1-b)^2}{b}}\right)}$. Then the iterates of Bernoulli-LoRA-MVR (Algorithm 5) satisfy

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^{T})\right\|_{\mathrm{F}}^{2}\right] \leq \frac{2(f(W^{0}) - f^{*})}{\lambda_{\min}^{p} \gamma T} + \frac{\left\|G^{0} - \nabla f(W^{0})\right\|_{\mathrm{F}}^{2}}{b(2 - b)T} \cdot \frac{\lambda_{\max}^{p}}{\lambda_{\min}^{p}} + \frac{2b\sigma^{2}}{2 - b} \cdot \frac{\lambda_{\max}^{p}}{\lambda_{\min}^{p}}, \quad (103)$$

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 Φ_t as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\text{max}}^p}{2b(2-b)} \|G^t - \nabla f(W^t)\|_F^2.$$
 (104)

By Lemma 5 and Lemma 6, we have

$$\mathbb{E}\left[\Phi_{t+1}\right] \leq \mathbb{E}\left[f(W^{t})\right] - f^{*} - \frac{\gamma \lambda_{\min}^{p}}{2} \mathbb{E}\left[\left\|\nabla f(W^{t})\right\|_{F}^{2}\right] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] \\
+ \frac{\gamma \lambda_{\max}^{p}}{2} \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}\right] + \frac{\gamma(1-b)^{2} \lambda_{\max}^{p}}{2b(2-b)} \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}\right] \\
+ \frac{\gamma(1-b)^{2} L^{2} \lambda_{\max}^{p}}{2b(2-b)} \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{2-b} \\
\leq \mathbb{E}\left[\Phi_{t}\right] - \frac{\gamma \lambda_{\min}^{p}}{2} \mathbb{E}\left[\left\|\nabla f(W^{t})\right\|_{F}^{2}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{2-b} \\
- \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-b)^{2} L^{2} \lambda_{\max}^{p}}{2b(2-b)}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right].$$

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

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

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

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

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

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

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

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

H.3.2 CONVERGENCE UNDER POLYAK-ŁOJASIEWICZ CONDITION

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

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

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

$$\mathbb{E}\left[f(W^T) - f^*\right] \le \left(1 - \gamma \mu \lambda_{\min}^p\right)^T \Phi_0 + \frac{b\sigma^2}{(2 - b)\mu} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},\tag{105}$$

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

Proof. Denote Lyapunov function Φ_t as follows

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

By Lemma 5 and Lemma 6, we have

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

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

$$\mathbb{E}\left[\Phi_{t+1}\right] \leq \left(1 - \gamma \mu \lambda_{\min}^{p}\right) \mathbb{E}\left[\Phi_{t}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{2 - b} \\
\leq \left(1 - \gamma \mu \lambda_{\min}^{p}\right)^{t+1} \mathbb{E}\left[\Phi_{0}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{2 - b} \sum_{\tau=0}^{t} \left(1 - \gamma \mu \lambda_{\min}^{p}\right)^{t-\tau} \\
\leq \left(1 - \gamma \mu \lambda_{\min}^{p}\right)^{t+1} \mathbb{E}\left[\Phi_{0}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{2 - b} \sum_{\tau=0}^{\infty} \left(1 - \gamma \mu \lambda_{\min}^{p}\right)^{\tau} \\
= \left(1 - \gamma \mu \lambda_{\min}^{p}\right)^{t+1} \mathbb{E}\left[\Phi_{0}\right] + \frac{\gamma \lambda_{\max}^{p} b \sigma^{2}}{(2 - b) \gamma \mu \lambda_{\min}^{p}},$$

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

H.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**
$$t = 0, 1, ..., T - 1$$
 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:
$$\hat{A}^t = -\eta \left(\left(B_S^t \right)^\top B_S^t \right)^\dagger \left(B_S^t \right)^\top G^t$$
7:
$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$$

7:
$$W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t$$

8:

9: Sample
$$A_S^t \sim \mathcal{D}_S^A$$

Right sketch

10:
$$\hat{B}^t = -\eta g(W^t) \left(A_S^t \right)^\top \left(A_S^t \left(A_S^t \right)^\top \right)^\dagger A_S^t$$

11:
$$W^{t+1} = W^t + \frac{\alpha}{r} \hat{B}^t A_S^t$$

12:

13: Sample i_{t+1} uniformly at random from [n]

