
Learning Provably Improves the Convergence of Gradient Descent

Qingyu Song*
Xiamen University
simmonssong96@gmail.com

Wei Lin, Hong Xu
The Chinese University of Hong Kong
wlin23@cse.cuhk.edu.hk, hongxu@cuhk.edu.hk

Abstract

Learn to Optimize (L2O) trains deep neural network-based solvers for optimization, achieving success in accelerating convex problems and improving non-convex solutions. However, L2O lacks rigorous theoretical backing for its own training convergence, as existing analyses often use unrealistic assumptions—a gap this work highlights empirically. We bridge this gap by proving the training convergence of L2O models that learn Gradient Descent (GD) hyperparameters for quadratic programming, leveraging the Neural Tangent Kernel (NTK) theory. We propose a deterministic initialization strategy to support our theoretical results and promote stable training over extended optimization horizons by mitigating gradient explosion. Our L2O framework demonstrates over 50% better optimality than GD and superior robustness over state-of-the-art L2O methods on synthetic datasets. The code of our method can be found from <https://github.com/NetX-lab/MathL2OProof-Official>.

1 Introduction

Learn to optimize (L2O) represents an increasingly influential paradigm for tackling optimization problems [6]. Numerous studies have demonstrated the efficacy of employing learning-based models to achieve superior performance across a spectrum of optimization tasks. These encompass convex problems, exemplified by LASSO [7, 8, 22] and logistic regression [23, 34], and non-convex scenarios such as MIMO sum-rate maximization [35] and network resource allocation [33].

Distinct from black-box approaches [5, 36, 41], which directly derive solutions to optimization problems from a neural network (NN), the so-called “white-box” methodologies are garnering increased attention. This heightened interest stems from their inherent advantages, such as enhanced trustworthiness [14] and theoretical guarantees [34]. A key characteristic of these white-box strategies is the integration of mechanisms to ensure the “controllability” of the generated solutions. For instance, Lv et al. [25] employ a NN to predict the step size for the gradient descent (GD) algorithm, where the inherent structure of GD stabilizes the optimization trajectory. Similarly, Heaton et al. [14] integrate a conventional solver within an L2O framework to act as a safeguard, thereby preventing the learning-based model from producing solutions with extreme violations. This principle of guided or constrained learning has also been extended to the training phase of L2O models [39].

Further, “unrolling” has emerged as a prominent technique within L2O [6], characterized by the strategic replacement of components of conventional optimization algorithms with neural network (NN) blocks [12, 15, 20]. For instance, Liu et al. [23] introduce Math-L2O that imposes architectural constraints on unrolled L2O models by deriving necessary conditions for their convergence. Their analysis revealed that for a L2O model to achieve optimality, its embedded NN must effectively perform a linear combination of input feature vectors, weighted by learnable parameter matrices.

*This work was done at The Chinese University of Hong Kong (CUHK) when Qingyu was a PhD candidate.

Empirical validation demonstrates that the proposed methods exhibit strong generalization capabilities when trained using a coordinate-wise input-to-output strategy. Subsequent research by Song et al. [34] further enhance this generalization performance by reducing the magnitude of input features.

Despite these advancements, to the best of our knowledge, a formal demonstration of the convergence for unrolling-based L2O methods in solving general optimization problems remains elusive. While LISTA-CPSS [7] establishes convergence for the well-known LISTA framework [12], its analysis is based on the assumption that neural network (NN) outputs are confined to a specific subspace, a condition that is often not met in practical implementations. Similarly, while Math-L2O [23] derives necessary conditions for convergence, the mechanisms by which the training process itself can guarantee such convergence are not elucidated. Subsequent analysis by Song et al. [34] investigates the inference-time convergence of Math-L2O. However, this work relies on a stringent training assumption, effectively constraining the L2O model to emulate the behavior of a conventional Gradient Descent (GD) algorithm.

This apparent deficiency in comprehensively demonstrating L2O convergence stems from two fundamental, unresolved technical challenges. First, unrolling-based L2O models [8, 12, 22] represent a specialized class of NN architectures. Despite much progress in understanding the training convergence of general neural networks (NNs), notably through the Neural Tangent Kernel (NTK) theory since 2019 [2, 3, 11, 24, 29, 30], a formal proof of training convergence remains conspicuously absent. Such a proof is an essential precursor to establishing the convergence of the L2O model in its primary task of solving optimization problems. Second, the precise relationship between the training convergence achieved during the L2O model’s training phase (i.e., optimizing the NN parameters) and the convergence of the L2O model when applied to the target optimization problem (i.e., finding the optimal solution) is not well understood. For instance, Math-L2O [23] is designed to learn the step size for an underlying GD algorithm. While the problem-solving efficacy of Math-L2O is naturally evaluated based on the progression of GD iterations, its training convergence is measured in terms of training steps (e.g., epochs). These two notions of convergence: one on model parameter optimization and the other on problem-solving iterations, are largely decoupled and operate on fundamentally different scales.

In this work, we present the first rigorous demonstration that an unrolling framework can achieve theoretical convergence in solving optimization problems. Our analysis focuses on the state-of-the-art (SOTA) Math-L2O framework, wherein a NN functions as a recurrent block, iteratively generating hyperparameters for an underlying optimization algorithm. The solution obtained at each iteration, which utilizes these generated hyperparameters, is then incorporated as an input feature for the subsequent iteration [23]. This inherent recurrence imparts RNN-like characteristics to Math-L2O, significantly complicating the analysis of its training convergence. Specifically, the recurrent structure causes the NN to manifest as a high-order polynomial function with respect to (w.r.t.) its input features [3]. This characteristic poses challenges for establishing tight analytical bounds, potentially leading to looser convergence rates compared to non-recurrent architectures, as highlighted in related NTK analyses for RNNs [3]. Moreover, the Math-L2O architecture introduces an additional layer of complexity: the emergence of high-order polynomial dependencies not only on the input features but also on the learnable parameters themselves. This distinct feature renders the convergence proof for Math-L2O arguably more intricate than those for conventional RNNs, where such parameter-dependent high-order terms are typically less pronounced.

We address the pivotal connection between the NN’s training convergence and the ultimate problem-solving convergence of the L2O model. Within the Math-L2O framework, we establish this critical linkage by explicitly demonstrating an alignment between the convergence dynamics exhibited during the NN’s training phase and the convergence characteristics of its underlying backbone optimization algorithm. This alignment provides a novel theoretical bridge, ensuring that a successfully trained L2O model translates to effective convergence when applied to optimization tasks. Our contributions are summarized as follows:

1. We provide a formal proof that the Math-L2O training framework substantially enhances the convergence performance of its underlying backbone algorithms. This is achieved by rigorously establishing an explicit alignment between the convergence rates of the training process and the iterative steps of the backbone algorithm.

2. We establish the first linear convergence rate for Math-L2O training. Inspired by [29], we employ a NN architecture with a single wide layer and utilize NTK to prove the boundedness of NN outputs, gradients, and the training loss function within the Math-L2O framework.
3. We introduce a novel deterministic parameter initialization scheme, coupled with a specific learning rate configuration strategy. This combined approach is proven to guarantee the training convergence of the Math-L2O model across all iterations.
4. We empirically validate our theoretical findings through comprehensive experiments. The results showcase significant performance advantages, including up to a 50% improvement in solution optimality over the standard GD algorithm post-training, and superior robustness compared to SOTA L2O models and the Adam optimizer [10]. Furthermore, ablation studies empirically confirm the practical efficacy and individual contributions of our proposed theorems.

2 Preliminary

This section first defines the optimization problem objective and the L2O framework. The L2O training loss is then formulated based on these definitions. Then, the NN’s computational graph is employed to detail the forward pass and the derivation of parameter gradients.

2.1 Definitions

Let $d > b$, suppose $x \in \mathbb{R}^{d \times 1}$, $y \in \mathbb{R}^{b \times 1}$, and $\mathbf{M} \in \mathbb{R}^{b \times d}$, we define the optimization objective as:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{2} \|\mathbf{M}x - y\|_2^2. \quad (1)$$

This objective function is commonly selected for convergence analysis [4]. The least-squares problem, a frequent subject in NN convergence studies [2, 3, 11, 21, 29], is a specific instance of the minimization in Equation (1) where $d = b$ and $\mathbf{M} = \mathbf{I}$.

We assume f to be β -smooth, such that $\|\mathbf{M}^\top \mathbf{M}\|_2 \leq \beta$, and \mathbf{M} to possess full row rank, with $\lambda_{\min}(\mathbf{M}\mathbf{M}^\top) = \beta_0 > 0$. This setting often favors numerical algorithms (e.g., GD) over analytical solutions due to computational complexity. GD’s $\mathcal{O}(bd)$ complexity is typically lower than the $\mathcal{O}(b^3)$ of analytical methods involving costly matrix inversions. The loss function is then defined as the sum of N objectives specified in Equation (1):

$$F(X) = \frac{1}{2} \|\mathbf{M}X - Y\|_2^2, \quad (2)$$

where F , $\mathbf{M} \in \mathbb{R}^{Nb \times Nd}$, $X \in \mathbb{R}^{Nd \times 1}$, and $Y \in \mathbb{R}^{Nb \times 1}$ represent the concatenated objectives, parameters, variables, and labels, respectively, from N optimization problems (see Appendix A.1 for details). F is also β -smooth, given that $\|\mathbf{M}^\top \mathbf{M}\|_2 \leq \max_{i=1, \dots, N} \{\|\mathbf{M}_i^\top \mathbf{M}_i\|_2\} = \beta$.

Learn to Optimize (L2O). Given an initial point X_0 , L2O takes X_0 as the input and generates a solution, denoted as X_t , with a machine learning model. Typically, let g_W denote an L -layer NN with parameters $W = \{W_1, \dots, W_L\}$, $W_\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$, $n_1, \dots, n_L \in \mathbb{R}$. Math-L2O [23] takes an iterative workflow to generate solutions. For each step $t \in [T]$ in solving the problem in Equation (1), the NN model in Math-L2O is defined as $g_W(X_{t-1}, \nabla F(X_{t-1}))$. The NN receives the current state variable X_{t-1} and its gradient $\nabla F(X_{t-1})$ as input. The update rule at step t , which employs the Hadamard product (denoted by \odot), is formulated as:

$$X_t = X_{t-1} - \frac{1}{\beta} P_t \odot \nabla F(X_{t-1}), P_t = g_W(X_{t-1}, \nabla F(X_{t-1})). \quad (3)$$

P_t represents a vector whose entries are learned step sizes. The NN g_W takes structured layer-wise architecture. It employs a coordinate-wise architecture, processing each input dimension independently, recognized for its robustness in L2O applications [23, 34]. Thus, output dimension of the NN is one, i.e., $n_L = 1$. Denote $[L] := \{1, \dots, L\}$, for layer $\ell \in [L]$, we denote $G_{\ell,t}$ as the (inner) output of layer ℓ at step t . Utilizing ReLU (ReLU) [1] and Sigmoid (σ) [27] activations, $G_{\ell,t}$ is defined as:

$$G_{\ell,t} = \begin{cases} [X_{t-1}, \nabla F(X_{t-1})]^\top & \ell = 0, \\ \text{ReLU}(W_\ell G_{\ell-1,t}) & \ell \in [L-1], \\ P_t = 2\sigma(W_L G_{L-1,t})^\top & \ell = L. \end{cases} \quad (4)$$

The L2O training problem is defined by:

$$F(W) = \frac{1}{2} \|\mathbf{M}X_T - Y\|_2^2, X_T = L2O_W(X_0, \nabla F(X_0)). \quad (5)$$

2.2 Layer-Wise Derivative of NN's Parameters

Let k denote a training iteration for loss Equation (5) minimization, which is distinct from an optimization step t for solving objective Equation (2). The computational graph in Figure 1 illustrates the Math-L2O forward and backward operations, which parallel those of Recurrent Neural Networks (RNNs) [13]. Figure 1a details the NN block (see Equation (4)). Figure 1b depicts the overall process: the block takes an input solution, performs T internal optimization steps to produce an updated solution (red dashed arrows), and each training iteration k triggers a full backward pass (blue bold lines). As per [23], the gradient flow from the input features to the NN block is detached.