14:
$$G^{t+1} = \begin{cases} \nabla f(W^{t+1}), & \text{with probability } q \\ G^t + \left(\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^t)\right), & \text{with probability } 1 - q \end{cases}$$
15: **end for**

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

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}$$
 (107)

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

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

$$(107)$$

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

$$\mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2}\right] \leq (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2}\right] + (1-q)L^{2}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right]. \tag{109}$$

Proof. Taking the full mathematical expectation, we obtain

$$\begin{split} \mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{\mathrm{F}}^{2}\right] &\stackrel{(108)}{=} (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t+1}) + \left(\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^{t})\right)\right\|_{\mathrm{F}}^{2}\right] \\ &\stackrel{(13)}{=} (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2}\right] \\ &+ (1-q)\mathbb{E}\left[\left\|\left(\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^{t})\right) - \left(\nabla f(W^{t+1}) - \nabla f(W^{t})\right)\right\|_{\mathrm{F}}^{2}\right] \\ &\leq (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2}\right] \\ &+ (1-q)\mathbb{E}\left[\left\|\nabla f_{i_{t+1}}(W^{t+1}) - \nabla f_{i_{t+1}}(W^{t})\right\|_{\mathrm{F}}^{2}\right] \\ &\leq (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathrm{F}}^{2}\right] + (1-q)L^{2}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathrm{F}}^{2}\right], \end{split}$$

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

H.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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \le \frac{2(f(W^0) - f^*)}{\lambda_{\min}^p \gamma T} + q \frac{\left\|G^0 - \nabla f(W^0)\right\|_{\mathrm{F}}^2}{T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},\tag{110}$$

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 Φ_t as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\text{max}}^p}{2q} \|G^t - \nabla f(W^t)\|_F^2.$$
 (111)

By Lemma 5 and Lemma 7, we have

$$\mathbb{E}\left[\Phi_{t+1}\right] \leq \mathbb{E}\left[f(W^{t})\right] - f^{*} - \frac{\gamma \lambda_{\min}^{p}}{2} \mathbb{E}\left[\left\|\nabla f(W^{t})\right\|_{F}^{2}\right] \\
- \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] + \frac{\gamma \lambda_{\max}^{p}}{2} \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}\right] \\
+ \frac{\gamma \lambda_{\max}^{p}(1-q)}{2q} \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}\right] + \frac{\gamma \lambda_{\max}^{p}(1-q)L^{2}}{2q} \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] \\
\leq \mathbb{E}\left[\Phi_{t}\right] - \frac{\gamma \lambda_{\min}^{p}}{2} \mathbb{E}\left[\left\|\nabla f(W^{t})\right\|_{F}^{2}\right] - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma(1-q)L^{2}\lambda_{\max}^{p}}{2q}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right].$$

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

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

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

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

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

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

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

H.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

$$\mathbb{E}\left[f(W^T) - f^*\right] \le (1 - \gamma \mu \lambda_{\min}^p)^T \Phi_0,\tag{112}$$

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

Proof. Denote Lyapunov function Φ_t as follows

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

By Lemma 5 and Lemma 7, we have

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

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

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

Unrolling the recursion, we obtain

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

23772378

2379

2380

2381

2382

2383

2384

2385

2386

2387 2388

2389

2390 2391

239223932394

241724182419

2420

2421

2422

2423

2424 2425

2426

2427

2428

2429

I 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; Sadiev et al., 2024; Yi et al., 2024).

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

I.1 ANALYSIS OF FED-BERNOULLI-LORA-QGD

Algorithm 7 Fed-Bernoulli-LoRA-QGD

```
2395
                 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}^B_S or \mathcal{D}^A_S, Bernoulli probabilities p and q
2396
2397
                 2: for t = 0, 1, \dots, T - 1 do
2398
                          for any client l \in [M] in parallel do
2399
                               Compute gradient \nabla f_l(W^{t+1}) and send compressed version G_l^t = \mathcal{Q}_l^t \left( \nabla f_l(W^{t+1}) \right) to the
                 4:
2400
2401
                 5:
                          end for
                         G^t = \frac{1}{M} \sum_{l=1}^{M} G_l^t
2402
                 6:
2403
                          Sample c^t \sim \text{Be}(p)
2404
                 7:
                                                                                                                                             Bernoulli random variable
                          if c^t = 1 then
                 8:
2405
                              Sample B_S^t \sim \mathcal{D}_S^B
                 9:
                                                                                                                                                                        Left sketch
2406
                             \hat{A}^t = -\eta \left( \left( B_S^t \right)^\top B_S^t \right)^\dagger \left( B_S^t \right)^\top G^t
W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t
2407
               10:
2408
               11:
2409
                          else
               12:
2410
                              Sample A_S^t \sim \mathcal{D}_S^A
               13:
                                                                                                                                                                     Right sketch
2411
                              \hat{B}^{t} = -\eta G^{t} \left( A_{S}^{t} \right)^{\top} \left( A_{S}^{t} \left( A_{S}^{t} \right)^{\top} \right)^{\dagger}
W^{t+1} = W^{t} + \frac{\alpha}{r} \hat{B}^{t} A_{S}^{t}
               14:
2412
2413
               15:
2414
               16:
2415
                          Broadcast W^{t+1} to each client l \in [M]
               17:
2416
               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:

2432
$$\min_{W \in \mathbb{R}^{m \times n}} \frac{1}{M} \sum_{l=1}^{M} f_l(W),$$
 2435

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 $Q: \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ 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:

$$\mathbb{E}[\mathcal{Q}(W)] = W, \quad and \quad \mathbb{E}\left[\|\mathcal{Q}(W) - W\|_{\mathcal{F}}^{2}\right] \le \omega \|W\|_{\mathcal{F}}^{2}. \tag{114}$$

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:

$$\Delta^* := f^* - \frac{1}{M} \sum_{l=1}^{M} f_l^* \ge 0. \tag{115}$$

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

$$G^{t} = \frac{1}{M} \sum_{l=1}^{M} \mathcal{Q}_{l}^{t} \left(\nabla f_{l}(W^{t}) \right); \tag{116}$$

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases}. \tag{117}$$

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 (116)) 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 $\nabla f(W^t)$:

$$\mathbb{E}\left[G^t|W^t\right] = \frac{1}{M} \sum_{l=1}^{M} \mathbb{E}\left[\mathcal{Q}_l^t\left(\nabla f_l(W^t)\right)|W^t\right] \stackrel{\text{(114)}}{=} \frac{1}{M} \sum_{l=1}^{M} \nabla f_l(W^t) = \nabla f(W^t).$$

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

$$\mathbb{E}\left[\left\|G^{t}\right\|_{F}^{2}|W^{t}\right] = \mathbb{E}\left[\left\|\frac{1}{M}\sum_{l=1}^{M}\mathcal{Q}_{l}^{t}\left(\nabla f_{l}(W^{t})\right) - \nabla f(W^{t}) + \nabla f(W^{t})\right\|_{F}^{2}|W^{t}\right]$$

$$\stackrel{(13)}{=} \mathbb{E}\left[\left\|\frac{1}{M}\sum_{l=1}^{M}\mathcal{Q}_{l}^{t}\left(\nabla f_{l}(W^{t})\right) - \nabla f(W^{t})\right\|_{F}^{2}|W^{t}\right] + \left\|\nabla f(W^{t})\right\|_{F}^{2}$$

$$= \frac{1}{M^{2}}\sum_{l=1}^{M}\mathbb{E}\left[\left\|\mathcal{Q}_{l}^{t}\left(\nabla f_{l}(W^{t})\right) - \nabla f_{l}(W^{t})\right\|_{F}^{2}|W^{t}\right] + \left\|\nabla f(W^{t})\right\|_{F}^{2}$$

$$\stackrel{(114)}{\leq} \frac{\omega}{M^{2}}\sum_{l=1}^{M}\left\|\nabla f_{l}(W^{t})\right\|_{F}^{2} + \left\|\nabla f(W^{t})\right\|_{F}^{2}$$

$$\stackrel{(*)}{\leq} \frac{2L\omega}{M^{2}}\sum_{l=1}^{M}\left(f_{l}(W^{t}) - f_{l}^{*}\right) + \left\|\nabla f(W^{t})\right\|_{F}^{2}$$

$$= 2\frac{L\omega}{M}\left(f(W^{t}) - f^{*}\right) + \left\|\nabla f(W^{t})\right\|_{F}^{2} + 2\frac{L\omega}{M}\underbrace{\left(f^{*} - \frac{1}{M}\sum_{l=1}^{M}f_{l}^{*}\right)}_{:=\Delta^{*}},$$

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

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

I.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

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

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 $\left\{W^0, W^1, \dots, W^{T-1}\right\}$ with probabilities $\left\{\frac{w_t}{W_{T-1}}\right\}_{t=0}^{T-1}$, where $w_t = \frac{w_{t-1}}{(1+\gamma^2L^2\lambda_{\max}^p\omega/M)}$, $\mathcal{W}_{T-1} = \sum_{t=0}^{T-1} w_t$, and $w^{-1} > 0$.

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

I.1.2 Convergence under Polyak-Łojasiewicz Condition

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

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

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

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

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

I.2 ANALYSIS OF FED-BERNOULLI-LORA-MARINA

Algorithm 8 Fed-Bernoulli-LoRA-MARINA

2592

259325942595

262226232624

2625

2626

2627

2629

26302631

2632

2633

2634

2635

2636

26372638

2639

2640 2641 2642

2643

2644 2645