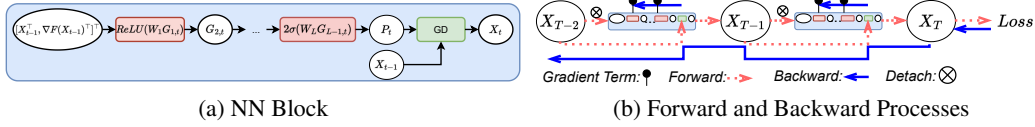


Figure 1: Computational Graph of Math-L2O

The derivative of an objective F w.r.t. the parameters W_ℓ of layer ℓ is determined via the computational graph, paralleling Back-Propagation-Through-Time (BPTT) for RNNs [26]:

$$\frac{\partial F}{\partial W_\ell} = \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{t=1}^T \left(\prod_{j=T}^{t+1} \frac{\partial X_j}{\partial X_{j-1}} \right) \frac{\partial X_t}{\partial P_t} \frac{\partial P_t}{\partial W_\ell} \right). \quad (6)$$

The summation aggregates gradients across T optimization steps. $\prod_{j=T}^{t+1} (\partial X_j / \partial X_{j-1})$ represents the chain rule application from the final output X_T to an intermediate state X_t .

Moreover, we derive two key gradients, instrumental for establishing the theoretical results in the ensuing section. Following Definition 2.2 in [2], the gradient of the ReLU is represented by a diagonal matrix \mathbf{D}_ℓ^t , where its i -th diagonal element is $[\mathbf{D}_\ell^t]_{i,i} := \mathbf{1}_{(W_\ell G_{\ell-1,t})_{i \geq 0}}$ for $i \in [n_\ell]$. Let $\Gamma_t := \mathbf{M}^\top (\mathbf{M}X_t - Y)$ and $\Xi_\ell := (\mathbf{I}_d \otimes W_L) (\prod_{j=L-1}^{\ell+1} \mathbf{D}_{j,t} (\mathbf{I}_d \otimes W_j)) \mathbf{I}_{n_\ell}$. Defining $\mathcal{D}(\cdot)$ as the operator that constructs a diagonal matrix from a vector, the gradients for an inner layer W_ℓ ($\ell < L$) and the final layer W_L are given by:

$$\frac{\partial F}{\partial W_\ell} = -\frac{1}{\beta} \Gamma_T^\top \sum_{t=1}^T \left(\prod_{j=T}^{t+1} (\mathbf{I}_d - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M} \mathcal{D}(P_j)) \right) \mathcal{D}(\Gamma_t) \mathcal{D}(P_t \odot (1 - P_t/2)) \Xi_\ell \otimes G_{\ell-1,t}^\top, \quad (7)$$

$$\frac{\partial F}{\partial W_L} = -\frac{1}{\beta} \Gamma_T^\top \sum_{t=1}^T \left(\prod_{j=T}^{t+1} (\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_j) \mathbf{M}^\top \mathbf{M}) \right) \mathcal{D}(\Gamma_T) \mathcal{D}(P_t \odot (1 - P_t/2)) G_{L-1,t}^\top, \quad (8)$$

where \otimes denotes the Kronecker product. Equation (8) (for W_L) differs from Equation (7) (for W_ℓ) in its final terms: $G_{L-1,t}^\top$ replaces $\Xi_L \otimes G_{\ell-1,t}^\top$. This simplification arises as W_L is the terminal layer, and $G_{L-1,t}$ is its direct input from layer $L-1$. Thus, its gradient calculation does not involve a subsequent layer propagation factor analogous to Ξ_L .

3 L2O Convergence Demonstration Framework

This section rigorously substantiates the convergence of the L2O framework, Math-L2O. We first expose theoretical and numerical instabilities prevalent in current SOTA L2O methods. Then, we demonstrate Math-L2O's accelerated training convergence compared to GD and then present a formal methodology to establish its convergence.

3.1 Limitations Analysis of Existing SOTA L2O Frameworks

We analyze limitations in the convergence guarantees of two SOTA L2O frameworks: LISTA-CPSS [7] and Math-L2O [23]. LISTA-CPSS [7] constructively proves that its predecessor, LISTA [12], can attain a linear convergence rate. However, this theoretical guarantee is contingent upon several stringent conditions. Math-L2O [23] proposes an L2O framework derived from the GD algorithm, incorporating necessary conditions for convergence. Both frameworks employ sequential solution updates and utilize BPTT for parameter optimization.

Initially, we assess training loss across varying optimization steps. This is pertinent due to the well-documented issue of gradient explosion of BPTT arising from long-term gradient accumulation [19].

Both models are trained on 10 randomly sampled optimization problems for 400 epochs. Figure 2 depicts training losses (y-axis) against optimization steps (x-axis) for several learning rates (distinguished by line color). Data points exhibiting numerical overflow (indicative of gradient explosion at the first training iteration) are excluded, resulting in plot lines terminating before 100 steps for affected configurations. The results demonstrate that both frameworks suffer from poor convergence at low learning rates (LRs) and training instability at high LRs.

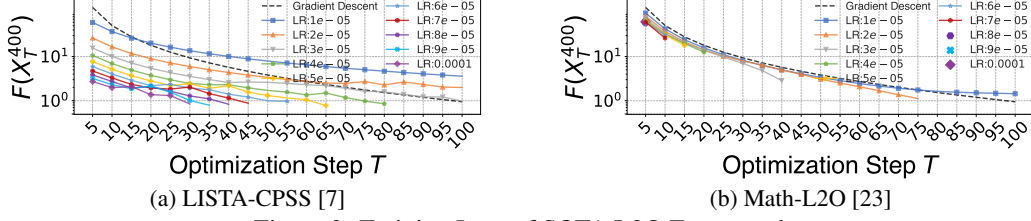


Figure 2: Training Loss of SOTA L2O Frameworks

Further, we examine the convergence conditions outlined for LISTA-CPSS [7], illustrating their propensity for violation during typical training procedures. The first condition mandates asymptotic sign consistency between iterates X_t and the solution X^* , requiring $\text{sign}(X_t) = \text{sign}(X^*)$ for all t . The second condition imposes constraints on the columns of the learned parameter matrix \mathbf{W} relative to the columns of the objective coefficient matrix \mathbf{M} . Specifically, denoting column indices by i and j , it necessitates that $\mathbf{W}_i^\top \mathbf{M}_i = 1$ and $\mathbf{W}_i^\top \mathbf{M}_j > 1$ for all $j \neq i$.

Following the experimental design in [23], we quantify the violation percentage of the aforementioned conditions during inference. The results are presented in Figure 3. We consider two settings: (i) shared parameters W across iterations (Figure 3a), and (ii) unique parameters W_t per step t (Figure 3b). Both scenarios reveal that the specified conditions are frequently violated post-training. For instance, in the shared W case (Figure 3a), while the conditions hold in later steps, substantial violations occur in early steps. The divergence contradicts the convergence rate analysis presented in [7].

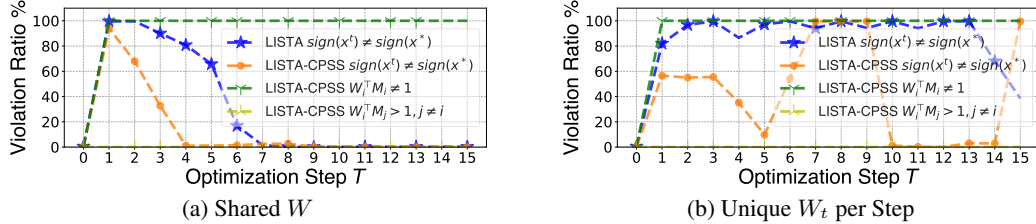


Figure 3: Violation Ratio of LISTA-CPSS Conditions During Inference

The preceding observations highlight that training is indispensable for L2O convergence analysis. Three fundamental questions arise in L2O: (i) *What is the impact of training on convergence?* (ii) *How can training be incorporated into the convergence analysis framework?* (iii) *What mechanisms ensure a stable training process?* We propose a concise approach to address these questions, establishing a direct alignment between the training’s convergence rate and an existing algorithm’s rate.

3.2 L2O Convergence Demonstration: Aligning L2O to An Algorithm

First, we introduce a general convergence analysis framework. Let X^* be the optimal solution, r_t represents an iteration-dependent rate term, and $C(X_0)$ be a constant that dependent on the initial point X_0 (and X^*), the convergence rate of an algorithm (either learned or classical) for minimizing an objective $F(X)$ (e.g., the objective in Equation (1) or the loss in Equation (2)) is often formulated as: $F(X_t) \leq r_t C(X_0)$. For example, standard GD has a rate of $F(X_t) \leq \frac{\beta}{t} \|X_0 - X^*\|_2^2$ [4].

The performance of L2O models, stabilized via training, is typically assessed after T iterations [23, 34]. We formulate the L2O training convergence rate w.r.t. training iteration k as:

$$F(X_T^k) \leq r_k C(X_T^0), \quad \text{where } X_T^0 = \text{L2O}_{W^k}(X_0^0), \quad (9)$$

with X_T^0 being the initial solution from the L2O model and a constant mapping C . Based on the proof in [38], non-learning GD algorithm’s convergence rate corresponding to the initial L2O state is:

$$F(X_T^0) \leq \frac{\beta}{T} \|X_0^0 - X^*\|_2^2. \quad (10)$$

Given the independence of training iteration k and optimization step T , we align the LHS of Equation (9) with the RHS of Equation (10) by setting $C(X_T^0) = F(X_T^0)$. Given initial point is constant that $X_T^0 = X_T^k$, this yields the combined training convergence rate:

$$F(X_T^k) \leq r_k \frac{\beta}{T} \|X_0^k - X^*\|_2^2. \quad (11)$$

Here, the LHS represents the objective value after k training iterations, while the RHS is a constant term dependent on the initial point X_0 . W.r.t. T , Equation (11) demonstrates a sub-linear convergence rate of at least $\mathcal{O}(1/T)$. The rate indicates that integrating L2O with an existing algorithm via training can enhance its convergence. Such integration is achieved by the Math-L2O framework [23], which utilizes a NN to learn hyperparameters for non-learning algorithms (e.g., step size for GD, step size and momentum for Nesterov Accelerated Gradient [4]).

Further, we construct the Math-L2O training rate r_k (see Equation (9)). Section 4 establishes its linear convergence. Subsequently, Section 5 proposes a deterministic initialization strategy to ensure the alignment ($C(X_T^0) = F(X_T^0)$) and uphold the theoretical conditions for this linear rate.

4 Linear Convergence of L2O Training

In this section, we establish the linear convergence rate for training a Math-L2O model employing an over-parameterized NN, w.r.t. the loss defined in Equation (2). By training the NN (Equation (4)) using GD, we establish its linear convergence rate via NTK theory. Classical NTK theory [16] requires infinite NN width to maintain a non-singular kernel matrix, which facilitates a gradient lower bound akin to the Polyak-Lojasiewicz condition [29, 32]. Applying the relaxation from [29] and the rigorous NN formalizations (Section 2), we demonstrate that an NN width of $\mathcal{O}(Nd)$ is sufficient.

To derive the rate, we first introduce a lemma to bound Math-L2O's gradients. We then prove that appropriate initialization leads to deterministic loss minimization in the initial training iteration. After that, we develop a strategy to maintain this property throughout training, thereby ensuring convergence. This approach culminates in a linear convergence rate for an $\mathcal{O}(Nd)$ -width NN. The main results are summarized herein, with detailed proofs deferred to Appendix A.5 and Appendix A.6.