```
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
2596
2597
2598
2599
                  2: for t = 0, 1, \dots, T - 1 do
2600
                            Sample c^t \sim \text{Be}(p)
                                                                                                                                                       Bernoulli random variable
                            if c^t = 1 then
                  4:
                                 Sample B_S^t \sim \mathcal{D}_S^B
                                                                                                                                                                                   Left sketch
2602
                               \hat{A}^t = -\eta \left( \left( B_S^t \right)^\top B_S^t \right)^\dagger \left( B_S^t \right)^\top G^t
W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t
2603
2604
                  7:
                  8:
2606
                                Sample A_S^t \sim \mathcal{D}_S^A
                  9:
                                                                                                                                                                                Right sketch
2607
                                \hat{B}^{t} = -\eta G^{t} \left( A_{S}^{t} \right)^{\top} \left( A_{S}^{t} \left( A_{S}^{t} \right)^{\top} \right)^{\dagger}
2608
2609
                                 W^{t+1} = W^t + \frac{\alpha}{\pi} \hat{B}^t A_S^t
                11:
2610
                12:
2611
                            Broadcast W^{t+1} to each client l \in [M]
                13:
2612
                            Sample s^t \sim \text{Be}(q)
2613
                15:
                            for any client l \in [M] in parallel do
2614
                                Compute gradient \nabla f_l(W^{t+1})
                               G_l^{t+1} = \begin{cases} \nabla f_l(W^{t+1}), & \text{with probability } q \\ G_l^t + \mathcal{Q}_l^t \left( \nabla f_l(W^{t+1}) - \nabla f_l(W^t) \right), & \text{with probability } 1 - q \end{cases}
2615
2616
2617
2618
                19:
                           G^{t+1} = \frac{1}{M} \sum_{l=1}^{M} G_l^{t+1}
2619
                20:
2620
2621
                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 & 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):

$$W^{t+1} = W^t - \gamma \hat{G}^t, \text{ where } \hat{G}^t = \begin{cases} H_B^t G^t, \text{ with probability } p \\ G^t H_A^t, \text{ with probability } 1 - p \end{cases}$$
 (118)

$$G_l^{t+1} = \begin{cases} \nabla f_l(W^{t+1}), & \text{with probability } q \\ G_l^t + \mathcal{Q}_l^t \left(\nabla f_l(W^{t+1}) - \nabla f_l(W^t) \right), & \text{with probability } 1 - q \end{cases}$$
(119)

$$G^{t+1} = \frac{1}{M} \sum_{l=1}^{M} G_l^{t+1}. \tag{120}$$

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

$$\mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{\mathcal{F}}^{2}\right] \leq (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{\mathcal{F}}^{2}\right] + (1-q)\frac{\omega L^{2}}{M}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{\mathcal{F}}^{2}\right]. \tag{121}$$

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

$$\mathbb{E}\left[\left\|G^{t+1} - \nabla f(W^{t+1})\right\|_{F}^{2} |W^{t+1}\right] = (1-q)\mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t}) + \frac{1}{M}\sum_{l=1}^{M} \mathcal{Q}_{l}^{t} \left(\nabla f_{l}(W^{t+1}) - \nabla f_{l}(W^{t})\right)\right\|_{F}^{2} |W^{t+1}|\right] \\
\stackrel{(13)}{=} (1-q)\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2} + (1-q)\mathbb{E}\left[\left\|\frac{1}{M}\sum_{l=1}^{M} \mathcal{Q}_{l}^{t} \left(D_{l}^{t+1}\right) - D^{t+1}\right\|_{F}^{2} |W^{t+1}|\right] \\
= (1-q)\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2} + \frac{1-q}{M^{2}}\sum_{m=1}^{M}\mathbb{E}\left[\left\|\mathcal{Q}_{l}^{t} \left(D_{l}^{t+1}\right) - D_{l}^{t+1}\right\|_{F}^{2} |W^{t+1}|\right] \\
\stackrel{(114)}{\leq} (1-q)\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2} + \frac{(1-q)\omega}{M^{2}}\sum_{l=1}^{M}\left\|\nabla f_{l}(W^{t+1}) - \nabla f_{l}(W^{t})\right\|_{F}^{2} \\
\leq (1-q)\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2} + \frac{(1-q)\omega L^{2}}{M}\left\|W^{t+1} - W^{t}\right\|_{F}^{2},$$

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

I.2.1 Convergence for Smooth Non-Convex Functions

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

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

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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \le \frac{2\left(f(W^0) - f^*\right)}{\gamma \lambda_{\min}^p T} + \frac{\left\|G^0 - \nabla f(W^0)\right\|_{\mathrm{F}}^2}{qT} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},\tag{122}$$

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 Φ_t as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\text{max}}^p}{2a} \|G^t - \nabla f(W^t)\|_F^2.$$
 (123)

By Lemma 5 and Lemma 9, we have

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

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

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

Summing over, we get

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

Finally, we derive

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

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

I.2.2 Convergence under Polyak-Łojasiewicz Condition

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

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

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

$$\mathbb{E}\left[f(W^T) - f^*\right] \le (1 - \gamma \mu \lambda_{\min}^p)^T \Phi_0, \tag{124}$$

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 $\Phi_0 = f(W^0) - f^* + \frac{\gamma\lambda_{\max}^{p}}{q} \|G^0 - \nabla f(W^0)\|_{\mathrm{F}}^2$.

Proof. Denote Lyapunov function Φ_t as follows

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

By Lemma 5 and Lemma 7, we have

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

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

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

Taking recursion, we have

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

I.3 ANALYSIS OF FED-BERNOULLI-LORA-EF21

Algorithm 9 Fed-Bernoulli-LoRA-EF21

20: **end for**

278027812782

2783

2784

2785

2786

2787

2788

27892790

2791

2792

2793

2794

2795

2796

2797

2798

2799

2801

2802

2804

2805

2806

2807

2754

27552756

```
2757
               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\},
2758
                    scaling factor \alpha > 0, chain length T, sketch distribution \mathcal{D}_S^B or \mathcal{D}_S^A, Bernoulli probability p
2759
               2: for t = 0, 1, \dots, T - 1 do
2760
                       Sample c^t \sim \text{Be}(p)
                                                                                                                             Bernoulli random variable
2761
                       if c^t = 1 then
               4:
2762
                           Sample B_S^t \sim \mathcal{D}_S^B
                                                                                                                                                     Left sketch
2763
                          \hat{A}^t = -\eta \left( (B_S^t)^\top B_S^t \right)^\dagger \left( B_S^t \right)^\top G^t
W^{t+1} = W^t + \frac{\alpha}{r} B_S^t \hat{A}^t
2764
2765
2766
               8:
2767
               9:
                           Sample A_S^t \sim \mathcal{D}_S^A
                                                                                                                                                   Right sketch
2768
                          \hat{B}^{t} = -\eta G^{t} \left( A_{S}^{t} \right)^{\top} \left( A_{S}^{t} \left( A_{S}^{t} \right)^{\top} \right)^{\dagger}
2769
                           W^{t+1} = W^t + \frac{\alpha}{2} \hat{B}^t A_S^t
2770
             12:
                       Broadcast W^{t+1} to each client l \in [M]
             13:
2772
                       for any client l \in [M] in parallel do
             14:
2773
                           Compute gradient \nabla f_l(W^{t+1})
             15:
2774
                           G_l^{t+1} = G_l^t + \mathcal{C}_l^t \left( \nabla f_l(W^{t+1}) - G_l^t \right)
             16:
2775
                           Send G_l^{t+1} to the server
             17:
2776
             18:
2777
                      G^{t+1} = \frac{1}{M} \sum_{l=1}^{M} G_l^{t+1}
2778
             19:
2779
```