4.1 Bound Outputs of Math-L2O

We define $\alpha_0 := \sigma_{\min}(G_{L-1,T}^0)$ and let $C_\ell > 0$ for $\ell \in [L]$ be any sequence of positive numbers. Moreover, for $t, j \in [T]$, we define the following quantities:

$$\begin{aligned} \bar{\lambda}_\ell &= \|W_\ell^0\|_2 + C_\ell, \Theta_L = \prod_{\ell=1}^L \bar{\lambda}_\ell, \Phi_j = \|X_0\|_2 + \frac{2j-1}{\beta} \|\mathbf{M}^\top Y\|_2, \\ \Lambda_j &= (1+\beta)\|X_0\|_2^2 + \frac{(4j-3)(1+\beta)+\beta}{\beta} \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 + \frac{(2j-1)(\beta(2j-1)+(2j-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2, \\ S_{\Lambda,T} &= \sum_{t=1}^T \Lambda_t, \quad \delta_1^t = \sum_{s=1}^t \left(\prod_{j=s+1}^t (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \Lambda_s, \\ S_{\bar{\lambda},L} &= \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2}, \quad \delta_2 = \sum_{s=1}^{T-1} \left(\prod_{j=s+1}^{T-1} (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \Lambda_s, \\ \zeta_1 &= \sqrt{\beta} \|X_0\|_2 + (2T+1) \|Y\|_2, \quad \delta_3 = (1+\beta) \|X_0\|_2 + (2T-1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2, \\ \zeta_2 &= \|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2, \quad \delta_4 = \sigma(\delta_3 \Theta_L) (1 - \sigma(\delta_3 \Theta_L)), \end{aligned} \quad (12)$$

where X_0 denotes the initial point, and \mathbf{M} (parameter matrix) and Y (labels) are input features from Equation (2). The defined quantities are positive under the conditions $j \geq 1$ and $\bar{\lambda}_\ell > 0$.

First, we derive a bound for the training gradients by considering them as perturbations from initialization. This bound relates the gradient magnitude to the objective function in Equation (2), as detailed in the following lemma. Despite the derivative for inner layers (Equation (7)) containing an additional term compared to that of the last layer (Equation (8)), a uniform bound as stated applies. The proof is provided in Appendix A.5.4.

Lemma 4.1. *Assuming $\max(\|W_\ell^{k+1}\|_2, \|W_\ell^k\|_2) \leq \bar{\lambda}_\ell$ for $\ell \in [L]$, for any training iteration k , the gradient of the ℓ -th layer parameters W_ℓ^k is bounded by: $\left\| \frac{\partial F}{\partial W_\ell^k} \right\|_2 \leq \frac{\sqrt{\beta} \Theta_L S_{\Lambda,T}}{2\bar{\lambda}_\ell} \|\mathbf{M} X_T^k - Y\|_2$.*

Building upon Lemmas 4.1 and A.6 and auxiliary results (see Appendix A.5), we analyze the dynamics of the final solution X_T w.r.t. parameter updates during training. The subsequent lemma

establishes a rigorous formulation for the fluctuation of X_T in response to changes in parameters between adjacent training iterations. This result demonstrates that Math-L2O, viewed as a function of its learnable parameters, exhibits semi-smoothness, aligning with findings for ReLU-Nets in [29]. The proof is provided in Appendix A.5.3.

The semi-smoothness of the Math-L2O NN is preserved despite its recurrent operations. The coefficient associated with $\|W_\ell^{k+1} - W_\ell^k\|_2$ exhibits $\mathcal{O}(e^{LT})$ scaling, where e is an initialization parameter detailed in Section 5. This represents a looser bound compared to that for ReLU-Nets [29], which is a consequence of Math-L2O’s greater architectural complexity, specifically the T -fold execution of an L -layer NN block (see Equation (8)). However, this scaling behavior is consistent with observations for other deep architectures [2].

Lemma 4.2. *For any training iteration k , assume there exist constants $\bar{\lambda}_\ell \in \mathbb{R}^+$ for $\ell \in [L]$ such that $\max_{k' \in \{k, k+1\}} \|W_\ell^{k'}\|_2 \leq \bar{\lambda}_\ell$. Let X_t^{k+1} and X_t^k be outputs of the Math-L2O (defined in Equations (3) and (4)) corresponding to parameters $W^{k+1} = \{W_\ell^{k+1}\}_{\ell=1}^L$ and $W^k = \{W_\ell^k\}_{\ell=1}^L$, respectively. Then, Math-L2O exhibits the following semi-smoothness property:*

$$\|X_t^{k+1} - X_t^k\|_2 \leq \frac{1}{2} \sum_{s=1}^{t-1} \left(\prod_{j=s+1}^t (1 + (1 + \beta)/2\Theta_L\Phi_j) \right) \Lambda_s \Theta_L \left(\sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right).$$

Lemma 4.2 demonstrates that Math-L2O solutions exhibit a bounded response to perturbations in its NN parameters. This finding, in conjunction with Lemma 4.1, facilitates a more nuanced analysis of the loss dynamics. Further, judicious selection of learning rates enables control over the evolution of NN parameters. Such control is instrumental in bounding the constant quantities from these lemmas, thereby establishing the desired convergence rate presented in the subsequent theorem.

4.2 Linear Training Convergence Rate of Math-L2O

Leveraging the bounds on Math-L2O’s output (Lemma A.6) and its gradient (Lemma 4.1), the following theorem establishes the linear convergence rate for training the Math-L2O model. The proof is provided in Appendix A.6.

Theorem 4.3. *Consider the NN defined in Equation (4), using quantities from Equation (12), suppose the following conditions hold at initialization:*

$$\alpha_0 \geq 8(1 + \beta)\zeta_2, \quad (13a) \quad \alpha_0^2 \geq \frac{\beta^3}{4\beta_0^2} \delta_4^{-2} \left(-\frac{1}{2} \Theta_{L-1}^2 \Lambda_T S_{\Lambda, T-1} + \Theta_L^2 (\Lambda_T + \delta_2) S_{\bar{\lambda}, L} S_{\Lambda, T} \right). \quad (13b)$$

$$\alpha_0^2 \geq \max_{\ell \in [L]} \frac{\Theta_L}{C_\ell \bar{\lambda}_\ell} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_4^{-2} \zeta_1 S_{\Lambda, T}, \quad (13c) \quad \alpha_0^3 \geq \frac{(1+\beta)\beta^2 \sqrt{\beta}}{2\beta_0^2} \delta_4^{-2} \Theta_L \Theta_{L-1} \zeta_1 \zeta_2 S_{\bar{\lambda}, L} S_{\Lambda, T}, \quad (13d)$$

Let the learning rate η satisfy:

$$\eta < \frac{8}{\beta} (\delta_2 + \Lambda_T) (\delta_2 + \Theta_L S_{\Lambda, T} S_{\bar{\lambda}, L})^{-1} S_{\Lambda, T}^{-2}, \quad (14a) \quad \eta < \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_4^{-2} \alpha_0^{-2}. \quad (14b)$$

Then, for weights $W^k = \{W_\ell^k\}_{\ell=1}^L$ at training iteration k , the loss function $F(W^k)$ converges linearly to a global minimum:

$$F(W^k) \leq (1 - 4\eta \frac{\beta_0^2}{\beta^2} \delta_4 \alpha_0^2)^k F(W^0).$$

Remark 1. $(1 - 4\eta \frac{\beta_0^2}{\beta^2} \delta_4 \alpha_0^2)^k$ is r_k in Equation (11), which is a less than one term since $\delta_4 = \sigma(\delta_3 \Theta_L)(1 - \sigma(\delta_3 \Theta_L)) > 0$ and $\alpha_0 := \sigma_{\min}(G_{L-1, T}^0) > 0$ ($G_{L-1, T}^0$ is a thin matrix), which ensure that the L2O converges at least as fast as GD.

Equations (14a) and (14b) are based on the quantities defined in Equation (12). Each quantity represents an inner formulation in the demonstration of lemmas and theorems. We use these quantities to simplify the formulations. The conditions specified in Equation (13) impose additional lower bounds on α_0 , the minimal singular value of the $(L-1)$ -th layer’s inner output. The bounds stipulated in Equations (13b) to (13d) are influenced by both the network depth L and the number of gradient descent (GD) iterations T . In contrast, the constraint in Equation (13a) primarily depends on T . An initialization strategy ensuring these conditions are met is proposed in Section 5. We provide a detailed interpretation in Appendix A.2.

4.3 Analysis of Learning Rate Magnitude

The bounds in Equations (14a) and (14b) indicate that the learning rate η diminishes as L and T increase. We argue that this requirement for a small η is not a significant limitation; it is consistent with the NTK framework, which does not rely on large learning rates for convergence. To quantify

this, we examine the scaling of η relative to T , L , and $\bar{\lambda}_{\max}$. Here, $\bar{\lambda}_{\max} = \max\{\bar{\lambda}_\ell\}, \ell \in [L]$ is the maximum constant upper bound on the singular values of the NN layers (Equation (12)). These bounds are parameters that can be directly influenced by the choice of initialization method.

First, analyzing Equation (14a), we derive the scaling of η as $\mathcal{O}(\frac{T\bar{\lambda}_{\max}^{LT} + T^2}{((T\bar{\lambda}_{\max}^{LT}) + \bar{\lambda}_{\max}^L T^3 L \bar{\lambda}_{\max}^{-2}) T^6})$, where constant factors independent of T and L are omitted. This expression highlights that the magnitude of η is strongly dependent on the bound $\bar{\lambda}_{\max}$. This dependence implies that the learning rate can be prevented from becoming extremely small by using a proper initialization method (such as our proposed method in Section 5) to control $\bar{\lambda}_{\max}$.

Moreover, Equation (14b) shows that η 's magnitude is highly correlated with the lower bound of α_0 (the penultimate layer's singular value, per Section 4.1). Given the four distinct lower bounds for α_0 derived in Equation (13), we now formulate the magnitude of η for each respective case. First, if Equation (13c) holds, $\eta = \mathcal{O}(\exp(2T\bar{\lambda}_{\max}^L)T^{-2})$, which is a non-restrictive bound due to the exponential term. Second, if Equation (13d) holds, $\eta = \mathcal{O}((\bar{\lambda}_{\max}^{2L} T^4 + \bar{\lambda}_{\max}^{2L} (T + T\bar{\lambda}_{\max}^{LT}) L \bar{\lambda}_{\max}^{-2} T^3)^{-1})$. This scales inversely with $\bar{\lambda}_{\max}$ and exponentially with L and T . Third, if Equation (13b) holds, $\eta = \mathcal{O}(\bar{\lambda}_{\max}^{-L} T^{-3})$, which also scales inversely with $\bar{\lambda}_{\max}$ and exponentially with L . Finally, if Equation (13a) holds, $\eta = \mathcal{O}(\exp(\frac{2}{3}T\bar{\lambda}_{\max}^L)(\bar{\lambda}_{\max}^{2L} T^2 L \bar{\lambda}_{\max}^{-2})^{-\frac{2}{3}})$, which, similar to the first case, is a non-restrictive bound due to the exponential term.

The foregoing results indicate that a larger $\bar{\lambda}_{\max}$ correlates with a smaller learning rate η . Nevertheless, this does not result in a degradation of convergence speed. This conclusion is supported by two observations: *Theoretical Consistency*: The requirement for a small η is permissible under NTK theory [16]. The NTK regime assumes infinitely wide networks, where convergence is achieved within a compact space around the initialization, thus obviating the need for large learning steps. *Empirical Insensitivity*: Our experimental results demonstrate that the convergence speed is robust to the learning rate. As depicted in Figure 4a, our method achieves similar convergence rates for η across a wide range (e.g., 10^{-3} to 10^{-7}).

Adopting a small learning rate is a pragmatic trade-off to avoid the requirement for an extremely wide NN. Existing analyses [3, 29] that remove the infinite-width assumption often impose a polynomial width dependency (e.g., $\mathcal{O}(N^3)$) on the sample size N . In our framework (Section 2), the coordinate-wise L2O treats d -dimensional features as independent inputs, leading to an effective sample size of Nd . A polynomial dependency on Nd would be impractical. Therefore, we opt for the alternative constraint of a smaller learning rate, which permits a feasible network width.

5 Deterministic Initialization

This section introduces an initialization strategy ensuring the alignment between Math-L2O and GD (see Section 3) while also satisfying the conditions presented in Section 4. The proposed initialization strategy first establishes Math-L2O to operate as a standard GD algorithm, and then guarantees the uniform convergence of Math-L2O throughout subsequent training iterations.

5.1 Initialization for Alignment

Following methodology in [29], we let $C_\ell = 1$ for $\ell \in [L]$. For parameter matrices initialization W (see Section 2), we randomly initialize parameter matrices of first $L - 1$ layers, i.e., $\{W_1^0, \dots, W_{L-1}^0\}$ from a standard Gaussian distribution and set the last layer's parameter matrix $W_L^0 = \mathbf{0}$. Through the 2σ activation detailed in Equation (4), it outputs a constant step size, i.e., $P_T = \mathbf{I}$. Consequently, the learning proceeds with a uniform step size of $1/\beta$ after initialization, emulating standard GD and its typical sub-linear convergence rate [38]. Moreover, this zero-initialization of W_L^0 ensures that initial gradients for the inner layers are all zero (as shown in Equation (7)), which serves to mitigate gradient explosion.