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

$$\mathbb{E}\left[\left\|\mathcal{C}\left(W\right) - W\right\|_{\mathrm{F}}^{2}\right] \le (1 - \beta) \left\|W\right\|_{\mathrm{F}}^{2}, \quad \forall W \in \mathbb{R}^{m \times n}.$$
 (126)

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

$$W^{t+1} = W^t - \gamma \hat{G}^t, \quad \text{where} \quad \hat{G}^t = \begin{cases} H_B^t G^t, & \text{with probability } p \\ G^t H_A^t, & \text{with probability } 1 - p \end{cases} \tag{127}$$

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

$$(128)$$

$$G^{t+1} = \frac{1}{M} \sum_{l=1}^{M} G_l^{t+1}. \tag{129}$$

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

$$\mathbb{E}\left[\left\|G_{l}^{t+1} - \nabla f_{l}(W^{t+1})\right\|_{F}^{2}\right] \leq \sqrt{1-\beta}\mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] + \frac{(1-\beta)L^{2}}{1-\sqrt{1-\beta}}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right]$$

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

$$\mathbb{E}\left[\left\|G_{l}^{t+1} - \nabla f_{l}(W^{t+1})\right\|_{\mathrm{F}}^{2}\right] \stackrel{(128),(129)}{=} \mathbb{E}\left[\mathbb{E}\left[\left\|C_{l}^{t}\left(\nabla f_{l}(W^{t+1}) - G_{l}^{t}\right) - \left(\nabla f_{l}(W^{t+1}) - G_{l}^{t}\right)\right\|_{\mathrm{F}}^{2}|G_{l}^{t+1}, W^{t+1}|\right]\right] \\ \stackrel{(126)}{\leq} \left(1 - \beta\right) \mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t+1})\right\|_{\mathrm{F}}^{2}\right] \\ \stackrel{(1-\beta)}{\leq} \left(1 - \beta\right) \left(1 + \theta\right) \mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{\mathrm{F}}^{2}\right] \\ + \left(1 - \beta\right) \left(1 + \frac{1}{\theta}\right) \mathbb{E}\left[\left\|\nabla f_{l}(W^{t+1}) - \nabla f_{l}(W^{t})\right\|_{\mathrm{F}}^{2}\right],$$

where in the last inequality we used $\|U+V\|_{\mathrm{F}}^2 \leq (1+\theta) \|U\|_{\mathrm{F}}^2 + \left(1+\frac{1}{\theta}\right) \|V\|_{\mathrm{F}}^2$ 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

$$\mathbb{E}\left[\left\|G_{l}^{t+1} - \nabla f_{l}(W^{t+1})\right\|_{F}^{2}\right] \leq \sqrt{1-\beta}\mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] + \frac{1-\beta}{1-\sqrt{1-\beta}}\mathbb{E}\left[\left\|\nabla f_{l}(W^{t+1}) - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] \\
\leq \sqrt{1-\beta}\mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] + \frac{(1-\beta)L^{2}}{1-\sqrt{1-\beta}}\mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right],$$

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.

I.3.1 Convergence for Smooth Non-Convex Functions

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

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

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

$$\mathbb{E}\left[\left\|\nabla f(\widetilde{W}^T)\right\|_{\mathrm{F}}^2\right] \le \frac{2(f(W^0) - f^*)}{\gamma \lambda_{\min}^p T} + \frac{\mathcal{G}^0}{(1 - \sqrt{1 - \beta})T} \cdot \frac{\lambda_{\max}^p}{\lambda_{\min}^p},\tag{130}$$

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 $\mathcal{G}^0 := \frac{1}{M} \sum_{l=1}^M \left\| \mathcal{G}_l^0 - \nabla f_l(W^0) \right\|_{\mathrm{F}}^2$.

Proof. Denote Lyapunov function Φ_t as follows

$$\Phi_t = f(W^t) - f^* + \frac{\gamma \lambda_{\text{max}}^p}{2(1 - \sqrt{1 - \beta})} \cdot \frac{1}{M} \sum_{l=1}^M \|G_l^t - \nabla f_l(W^t)\|_F^2.$$
 (131)

By Lemma 5 and Lemma 10, we have

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

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

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

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

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

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

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

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

I.3.2 CONVERGENCE UNDER POLYAK-ŁOJASIEWICZ CONDITION

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

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

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

$$\mathbb{E}\left[f(W^T) - f^*\right] \le (1 - \gamma \mu \lambda_{\min}^p)^T \Phi_0, \tag{132}$$

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 $\Phi_0 = f(W^0) - f^* + \frac{\gamma \lambda_{\max}^p}{1 - \sqrt{1 - \beta}} \frac{1}{M} \sum_{l=1}^M \left\| G_l^0 - \nabla f_l(W^0) \right\|_{\mathrm{F}}^2$.

Proof. Denote Lyapunov function Φ_t as follows

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

By Lemma 5 and Lemma 10, we have

$$\mathbb{E}\left[\Phi_{t+1}\right] \leq \mathbb{E}\left[f(W^{t})\right] - f^{*} - \frac{\gamma \lambda_{\min}^{p}}{2} \mathbb{E}\left[\left\|\nabla f(W^{t})\right\|_{F}^{2}\right] - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] \\
+ \frac{\gamma \lambda_{\max}^{p}}{2} \cdot \mathbb{E}\left[\left\|G^{t} - \nabla f(W^{t})\right\|_{F}^{2}\right] + \frac{\gamma \lambda_{\max}^{p} \sqrt{1-\beta}}{1-\sqrt{1-\beta}} \cdot \frac{1}{M} \sum_{l=1}^{M} \mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] \\
+ \frac{\gamma \lambda_{\max}^{p} (1-\beta) L^{2}}{(1-\sqrt{1-\beta})^{2}} \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right] \\
\leq (1-\gamma\mu\lambda_{\min}^{p}) \mathbb{E}\left[f(W^{t}) - f^{*}\right] + \frac{\gamma\lambda_{\max}^{p} (1+\sqrt{1-\beta})}{2(1-\sqrt{1-\beta})} \cdot \frac{1}{M} \sum_{l=1}^{M} \mathbb{E}\left[\left\|G_{l}^{t} - \nabla f_{l}(W^{t})\right\|_{F}^{2}\right] \\
- \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma\lambda_{\max}^{p} (1-\beta) L^{2}}{(1-\sqrt{1-\beta})^{2}}\right) \mathbb{E}\left[\left\|W^{t+1} - W^{t}\right\|_{F}^{2}\right],$$

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

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

Taking the recursion, we have

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

complete it was that from new reps

EXPERIMENTS: MISSING DETAILS

In this section, we provide additional details regarding the experimental setting from Section 7.

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.

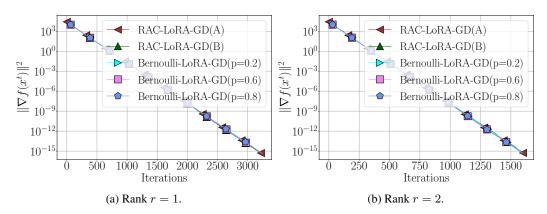


Figure 2: Comparison of RAC-LoRA-GD and Bernoulli-LoRA-GD on linear regression fine-tuning. Curves with $p = 0.01, 0.2, \dots$ indicate Bernoulli-LoRA-GD sampling parameters. RAC-LoRA-GD(A) trains B after resampling A, while RAC-LoRA-GD(B) does the reverse. All methods use $\gamma = c/\hat{L}$ with $c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

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

sklearn.datasets.make_regression

with moderate noise and a controlled rank structure:

```
2970
       wt_D, wt_b = make_regression(n_samples=90000, n_features=4096,
2971 2
                                        n_informative=4096, noise=20.0,
2972 <sub>3</sub>
                                        bias=0.0, tail_strength=0.8,
2973 4
                                        effective_rank=64, random_state=42)
2974
       followed by standard scaling. The fine-tuning dataset (\hat{D}, \hat{b}) was produced similarly:
2975
       h_D, h_b = make_regression(n_samples=10000, n_features=4096,
2977 2
                                      n_{informative=4096//2, noise=50.0,
2978 3
                                      bias=10.0, tail_strength=0.9,
                                      effective_rank=32, random_state=84)
2979 4
2980
       and subsequently adjusted with a biased scaling (mean 1, standard deviation 2).
2981
```

LLM USE ACKNOWLEDGMENT

In this paper, we used large language models (LLMs) to assist with grammar and wording during the preparation of the manuscript. We did not use LLMs to derive convergence theorems, generate empirical plots, or search for citations. This usage is in accordance with two primary LLM-related policies.