The condition $\alpha_0 > 0$ (see Theorem 4.3) is fulfilled by randomly sampling the initial weight matrices $\{W_k^0\}_{k=1}^{L-1}$ from a standard Gaussian distribution. This approach generally ensures full row rank for fat matrices (more columns than rows) [37]. Each matrix W_k^0 then undergoes QR decomposition. Non-negativity is subsequently enforced upon the elements of the resulting upper triangular factor (e.g., via its element-wise absolute value, achieved in PyTorch using its `sign` function).

5.2 Enhancing Singular Values for Linear Convergence of Training

Motivated by properties of minimal singular values in ReLU-Nets identified in [29], we analyze the order-gap for α_0 between the left-hand side (LHS) and right-hand side (RHS) of the inequalities in Equation (13). To satisfy these inequalities, we propose increasing α_0 . This is achieved by applying a constant *expansion coefficient* $e \geq 1$ to the initial NN parameters $\{W_1^0, \dots, W_{L-1}^0\}$, transforming them to $\{eW_1^0, \dots, eW_{L-1}^0\}$. This parameter expansion scales the minimal singular value α_0 to $e^{L-1}\alpha_0$, reflecting the cumulative impact across $L - 1$ layers. However, other terms on the RHS of Equation (13) also depend on e . We then establish four lemmas to demonstrate that the conditions for linear convergence, as specified in Theorem 4.3, are met for an appropriately chosen value of e .

First, we set the initial point to the origin, $X_0 = \mathbf{0}$, a choice commonly adopted in L2O literature [23, 34]. Then, with $C_\ell = 1$ for $\ell \in [L]$, we present four lemmas demonstrating that the conditions for linear convergence (see Theorem 4.3) are satisfied for an appropriately chosen constant e . The lemmas indicate that a larger e is required as the number of optimization steps (T) increases. Specifically, Lemma 5.2 establishes that e scales exponentially with T . Conversely, increasing the network depth (L) alleviates the need for a large e . The proofs are provided in Appendix B.

Lemma 5.1. *Assuming $X_0 = \mathbf{0}$, if $e = \Omega(T^{\frac{1}{L-1}})$, then the inequality Equation (13a) holds.*

Lemma 5.2. *If $e = \Omega(T^{\frac{3T+6}{TL-T-4L+6}})$, then the inequality Equation (13b) holds.*

Lemma 5.3. *Assuming $X_0 = \mathbf{0}$, if $e = \Omega(T^{\frac{4}{L-1}})$, then the inequality Equation (13c) holds.*

Lemma 5.4. *Assuming $X_0 = \mathbf{0}$, if $e = \Omega(T^{\frac{5}{L-1}} L^{\frac{1}{L-1}})$, then the inequality Equation (13d) holds.*

6 Empirical Evaluation

This section presents an empirical evaluation of the framework proposed in Section 3 and the theoretical results from Section 4. Experiments are conducted using Python 3.9 and PyTorch 1.12.0 on an Ubuntu 20.04 system equipped with 128GB of RAM and two NVIDIA RTX 3090 GPUs.

Data Generation. Due to GPU memory constraints, vectors $X \in \mathbb{R}^{5120 \times 1}$ and $Y \in \mathbb{R}^{4000 \times 1}$ for Equation (2) are generated by sampling from a standard Gaussian distribution. These represent ten problem instances with respective dimensional components of 512 (for X) and 400 (for Y). Following the coordinate-wise approach in [23], we formed an input feature matrix of 5120×2 . This setup is equivalent to a training batch of 5120 two-feature samples.

Math-L2O Model Architecture. The Math-L2O model is configured with $T = 100$ optimization steps (Equation (2)). Its architecture comprises a $L = 3$ -layer DNN, as formulated in Equation (4). The first layer has an output dimension of 2. To ensure over-parameterization, the $(L - 1)$ -th (i.e., second) layer’s output dimension is set to $512 \times 10 = 5120$. The final layer produces a scalar output (dimension 1). Three specific model configurations are designed for ablation studies, foundational experiments, and robustness evaluations. These are detailed in Appendix C.1.

Training and Initialization Configurations. L2O models are trained using the Stochastic Gradient Descent (SGD) optimizer. For the $L = 3$ -layer network configuration, parameters for the initial two layers ($l = 1, 2$) are initialized according to the methodology presented in Section 5.1, while parameters for the final layer ($l = 3$) are zero-initialized.

6.1 Training Performance

We evaluated the mean training loss in Equation (2) across all samples. Figure 4a illustrates this loss at $T = 100$, benchmarked against the standard GD objective (black dashed line). The results demonstrate that Math-L2O consistently achieves fast training convergence, corroborating the theoretical linear convergence established in Theorem 4.3.

Further, we investigated the robustness of our proposed L2O method to variations in optimization steps and learning rates (LRs). Models corresponding to different step/LR configurations are trained for 400 epochs. Figure 4b presents the training objectives for these configurations, benchmarked against standard GD (black dashed line). In contrast to the instability observed for Math-L2O [23]

and LISTA-CPSS [7] under certain settings (Figure 2), the consistent convergence across all tested configurations in Figure 4b demonstrates the robustness of our proposed L2O approach.

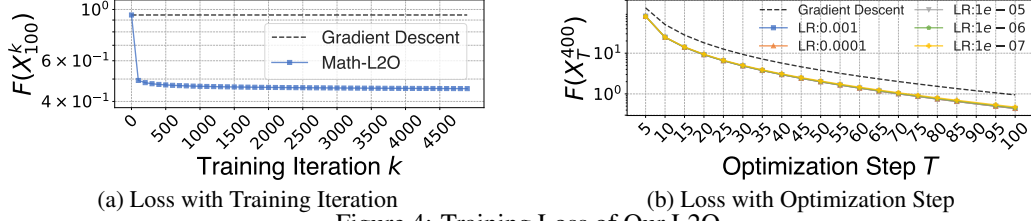


Figure 4: Training Loss of Our L2O

Moreover, we evaluate the inference performance of our framework against baseline methods. Experimental results (in Appendix C.4) demonstrate the framework’s robustness to hyperparameters.

6.2 Ablation Studies for Learning Rate η and Expansion Coefficient e

We conduct ablation studies to assess the impact of the LR η , theoretically bounded in Equations (14a) and (14b) (Theorem 4.3), and the initialization coefficient e , defined in Section 5. The experimental configuration employs $T = 20$, input $X \in \mathbb{R}^{32 \times 32}$, output $Y \in \mathbb{R}^{32 \times 20}$, and a neural network width of 1024. Performance is measured by the relative improvement of the proposed L2O method over standard GD at iteration $T = 20$, calculated as $\frac{\text{obj}_{\text{GD}} - \text{obj}_{\text{L2O}}}{\text{obj}_{\text{GD}}}$. These studies further validate Corollary C.1, which establishes an inverse relationship between the viable LR η and the coefficient e , implying that a larger e necessitates a smaller η to ensure convergence.

With the initialization coefficient fixed at $e = 50$, we evaluate the impact of varying the LR η on the relative objective improvement. The results in Figure 5a demonstrate that while LR’s such as 10^{-4} and smaller achieve convergence, $\eta = 10^{-3}$ leads to unstable behavior or divergence. This finding empirically supports the existence of an operational upper bound on the LR, consistent with the theoretical constraints outlined in Equations (14a) and (14b). Moreover, reducing the LR below this stability threshold results in slower convergence rates. This observation aligns with the implication of Theorem 4.3 that, under the specified conditions, larger permissible LR’s yield faster convergence.

Fixing the LR at $\eta = 10^{-7}$, we examine the influence of the initialization coefficient e on performance. The results, presented in Figure 5b, demonstrate that the relative objective improvement consistently increases with larger values of e . Additional results exploring different e and LR combinations are deferred to Appendix C owing to space constraints. These findings validate the proposed strategies for selecting the initialization coefficient and learning rate.

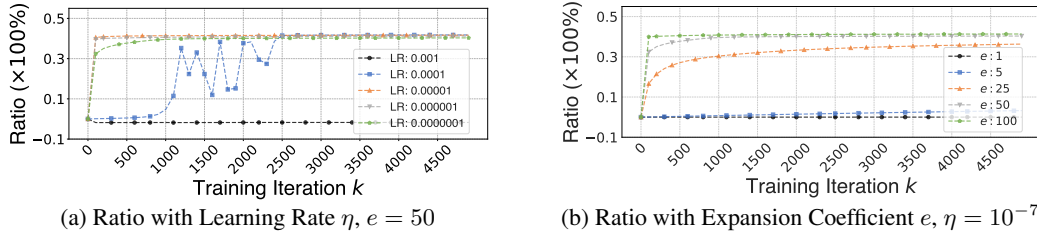


Figure 5: Ablation Studies of Improve Ratio to Learning Rate and Expansion Coefficient

7 Conclusion

This work analyzes a Learning-to-Optimize (L2O) framework that accelerates Gradient Descent (GD) through adaptive step-size learning. We theoretically prove that the L2O training enhances GD’s convergence rate by linking network training bounds to GD’s performance. Leveraging Neural Tangent Kernel (NTK) theory and the over-parameterization scheme via wide layers, we establish convergence guarantees for the complete L2O system. A principled initialization strategy is introduced to satisfy the theoretical requirements for these guarantees. Empirical results across various optimization problems validate our theory and demonstrate substantial practical efficacy.

Acknowledgements

This work is supported in part by funding from CUHK (4937007, 4937008, 5501329, 5501517) and Science and Technology Project of State Grid Corporation of China (No.5700-202440239A-1-1-ZN).

References

- [1] AF Agarap. Deep Learning Using Rectified Linear Units (ReLU). *arXiv preprint arXiv:1803.08375*, 2018.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A Convergence Theory for Deep Learning via Over-Parameterization. In *International conference on machine learning*, pages 242–252. PMLR, 2019.
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the Convergence Rate of Training Recurrent Neural Networks. *Advances in neural information processing systems*, 32, 2019.
- [4] Amir Beck and Marc Teboulle. A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [5] Yanmei Cao, Guomei Zhang, Guobing Li, and Jia Zhang. A Deep Q-Network Based-Resource Allocation Scheme for Massive MIMO-NOMA. *IEEE Communications Letters*, 25(5):1544–1548, 2021.
- [6] Tianlong Chen, Xiaohan Chen, Wuyang Chen, Zhangyang Wang, Howard Heaton, Jialin Liu, and Wotao Yin. Learning to Optimize: A Primer and a Benchmark. *The Journal of Machine Learning Research*, 23(1):8562–8620, 2022.
- [7] Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Theoretical Linear Convergence of Unfolded ISTA and Its Practical Weights and Thresholds. *Advances in Neural Information Processing Systems*, 31, 2018.
- [8] Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Hyperparameter Tuning is All You Need for LISTA. *Advances in Neural Information Processing Systems*, 34:11678–11689, 2021.
- [9] Chizat, Lenaïc and Bach, Francis. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *COLT*, 2020.
- [10] Kingma Diederik. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [11] Simon S Du, Xiyu Zhai, Barnabas Póczos, and Aarti Singh. Gradient Descent Provably Optimizes Over-Parameterized Neural Networks. *arXiv preprint arXiv:1810.02054*, 2018.
- [12] Karol Gregor and Yann LeCun. Learning Fast Approximations of Sparse Coding. In *Proceedings of the 27th international conference on machine learning*, pages 399–406, 2010.
- [13] Stephen Grossberg. Recurrent neural networks. *Scholarpedia*, 8(2):1888, 2013.
- [14] Howard Heaton, Xiaohan Chen, Zhangyang Wang, and Wotao Yin. Safeguarded Learned Convex Optimization. In *AAAI*, pages 7848–7855, 2023.
- [15] Qiyu Hu, Yunlong Cai, Qingjiang Shi, Kaidi Xu, Guanding Yu, and Zhi Ding. Iterative Algorithm Induced Deep-Unfolding Neural Networks: Precoding Design for Multiuser MIMO Systems. *IEEE TWC*, 20(2):1394–1410, 2020.
- [16] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [17] Sungyoon Kim and Mert Pilanci. Convex Relaxations of ReLU Neural Networks Approximate Global Optima in Polynomial Time. *arXiv preprint arXiv:2402.03625*, 2024.

- [18] Sungyoon Kim, Aaron Mishkin, and Mert Pilanci. Exploring the Loss Landscape of Regularized Neural Networks via Convex Duality. *arXiv preprint arXiv:2411.07729*, 2024.
- [19] Timothy P Lillicrap and Adam Santoro. Backpropagation Through Time and the Brain. *Current Opinion in Neurobiology*, 55:82–89, 2019.
- [20] Wei Lin, Qingyu Song, and Hong Xu. Adaptive Coordinate-Wise Step Sizes for Quasi-Newton Methods: A Learning-to-Optimize Approach. *arXiv preprint arXiv:2412.00059*, 2024.
- [21] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss Landscapes and Optimization in Over-Parameterized Non-Linear Systems and Neural Networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- [22] Jialin Liu and Xiaohan Chen. ALISTA: Analytic Weights Are As Good As Learned Weights in LISTA. In *International Conference on Learning Representations (ICLR)*, 2019.
- [23] Jialin Liu, Xiaohan Chen, Zhangyang Wang, Wotao Yin, and HanQin Cai. Towards Constituting Mathematical Structures for Learning to Optimize. In *International Conference on Machine Learning*, pages 21426–21449. PMLR, 2023.
- [24] Xin Liu, Zhisong Pan, and Wei Tao. Provable Convergence of Nesterov’s Accelerated Gradient Method for Over-Parameterized Neural Networks. *Knowledge-Based Systems*, 251:109277, 2022.
- [25] Kaifeng Lv, Shunhua Jiang, and Jian Li. Learning Gradient Descent: Better Generalization and Longer Horizons. In *ICML*, pages 2247–2255. PMLR, 2017.
- [26] Michael C Mozer. A Focused Backpropagation Algorithm for Temporal Pattern Recognition. In *Backpropagation*, pages 137–169. Psychology Press, 2013.
- [27] Sridhar Narayan. The Generalized Sigmoid Activation Function: Competitive Supervised Learning. *Information sciences*, 99(1-2):69–82, 1997.
- [28] Nathanson, Melvyn B. *Weyl’s inequality*, pages 97–119. Springer New York, New York, NY, 1996. ISBN 978-1-4757-3845-2. doi: 10.1007/978-1-4757-3845-2_4.
- [29] Quynh Nguyen. On the Proof of Global Convergence of Gradient Descent for Deep ReLU Networks with Linear Widths. In *International Conference on Machine Learning*, pages 8056–8062. PMLR, 2021.
- [30] Quynh N Nguyen and Marco Mondelli. Global Convergence of Deep Networks with One Wide Layer Followed by Pyramidal Topology. *Advances in Neural Information Processing Systems*, 33:11961–11972, 2020.
- [31] Mert Pilanci. From Complexity to Clarity: Analytical Expressions of Deep Neural Network Weights via Clifford Algebra and Convexity. *Transactions on Machine Learning Research*, 2024.
- [32] Boris Teodorovich Polyak. Minimization of Unsmooth Functionals. *USSR Computational Mathematics and Mathematical Physics*, 9(3):14–29, 1969.
- [33] Yifei Shen, Yuanming Shi, Jun Zhang, and Khaled B. Letaief. Graph Neural Networks for Scalable Radio Resource Management: Architecture Design and Theoretical Analysis. *IEEE JSAC*, 39(1):101–115, 2021.
- [34] Qingyu Song, Wei Lin, Juncheng Wang, and Hong Xu. Towards Robust Learning to Optimize with Theoretical Guarantees. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27498–27506, 2024.
- [35] Qingyu Song, Juncheng Wang, Jingzong Li, Guochen Liu, and Hong Xu. A Learning-only Method for Multi-Cell Multi-User MIMO Sum Rate Maximization. In *International Conference on Computer Communications*. IEEE, 2024.

- [36] Haoran Sun, Xiangyi Chen, Qingjiang Shi, Mingyi Hong, Xiao Fu, and Nicholas D Sidiropoulos. Learning to Optimize: Training Deep Neural Networks for Interference Management. *IEEE TSP*, 66(20):5438–5453, 2018.
- [37] Terence Tao. *Topics in Random Matrix Theory*, volume 132. American Mathematical Soc., 2012.
- [38] Ryan Tibshirani. 10-725: Convex Optimization — lecture 6: Gradient descent, convergence analysis, and subgradients. Course lecture notes, Carnegie Mellon University, September 2013. URL <https://www.stat.cmu.edu/~ryantibs/convexopt-F13/scribes/lec6.pdf>. Scribed by Micol Marchetti-Bowick. Lecture date: 12 Sep 2013. Accessed 11 May 2025.
- [39] Junjie Yang, Tianlong Chen, Mingkang Zhu, Fengxiang He, Dacheng Tao, Yingbin Liang, and Zhangyang Wang. Learning to Generalize Provably in Learning to Optimize. In *International Conference on Artificial Intelligence and Statistics*, pages 9807–9825, 2023.
- [40] Yu, Zixiong and Tian, Songtao and Chen, Guhan. Divergence of empirical neural tangent kernel in classification problems. In *ICLR*, 2025.
- [41] Yu Zhao, Ignas G. Niemegheers, and Sonia M. Heemstra De Groot. Dynamic Power Allocation for Cell-Free Massive MIMO: Deep Reinforcement Learning Methods. *IEEE Access*, 9:102953–102965, 2021.

A Appendix

A.1 Details for Definitions

General L2O. Given X_0 , we have the following L2O update with NN g to generate X_T :

$$X_t = X_{t-1} + g(W_1, W_2, \dots, W_L, X_{t-1}, \nabla F(X_{t-1})), t \in [T]. \quad (15)$$

Concatenation of N Problems. For $t \in [T]$, we make the following denotations to represent the concatenation of N samples (each is a unique optimization problem):

$$\mathbf{M} := \begin{bmatrix} \mathbf{M}_1 & & \\ & \dots & \\ & & \mathbf{M}_N \end{bmatrix}, X_t := [x_{1,t}^\top | x_{2,t}^\top | \dots | x_{N,t}^\top]^\top, Y := [y_1^\top | y_2^\top | \dots | y_N^\top]^\top.$$

X_t and Y are still column vectors since we take the coordinate-wise setting from [23].

A.2 Detailed Interpretation of Quantities in Theorem 4.3

We elaborate on the quantities introduced in Theorem 4.3. Our notational convention is as follows: subscripts T and k identify constant terms that are dependent on the total steps T and the training iteration k . Conversely, indices j and t (appearing as superscripts or subscripts) are used to reference scalar-valued functions at a specific step t or index j .

- $\bar{\lambda}_\ell$ is a positive constant upper bound for each ℓ -th layer in NN g_W (see Section 2), which is constructed in the proof of Theorem 4.3.
- Θ_L is a positive constant w.r.t. $\bar{\lambda}_\ell$.
- Denote $\bar{\lambda}_{\min}, \bar{\lambda}_{\max} := \min\{\bar{\lambda}_\ell\}, \max\{\bar{\lambda}_\ell\}, \ell \in [L]$, Θ_L is lower and upper bounded by $\Omega(\bar{\lambda}_{\min}^L)$ and $\mathcal{O}(\bar{\lambda}_{\max}^L)$, respectively. Moreover, Θ_L^{-1} is $\Omega(\bar{\lambda}_{\max}^{-L})$ and $\mathcal{O}(\bar{\lambda}_{\min}^{-L})$.
- Φ_j is a scalar-valued function w.r.t. step j . The constant coefficients are given by initial point X_0 , coefficient matrix \mathbf{M} , and coefficient vector Y from problem defined in Equation (2). β is the smoothness extent of objective. We use two denotations, j and t , for step, which are used to formulate different computations in formulations. This formulation is derived by the upper bound relaxation of L_2 -norm of gradient at X_0 . Φ_j is $\mathcal{O}(j)$ and $\Omega(j)$.
- Λ_j is a scalar-valued function w.r.t. step j , which is identical to those in Φ_j . Λ_j is $\mathcal{O}(j^2)$ and $\Omega(j^2)$.
- $S_{\Lambda,T}$ and $S_{\bar{\lambda},L}$ are positive constants, which represents the summation of Λ of T steps and summation of $\bar{\lambda}$ of L -th NN layers, respectively.
- $S_{\Lambda,T}$ is used in the demonstration for Lemma 4.1 (bound of gradient of NN training), line 625, page 22. The proof is achieved by upper bound relaxation of L_2 -norm. $S_{\bar{\lambda},L}$ is used in Theorem 4.3 and related auxiliary lemmas. $S_{\Lambda,T}$ is $\mathcal{O}(T^3)$ and $\Omega(T^3)$. $S_{\bar{\lambda},L}^{-1}$ is $\mathcal{O}(T^{-3})$ and $\Omega(T^{-3})$. Denote $\bar{\lambda}_{\min}, \bar{\lambda}_{\max} = \min\{\bar{\lambda}_\ell\}, \max\{\bar{\lambda}_\ell\}, \ell \in [L]$, $S_{\bar{\lambda},L}$ is $\Omega(L\bar{\lambda}_{\max}^{-2})$ and $\mathcal{O}(L\bar{\lambda}_{\min}^{-2})$. Moreover, $S_{\bar{\lambda},L}^{-1}$ is $\Omega(L^{-1}\bar{\lambda}_{\max}^2)$ and $\mathcal{O}(L^{-1}\bar{\lambda}_{\min}^2)$.
- ζ_1 and ζ_2 are two positive constants scale linearly w.r.t., X_0, \mathbf{M}, Y, T , and β . ζ_1 and ζ_2 are both $\Omega(T)$ and $\mathcal{O}(T)$.
- ζ_1 and ζ_2 : Two positive constants scale linearly w.r.t., X_0, \mathbf{M}, Y, T , and β . ζ_1 and ζ_2 are both $\Omega(T)$ and $\mathcal{O}(T)$.
- δ_1^t : A scalar-valued function w.r.t. step t . The constant coefficients are Θ_L, Φ_j , and Λ_s , where s denotes an step. Denote $\bar{\lambda}_{\min}, \bar{\lambda}_{\max} = \min\{\bar{\lambda}_\ell\}, \max\{\bar{\lambda}_\ell\}, \ell \in [L]$, δ_1^t is $\Omega(t\bar{\lambda}_{\min}^{Lt})$ and $\mathcal{O}(t\bar{\lambda}_{\max}^{Lt})$.
- δ_2 : Positive constant scales with T . Denote $\bar{\lambda}_{\min}, \bar{\lambda}_{\max} = \min\{\bar{\lambda}_\ell\}, \max\{\bar{\lambda}_\ell\}, \ell \in [L]$, δ_2 is $\Omega(T\bar{\lambda}_{\min}^{LT})$ and $\mathcal{O}(T\bar{\lambda}_{\max}^{LT})$. Moreover, δ_2^{-1} is $\Omega(T\bar{\lambda}_{\max}^{-LT})$ and $\mathcal{O}(T\bar{\lambda}_{\min}^{-LT})$.
- δ_3 : Positive constant scales linearly w.r.t., X_0, \mathbf{M}, Y, T , and β . δ_3 is both $\Omega(T)$ and $\mathcal{O}(T)$.
- δ_4 : Denote $\bar{\lambda}_{\min}, \bar{\lambda}_{\max} = \min\{\bar{\lambda}_\ell\}, \max\{\bar{\lambda}_\ell\}, \ell \in [L]$, δ_4 is $\Omega(\exp(-T\bar{\lambda}_{\max}^L))$ and $\mathcal{O}(\exp(-T\bar{\lambda}_{\min}^L))$. Moreover, δ_4^{-1} is $\mathcal{O}(\exp(T\bar{\lambda}_{\max}^L))$ and $\Omega(\exp(T\bar{\lambda}_{\min}^L))$.

A.3 Derivative of General L2O

In this section, we derive a general framework for any L2O models by the chain rule, which gives us a complete workflow of each component in the derivatives within the chain. Then, we apply it to the Math-L2O framework [23] to get the formulation for the L2O model defined in Equation (4).

Due to the chain rule, we derive the following general formulation of the derivative in L2O model:

$$\frac{\partial F(X_T)}{\partial W_\ell} = \frac{\partial F(X_T)}{\partial X_T} \left(\frac{\partial X_T}{\partial X_{T-1}} \frac{\partial X_{T-1}}{\partial W_\ell} + \frac{\partial X_T}{\partial G_{L,t}} \frac{\partial G_{L,t}}{\partial W_\ell} \right).$$

We then calculate each term in the right-hand side (RHS) in the above formulation. First, we calculate $\frac{\partial X_{T-1}}{\partial W_\ell}$ as:

$$\frac{\partial X_{T-1}}{\partial W_\ell} = \frac{\partial X_{T-1}}{\partial X_{T-2}} \frac{\partial X_{T-2}}{\partial W_\ell} + \frac{\partial X_{T-1}}{\partial G_{L,T-1}} \frac{\partial G_{L,T-1}}{\partial W_\ell}.$$

Thus, we can iteratively derive the gradient until X_1 . After rearranging terms, we have the following complete formulation of $\frac{\partial F}{\partial W_\ell}$:

$$\frac{\partial F(X_T)}{\partial W_\ell} = \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{t=1}^T \left(\prod_{j=T}^{t+1} \frac{\partial X_j}{\partial X_{j-1}} \right) \frac{\partial X_t}{\partial G_{L,t}} \frac{\partial G_{L,t}}{\partial W_\ell} \right). \quad (16)$$

We note that $\frac{\partial X_j}{\partial X_{j-1}}$ relies on different implementations. For example, for general L2O model that the update in each step is directly the output of neural networks (NNs), we have $\frac{\partial X_j}{\partial X_{j-1}} := \mathbf{I} + \frac{\partial G_{L,j}}{\partial X_{j-1}}$. Then, Equation (16) is derived by:

$$\frac{\partial F}{\partial W_\ell} = \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{t=1}^T \left(\prod_{j=T}^{t+1} \left(\mathbf{I} + \frac{\partial G_{L,j}}{\partial X_{j-1}} \right) \right) \frac{\partial X_t}{\partial G_{L,t}} \frac{\partial G_{L,t}}{\partial W_\ell} \right). \quad (17)$$

$\frac{\partial G_{L,j}}{\partial X_{j-1}}$ depends on specific implementation of NNs. Liu et al. [23] simplify $\frac{\partial G_{L,j}}{\partial X_{j-1}}$ by detaching input tensor from the back-propagation process, which truncate the branches in the chain from $F(X_T)$ to W_ℓ . The detaching operation yields simpler $\frac{\partial X_j}{\partial X_{j-1}}$. As will be introduced in the following sections, $\frac{\partial X_j}{\partial X_{j-1}}$ depends only on NN's output.

Further, the definition of $\frac{\partial X_T}{\partial G_{L,t}}$ is framework-dependent. In the general L2O model, $\frac{\partial X_T}{\partial G_{L,t}} := \mathbf{I}$, whereas in Math-L2O [23], it is defined based on the FISTA algorithm [4]. Subsequently, we perform a layer-by-layer computation for each derivative $\frac{\partial G_{L,j}}{\partial X_{t-1}}$ and $\frac{\partial G_{L,t}}{\partial W_\ell}$.

First, we derive $\frac{\partial G_{L,t}}{\partial G_{L-1,t}}$ by:

$$\frac{\partial G_{L,t}}{\partial G_{L-1,t}} = \begin{cases} \nabla \text{ReLU}(G_{L-1,t}) W_\ell & \ell \in [L-1], \\ \nabla 2\sigma(G_{\ell,t}) W_\ell & \ell = L. \end{cases}$$

For simplification, we use ∇ReLU and $\nabla 2\sigma$ to represent derivatives $\nabla \text{ReLU}(G_{L-1,t})$ and $\nabla 2\sigma(G_{\ell,t})$, respectively, which are corresponding diagonal matrices of coordinate-wise activation function's derivatives. Next, $\frac{\partial G_{L,t}}{\partial X_{t-1}}$ is given by:

$$\frac{\partial G_{L,j}}{\partial X_{T-1}} = \left(\prod_{\ell=L}^2 \frac{\partial G_{\ell,j}}{\partial G_{\ell-1,j}} \right) \frac{\partial G_{1,j-1}}{\partial X_{T-1}} = \nabla 2\sigma w_L \left(\prod_{\ell=L-1}^2 \nabla \text{ReLU} W_\ell \right) [\mathbf{I}, \mathbf{H}^\top], \quad (18)$$

where $\mathbf{H} := \mathbf{M}^\top \mathbf{M}$ denotes the Hessian matrix of the loss function in Equation (2).

Second, $\frac{\partial G_{L,t}}{\partial W_\ell}$ is given by:

$$\begin{aligned} \frac{\partial G_{L,t}}{\partial W_\ell} &= \left(\prod_{j=L}^{\ell+1} \frac{\partial G_{j,t}}{\partial G_{j-1,t}} \right) \frac{\partial G_{1,t}}{\partial W_\ell} \\ &= \begin{cases} \nabla 2\sigma w_L \left(\prod_{j=L-1}^{\ell+1} \nabla \text{ReLU} W_j \right) \nabla \text{ReLU}(\mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^\top) & \ell \in [L-1], \\ \nabla 2\sigma(\mathbf{I}_{n_\ell} \otimes G_{L-1,t}^\top) & \ell = L, \end{cases} \end{aligned} \quad (19)$$

where $\mathbf{I}_{n_\ell} \in \mathbb{R}^{n_\ell \times n_\ell}$, \otimes denotes Kronecker Product, and $\mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^\top \in \mathbb{R}^{n_\ell \times n_\ell n_{\ell-1}}$.

Substituting Equation (18) and Equation (19) into Equation (17) yields following final derivative formulation of general L2O model:

$$\begin{aligned}
& \frac{\partial F}{\partial W_\ell} \\
&= \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{t=1}^T \left(\prod_{j=T}^{t+1} \left(\mathbf{I} + \frac{\partial G_{L,j}}{\partial X_{j-1}} \right) \right) \frac{\partial X_T}{\partial G_{L,t}} \frac{\partial G_{L,t}}{\partial W_\ell} \right), \\
&= \begin{cases} \mathbf{K}_{n_\ell, n_{\ell-1}} \left(\left(X_T^k \mathbf{M}^\top - Y^\top \right) \mathbf{M} \right. \\ \quad \left(\sum_{t=1}^T \left(\mathbf{I} + \nabla 2\sigma w_L \left(\prod_{\ell=L-1}^2 \nabla \text{ReLU} W_\ell \right) [\mathbf{I}, \mathbf{H}^\top] \right)^{T-t} \right. \\ \quad \left. \left. \nabla 2\sigma w_L^\top \left(\prod_{j=L-1}^{\ell+1} \nabla \text{ReLU} W_j \right) \nabla \text{ReLU} (\mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^\top) \right) \right)^\top & \ell \in [L-1], \\ \mathbf{K}_{n_\ell, n_{\ell-1}} \left(\left(X_T^k \mathbf{M}^\top - Y^\top \right) \mathbf{M} \right. \\ \quad \left(\sum_{t=1}^T \left(\mathbf{I} + \nabla 2\sigma w_L \left(\prod_{\ell=L-1}^2 \nabla \text{ReLU} W_\ell \right) [\mathbf{I}, \mathbf{H}^\top] \right)^{T-t} \right. \\ \quad \left. \left. \nabla 2\sigma (\mathbf{I}_{n_\ell} \otimes G_{L-1,t}^\top) \right) \right)^\top & \ell = L, \end{cases} \quad (20)
\end{aligned}$$

where $\mathbf{K}_{n_\ell, n_{\ell-1}}$ denotes a commutation matrix, which is a $n_\ell * n_{\ell-1} \times n_\ell * n_{\ell-1}$ permutation matrix that swaps rows and columns in the vectorization process.

A.4 Derivative of Coordinate-Wise Math-L2O

Based on the results in Appendix A.3, in this section, we construct the gradient formulations for Math-L2O model. We present the results in Equation (7) and Equation (8).

As defined in Equations (3) and (4), Math-L2O [23] learns to choose hyperparameters of existing non-learning algorithms [23, 34]. Suppose $P_i \in \mathbb{R}^{N*d}, i \in [0, \dots, T]$ is the hyperparameter vector generated by NNs. Suppose $X_{-1} := X_0$, based on Equation 3, the solution update process from the initial step is defined by:

$$\begin{aligned}
X_1 &= X_0 - \frac{1}{\beta} P_1 \odot \nabla F(X_0), \\
X_2 &= X_1 - \frac{1}{\beta} P_2 \odot \nabla F(X_1), \\
&\dots, \\
X_T &= X_{T-1} - \frac{1}{\beta} P_T \odot \nabla F(X_{T-1}),
\end{aligned} \quad (21)$$

We re-use the definition in Section 2 that defines $\mathcal{D}(\cdot)$ as the operator that constructs a diagonal matrix from a vector, we calculate the following one-line and linear-like formulation of X_T with X_0 :

$$X_T = \prod_{t=T}^1 \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t) \mathbf{M}^\top \mathbf{M} \right) X_0 + \frac{1}{\beta} \sum_{t=1}^T \prod_{s=T}^{t+1} \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s) \mathbf{M}^\top \mathbf{M} \right) \mathcal{D}(P_t) \mathbf{M}^\top Y. \quad (22)$$

Given that P_t is generated by a non-linear neural network with X_{t-1} as input, the resulting system dynamics are inherently non-linear. Consequently, this system cannot be formulated as the aforementioned linear dynamic system. Moreover, we note that for non-smooth problems, the uncertain sub-gradient can be replaced by the gradient map to obtain analogous formulations [34].

Due to the above computational graph in Figure 1, the gradient of X_t comes from X_{t-1} and P_t , which yields the following framework of each layer's derivative (Equation (6)):

$$\frac{\partial F}{\partial W_\ell} = \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{j=T}^{t+1} \left(\prod_{j=T}^{t+1} \frac{\partial X_j}{\partial X_{j-1}} \right) \frac{\partial X_t}{\partial P_t} \frac{\partial P_t}{\partial W_\ell} \right). \quad (23)$$

We obtain the above equation by counting the number of formulations from F to W_ℓ . From the Figure 1, we conclude that each timestamp t leads to the gradient of $\frac{\partial X_T}{\partial X_{T-1}}$. Thus, there are $\prod_{j=T}^{t+1} \frac{\partial X_j}{\partial X_{j-1}}$ blocks of formulation in total.

We start with deriving the formulation of gradient w.r.t. the GD algorithm, which yields the gradient of $\frac{\partial X_T}{\partial P_T}$. Due to the GD formulation in Equation (21), we derive $\frac{\partial X_t}{\partial X_{t-1}}$ as:

$$\begin{aligned}
\frac{\partial X_t}{\partial X_{t-1}} &= \mathbf{I}_d - \frac{1}{\beta} \frac{\partial (P_t \odot \nabla F(X_{t-1}))}{\partial X_{t-1}} \\
&= \mathbf{I}_d - \frac{1}{\beta} \frac{\partial P_t \odot (\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y))}{\partial X_{t-1}}, \\
&= \mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_t) \mathbf{M}^\top \mathbf{M} - \frac{1}{\beta} \frac{\partial P_t \odot (\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y))}{\partial P_t} \frac{\partial P_t}{\partial X_{t-1}}, \\
&= \mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_t) \mathbf{M}^\top \mathbf{M} - \frac{1}{\beta} \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y)) \frac{\partial P_t}{\partial X_{t-1}}.
\end{aligned} \tag{24}$$

Next, we calculate $\frac{\partial P_t}{\partial X_{t-1}}$. Similarly, we derive $\frac{\partial \text{vec}(G_{L,t})}{\partial W_\ell}$ and each $\frac{\partial \text{vec}(G_{L,j})}{\partial X_{j-1}}$ of Math-L2O layer-by-layer. $\frac{\partial \text{vec}(G_{L,t})}{\partial \text{vec}(G_{L-1,t})}$ in Math-L2O is similar to Equation (19). We calculate:

$$\begin{cases} \frac{\partial P_t}{\partial W_\ell} = \mathcal{D}(P_t \odot (1 - P_t/2)) (\mathbf{I}_d \otimes W_L) \prod_{j=L-1}^{\ell+1} \mathbf{D}_{j,t} \mathbf{I}_d \otimes W_j \mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^\top & \ell \in [L-1], \\ \frac{\partial P_t}{\partial W_L} = \mathcal{D}(P_t \odot (1 - P_t/2)) G_{L-1,t}^\top & \ell = L. \end{cases} \tag{25}$$

Similarly, we calculate the following derivative of output of Math-L2O w.r.t. its input at step t :

$$\frac{\partial P_t}{\partial X_{t-1}} = \mathcal{D}(P_t \odot (1 - P_t/2)) W_L (\prod_{\ell=L-1}^2 \mathbf{D}_{\ell,t} W_\ell) [\mathbf{I}, \mathbf{H}^\top]^\top. \tag{26}$$

Substituting Equation (26) into Equation (24) yields $\frac{\partial X_t}{\partial X_{t-1}}$:

$$\begin{aligned}
\frac{\partial X_t}{\partial X_{t-1}} &= \mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_t) \mathbf{M}^\top \mathbf{M} \\
&\quad - \frac{1}{\beta} \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y)) \mathcal{D}(P_t \odot (1 - P_t/2)) W_L (\prod_{\ell=L-1}^2 \mathbf{D}_{\ell,t} W_\ell) [\mathbf{I}, \mathbf{H}^\top]^\top.
\end{aligned} \tag{27}$$

We note that in [23], the gradient formulations are simplified in the implementation by detaching the input feature from the computational graph. Thus, we can eliminate the complicated last term in the above formulation, which leads to the following compact version:

$$\frac{\partial X_t}{\partial X_{t-1}} = \mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_t) \mathbf{M}^\top \mathbf{M}. \tag{28}$$

In this paper, we take the gradient formulation in Equation (28).

Next, we calculate the $\frac{\partial X_t}{\partial P_t}$ component in Equation (23). We calculate the derivative of GD's output w.r.t. its input hyperparameter P (generated by NNs) as:

$$\frac{\partial X_t}{\partial P_t} = -\frac{1}{\beta} \mathcal{D}(\nabla F(X_{t-1})) = -\frac{1}{\beta} \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y)), \tag{29}$$

where $\nabla F(X_{t-1}) := \mathbf{M}^\top (\mathbf{M} X_{t-1} - Y)$ is the first-order derivative of the objective in Equation 1.

Substituting Equation (25), Equation (28), and Equation (29) into Equation (23) yields the final derivative of all layers' parameters.

First, for $\ell = L$, since there is no cumulative gradients of later layers, Equation 8 is directly calculated by:

$$\begin{aligned}
\frac{\partial F}{\partial W_L} &= -\frac{1}{\beta} \sum_{t=1}^T (\mathbf{M}^\top (\mathbf{M} X_T - Y))^\top (\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j) \mathbf{M}^\top \mathbf{M}) \\
&\quad \mathcal{D}((\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y))) \mathcal{D}(P_t \odot (1 - P_t/2)) G_{L-1,t}^\top.
\end{aligned}$$

And its transpose is given by:

$$\begin{aligned}
\frac{\partial F}{\partial W_L}^\top &= -\frac{1}{\beta} \sum_{t=1}^T G_{L-1,t} \mathcal{D}(P_t \odot (1 - P_t/2)) \mathcal{D}((\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y))) \\
&\quad (\prod_{j=t+1}^T \mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M} \mathcal{D}(P_j)) \mathbf{M}^\top (\mathbf{M} X_T - Y).
\end{aligned} \tag{30}$$

When $\ell \in [L - 1]$, the derivative is calculated by:

$$\begin{aligned}\frac{\partial F}{\partial W_\ell} &= \frac{\partial F(X_T)}{\partial X_T} \left(\sum_{t=1}^T \left(\prod_{j=T}^{t+1} \frac{\partial X_j}{\partial X_{j-1}} \right) \frac{\partial X_t}{\partial P_t} \frac{\partial P_t}{\partial W_\ell} \right), \\ &= -\frac{1}{\beta} \sum_{t=1}^T (\mathbf{M}^\top (\mathbf{M} X_T - Y))^\top \left(\prod_{j=T}^{t+1} \mathbf{I}_d - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M} \mathcal{D}(P_j) \right) \\ &\quad \mathcal{D}((\mathbf{M}^\top (\mathbf{M} X_{t-1} - Y))) \mathcal{D}(P_t \odot (1 - P_t/2)) \\ &\quad (\mathbf{I}_d \otimes W_L) \prod_{j=L-1}^{\ell+1} \mathbf{D}_{j,t} \mathbf{I}_d \otimes W_j \mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^\top.\end{aligned}$$

Remark 2. The only difference between Equation (8) and Equation (7) lies in the last term, where Equation (7) is more complicated due to the accumulated gradients from later layers.

The above two formulations are used in the next section to derive the gradient bound for each layer.

A.5 Tools

In this section, prior to constructing the convergence bounds, we first derive several analytical tools. These tools are foundational for the convergence rate analysis and also establish key properties of the L2O models. We use superscript k to denote parameters and variables at training iteration k , and subscript t to denote the optimization step.

A.5.1 NN's Outputs are Bounded

First, we demonstrate that the outputs and inner outputs of NN layers within the L2O model are bounded.

Bound $\|\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}\|_2, \forall k, t.$

Lemma A.1. Suppose $\|\mathbf{M}^\top \mathbf{M}\|_2 \leq \beta$ and $0 < P_t^k < 2$, we have the following bound:

$$\|\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}\|_2 < 1. \quad (31)$$

Proof. Suppose eigenvalues and eigenvectors of $\mathbf{M}^\top \mathbf{M}$ are σ_i and $v_i, i \in [1, \dots, N * d]$ respectively, we calculate:

$$\frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M} v_i = \frac{\sigma_i}{\beta} \mathcal{D}(P_t^k) v_i.$$

Due to $0 < P_t^k < 2$, we have following spectral norm definition:

$$\|\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}\|_2 = \max_{x \in \mathbb{R}^d} \frac{x^\top (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}) x}{x^\top x}$$

Then, by taking $x = v_i$, we calculate:

$$v_i^\top (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}) v_i = 1 - \frac{1}{\beta} v_i^\top \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M} v_i = 1 - \frac{\sigma_i}{\beta} v_i^\top \mathcal{D}(P_t^k) v_i \stackrel{\textcircled{1}}{\leq} 1,$$

where $\textcircled{1}$ is due to $0 < P_t^k < 2$. □

Remark 3. In our design, we ensure $0 < P_t^k < 2$ by an activation function 2σ at the output layer.

Bound $\|\mathcal{D}(P_t^k)\|_2, \forall k, t.$ Similar to the bound of $\|\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k) \mathbf{M}^\top \mathbf{M}\|_2, \forall k, t$, due to the Sigmoid function, we directly have:

Lemma A.2. Suppose $0 < P_t^k < 2$, we have the following bound:

$$\|\mathcal{D}(P_t^k)\|_2 < 2. \quad (32)$$

Proof. Since \mathcal{D} is the diagonalization operation and $0 < P_t^k < 2$, we directly have $\|\mathcal{D}(P_t^k)\|_2 < 2$. □

Besides, we can derive another bound from the Lipschitz property for the Sigmoid activation function:

$$\begin{aligned}
\|\mathcal{D}(P_t^k)\|_2 &= \|2\sigma(\text{ReLU}(\text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^{k\top}) \cdots W_{L-1}^{k\top})W_L^{k\top})\|_\infty, \\
&\stackrel{\textcircled{1}}{\leq} \frac{1}{2} \| [X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)] \|_2 \prod_{s=1}^{L-1} \|W_s^k\|_2 + 1, \\
&\stackrel{\textcircled{2}}{\leq} \frac{1}{2} (\|X_t^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_t^k - Y)\|_2) \prod_{s=1}^{L-1} \|W_s^k\|_2 + 1.
\end{aligned} \tag{33}$$

① is from equation (17), Lemma 4.2 of [30]. ② is from triangle inequality.

Remark 4. In contrast to the Lipschitz continuous property of ReLU, the aforementioned bound associated with the Sigmoid function prevents the derivation of meaningful numerical results. To analyze the convergence rate of Gradient Descent (GD), a tighter bound on the neural network's output is required. One potential alternative is the convex cone defined by W_L^k for the last hidden layer. However, such a cone spans an unbounded space for the set of learnable parameters.

Bound Semi-Smoothness of NN's Output, i.e., $\|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2, \forall k, t$. Since our L2O model is a coordinate-wise model [23], suppose $P_i = \alpha_i(P_t^{k+1})_i + (1 - \alpha_i)(P_t^k)_i$, $\alpha_i \in [0, 1]$, based on Mean Value Theorem, we have $(\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k))_i = \frac{\partial F}{\partial P_i}((P_t^{k+1})_i - (P_t^k)_i)$. Thus, we bound $\|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2$ by the following lemma:

Lemma A.3. Denote $j \in [L]$, for some $\bar{\lambda}_j \in \mathbb{R}$, we assume $\|W_j^{k+1}\|_2 \leq \bar{\lambda}_j$. Using quantities from Equation (12), we have:

$$\begin{aligned}
&\|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2 \\
&\leq \frac{1}{2}(1 + \beta)\|X_{t-1}^{k+1} - X_{t-1}^k\|_2 \Theta_L \\
&\quad + \frac{1}{2}(\|X_{t-1}^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)\|_2) \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2.
\end{aligned} \tag{34}$$

Remark 5. The above lemma shows the output of NN is a ‘‘mixed’’ Lipschitz continuous on input feature and learnable parameters. The first term illustrates the Lipschitz property on input feature. The second term can be regarded as a Lipschitz property on learnable parameters with a stable input feature.

Proof. Due to Mean Value Theorem, we have:

$$\begin{aligned}
& \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2 \\
&= \|\mathcal{D}(2\sigma(\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-1}^{k+1\top})W_L^{k+1})) \\
&\quad - \mathcal{D}(2\sigma(\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-1}^k\top)W_L^k))\|_2, \\
&\leq (2\sigma(P_i)(1 - \sigma(P_i)))_{\max} \\
&\quad \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-1}^{k+1\top})W_L^{k+1} \\
&\quad - \text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-1}^k\top)W_L^k\|_\infty, \\
&\leq \frac{1}{2} \|\text{ReLU}(\text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-1}^{k+1\top})W_L^{k+1} \\
&\quad - \text{ReLU}(\text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-1}^k\top)W_L^k\|_\infty, \\
&\stackrel{\textcircled{1}}{\leq} \frac{1}{2} \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-1}^{k+1\top}) \\
&\quad - \text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-1}^k\top)\|_\infty \|W_L^{k+1} - W_L^k\|_2 \\
&\quad + \frac{1}{2} \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_{L-1}^k\top)\|_2 \|W_L^{k+1} - W_L^k\|_2, \\
&\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-2}^{k+1\top})W_{L-1}^{k+1\top} \\
&\quad - \text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-2}^k\top)W_{L-1}^k\|_\infty \bar{\lambda}_L \\
&\quad + \frac{1}{2} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \prod_{j=1}^{L-1} \bar{\lambda}_j \|W_L^{k+1} - W_L^k\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \frac{1}{2} \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-2}^{k+1\top}) \\
&\quad - \text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-2}^k\top)\|_\infty \bar{\lambda}_{L-1} \bar{\lambda}_L \\
&\quad + \frac{1}{2} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \prod_{j=1}^{L-1} \bar{\lambda}_j \|W_L^{k+1} - W_L^k\|_2, \\
&\quad + \frac{1}{2} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \prod_{j=1}^{L-2} \bar{\lambda}_j \bar{\lambda}_L \|W_{L-1}^{k+1} - W_{L-1}^k\|_2, \\
&\stackrel{\textcircled{4}}{=} \frac{1}{2} \|\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)]W_1^{k+1\top}) \cdots W_{L-2}^{k+1\top}) \\
&\quad - \text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^k\top) \cdots W_{L-2}^k\top)\|_\infty \bar{\lambda}_{L-1} \bar{\lambda}_L \\
&\quad + \frac{1}{2} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \Theta_L (\bar{\lambda}_L^{-1} \|W_L^{k+1} - W_L^k\|_2 + \bar{\lambda}_{L-1}^{-1} \|W_{L-1}^{k+1} - W_{L-1}^k\|_2), \\
&\quad \dots, \\
&\stackrel{\textcircled{5}}{\leq} \frac{1}{2} \|[X_{t-1}^{k+1}, \mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)] - [X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \Theta_L \\
&\quad + \frac{1}{2} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \Theta_L \left(\sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right), \\
&\stackrel{\textcircled{6}}{\leq} \frac{1}{2} (1 + \beta) \|X_{t-1}^{k+1} - X_{t-1}^k\|_2 \Theta_L \\
&\quad + \frac{1}{2} (\|X_{t-1}^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)\|_2) \Theta_L \left(\sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right).
\end{aligned}$$

① is due to triangle and Cauchy Schwarz inequalities, where we make a upper bound relaxation from ∞ -norm to 2-norm. ② is due to 1-Lipschitz property of ReLU and $\max(\|W_L^{k+1}\|_2, \|W_L^k\|_2) \leq \bar{\lambda}_L$ in the definition. It is note-worthy that any activations with constant-Lipchitz properties can be applied. ③ is due to triangle and Cauchy Schwarz inequalities as well. We make a arrangement in ④ and eliminate inductions in \dots . In ⑤, we make another upper bound relaxation from ∞ -norm to 2-norm. ⑥ is due to triangle inequality, the definition of Frobenius norm, and $\|\mathbf{M}^\top \mathbf{M}\|_2 \leq L$ of objective's L-smooth property. \square

Semi-Smoothness of Inner Output of NN, i.e., Bound $\|G_{\ell,t}^a - G_{\ell,t}^b\|_2, \ell \in [L-1], \forall a, b, t$.

Lemma A.4. Denote $\ell \in [L - 1]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$, we assume $\max(\|W_\ell^a\|_2, \|W_\ell^b\|_2) \leq \bar{\lambda}_\ell$. Using quantities from Equation (12), we have:

$$\begin{aligned} \|G_{\ell,t}^a - G_{\ell,t}^b\|_2 &\leq (1 + \beta) \|X_{t-1}^a - X_{t-1}^b\|_2 \prod_{j=1}^\ell \bar{\lambda}_j \\ &\quad + (\|X_{t-1}^b\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^b - Y)\|_2) \prod_{j=1}^\ell \bar{\lambda}_j \sum_{s=1}^\ell \bar{\lambda}_s^{-1} \|W_s^a - W_s^b\|_2. \end{aligned}$$

Proof. Since the bounding target in Lemma A.4 is a degenerated version of that in Lemma A.3. Similar to the proof of Lemma A.3, we calculate:

$$\begin{aligned} &\|G_{\ell,t}^a - G_{\ell,t}^b\|_2 \\ &= \|\text{ReLU}(\text{ReLU}([X_{t-1}^a, \mathbf{M}^\top(\mathbf{M}X_{t-1}^a - Y)]W_1^{a^\top}) \cdots W_\ell^{a^\top}) \\ &\quad - \text{ReLU}(\text{ReLU}([X_{t-1}^b, \mathbf{M}^\top(\mathbf{M}X_{t-1}^b - Y)]W_1^{b^\top}) \cdots W_\ell^{b^\top})\|_2, \\ &\leq \|[X_{t-1}^a, \mathbf{M}^\top(\mathbf{M}X_{t-1}^a - Y)] - [X_{t-1}^b, \mathbf{M}^\top(\mathbf{M}X_{t-1}^b - Y)]\|_2 \prod_{j=1}^\ell \bar{\lambda}_j \\ &\quad + \|[X_{t-1}^b, \mathbf{M}^\top(\mathbf{M}X_{t-1}^b - Y)]\|_2 \prod_{j=1}^\ell \bar{\lambda}_j \sum_{s=1}^\ell \bar{\lambda}_s^{-1} \|W_s^a - W_s^b\|_2, \\ &\leq (1 + \beta) \|X_{t-1}^a - X_{t-1}^b\|_2 \prod_{j=1}^\ell \bar{\lambda}_j \\ &\quad + (\|X_{t-1}^b\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^b - Y)\|_2) \prod_{j=1}^\ell \bar{\lambda}_j \sum_{s=1}^\ell \bar{\lambda}_s^{-1} \|W_s^a - W_s^b\|_2. \end{aligned}$$

□

Bound NN's Inner Output $G_{l,t}^k, l = [L - 1], \forall k, t$.

Lemma A.5. Denote $\ell \in [L - 1]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$, we assume $\|W_\ell^k\|_2 \leq \bar{\lambda}_\ell$. Using quantities from Equation (12), we have:

$$\|G_{\ell,t}^k\|_2 \leq ((1 + \beta)\|X_0\|_2 + (2t - 1 + \frac{2t-2}{\beta}))\|\mathbf{M}^\top Y\|_2 \prod_{s=1}^\ell \bar{\lambda}_s.$$

Proof.

$$\begin{aligned} \|G_{\ell,t}^k\|_2 &= \|\text{ReLU}(\text{ReLU}([X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]W_1^{k^\top}) \cdots W_\ell^{k^\top})\|_2, \\ &\stackrel{\textcircled{1}}{\leq} \|[X_{t-1}^k, \mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)]\|_2 \prod_{s=1}^\ell \|W_s^k\|_2, \\ &\stackrel{\textcircled{2}}{\leq} (\|X_{t-1}^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)\|_2) \prod_{s=1}^\ell \|W_s^k\|_2, \\ &\stackrel{\textcircled{3}}{\leq} ((1 + \beta)\|X_0\|_2 + \left(\frac{(1+\beta)2(t-1)}{\beta} + 1\right)\|\mathbf{M}^\top Y\|_2) \prod_{s=1}^\ell \|W_s^k\|_2, \\ &\leq ((1 + \beta)\|X_0\|_2 + (2t - 1 + \frac{2t-2}{\beta}))\|\mathbf{M}^\top Y\|_2 \prod_{s=1}^\ell \bar{\lambda}_s. \end{aligned}$$

① is from equation (17), Lemma 4.2 of [30]. ② is from triangle inequality. ③ is due to definition of β -smoothness of objective and upper bound of $\|X_t\|_2$ in Lemma A.6. □

A.5.2 Outputs of L2O are Bounded

Next, we establish bounds for the Math-L2O's outputs. Leveraging the momentum-free setting, we formulate the dynamics from X_0 to X_t as a *semi-linear* system, where parameters are non-linearly generated by the NN block (see Figure 1a). Application of the Cauchy-Schwarz and triangle inequalities to this system yields the following explicit bound.

Lemma A.6 (Bound on Math-L2O Output). *For any training iteration k , the t -th output X_t^k of Math-L2O (as per Equation (3)) is bounded by: $\|X_t^k\|_2 \leq \|X_0\|_2 + \frac{2t}{\beta}\|\mathbf{M}^\top Y\|_2$.*

Proof. We calculate the upper bound based on the one-line formulation from X_0 in Equation (22).

$$\begin{aligned}
& \|X_t^k\|_2 \\
&= \left\| \prod_{s=t}^1 (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M}) X_0 + \frac{1}{\beta} \sum_{s=1}^t \prod_{j=t}^{s+1} (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M}) \mathcal{D}(P_s^k) \mathbf{M}^\top Y \right\|_2 \\
&\stackrel{\textcircled{1}}{\leq} \left\| \prod_{s=1}^t (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M}) X_0 \right\|_2 + \left\| \frac{1}{\beta} \sum_{s=1}^t \prod_{j=t}^{s+1} (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M}) \mathcal{D}(P_s^k) \mathbf{M}^\top Y \right\|_2 \\
&\stackrel{\textcircled{2}}{\leq} \prod_{s=1}^t \left\| \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M} \right\|_2 \|X_0\|_2 \\
&\quad + \frac{1}{\beta} \sum_{s=1}^t \prod_{j=t}^{s+1} \left\| \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^k) \mathbf{M}^\top \mathbf{M} \right\|_2 \|\mathcal{D}(P_s^k)\|_2 \|\mathbf{M}^\top Y\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \|X_0\|_2 + \frac{2}{\beta} \sum_{s=1}^t \|\mathbf{M}^\top Y\|_2 = \|X_0\|_2 + \frac{2t}{\beta} \|\mathbf{M}^\top Y\|_2,
\end{aligned}$$

where $\textcircled{1}$ is from the triangle inequality, $\textcircled{2}$ is due to Cauchy Schwarz inequalities, and $\textcircled{3}$ is due to Lemma A.1 and Lemma A.2. \square

This lemma demonstrates that Math-L2O outputs remain bounded independently of the training iteration k and the specific learnable parameters.

A.5.3 L2O is Semi-Smooth to Its Parameters

In this section, we treat the L2O model defined in Equation (21) and its corresponding neural network as functions of their learnable parameters. We then prove that these functions are semi-smooth with respect to these parameters. This property is foundational for establishing the convergence of the gradient descent algorithm, as its analysis inherently involves the relationship between parameters at adjacent iterations.

First, we give the following explicit formulation of P :

$$\begin{aligned}
P_t^k &= 2\sigma(W_L^k \text{ReLU}(W_{L-1}^k (\cdots \text{ReLU}(W_1^k [X_{t-1}^k, \mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)]^\top) \cdots)))^\top, \\
&= 2\sigma(\text{ReLU}(\cdots \text{ReLU}([X_{t-1}^k, \mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)]W_1^\top) \cdots W_{L-1}^{k\top})W_L^k).
\end{aligned}$$

Moreover, we present ReLU activation function with signal matrices defined in Section 2. We denote \cdot_K as the entry-wise product to the matrices, which is also equivalent to reshape a matrix to a vector then product a diagonal signal matrix and reshape back afterward.

$$\begin{aligned}
P_t^k &= 2\sigma(W_L^k \mathbf{D}_{L-1} \cdot_K W_{L-1}^k (\cdots \mathbf{D}_1 \cdot_K (W_1^k [X_{t-1}^k, \mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)]^\top) \cdots))^\top, \\
&= 2\sigma((\cdots ([X_{t-1}^k, \mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)]W_1^\top) \cdot_K \mathbf{D}_1 \cdots)W_{L-1}^{k\top} \cdot_K \mathbf{D}_{L-1} W_L^k).
\end{aligned}$$

Proof for Lemma 4.2. We demonstrate the semi-smoothness of Math-L2O's output, i.e., bound $\|X_t^{k+1} - X_t^k\|_2, \forall k, t$

Proof. Diverging from the approach in [30], X_T^{k+1} and X_T^k are the outputs of a non-linear neural network corresponding to different inputs. A direct subtraction between these terms, as would be feasible in a linear-like system, is therefore intractable. Consequently, we must construct an upper bound for this difference. By applying a norm-based relaxation and utilizing the quantities defined in

Equation (12), we proceed with the following calculation:

$$\begin{aligned}
& \|X_t^{k+1} - X_t^k\|_2 \\
&= \|X_{t-1}^{k+1} - \frac{1}{\beta} \mathcal{D}(P_t^{k+1})(\mathbf{M}^\top(\mathbf{M}X_{t-1}^{k+1} - Y)) - (X_{t-1}^k - \frac{1}{\beta} \mathcal{D}(P_t^k)(\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)))\|_2, \\
&= \left\| \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^{k+1})\mathbf{M}^\top\mathbf{M} \right) X_{t-1}^{k+1} - \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k)\mathbf{M}^\top\mathbf{M} \right) X_{t-1}^k \right. \\
&\quad \left. + \frac{1}{\beta} (\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k))\mathbf{M}^\top Y \right\|_2 \\
&\stackrel{\textcircled{1}}{\leq} \left\| \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^{k+1})\mathbf{M}^\top\mathbf{M} \right) - \left(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k)\mathbf{M}^\top\mathbf{M} \right) \right\|_2 \|X_{t-1}^{k+1}\|_2 \\
&\quad + \left\| \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_t^k)\mathbf{M}^\top\mathbf{M} \right\|_2 \|X_{t-1}^{k+1} - X_{t-1}^k\|_2 + \frac{1}{\beta} \|\mathbf{M}^\top Y\|_2 \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2, \\
&\stackrel{\textcircled{2}}{\leq} \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2 \|X_{t-1}^{k+1}\|_2 + \|X_{t-1}^{k+1} - X_{t-1}^k\|_2 + \frac{1}{\beta} \|\mathbf{M}^\top Y\|_2 \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2 (\|X_0\|_2 + \frac{2t-2}{\beta} \|\mathbf{M}^\top Y\|_2) + \|X_{t-1}^{k+1} - X_{t-1}^k\|_2 \\
&\quad + \frac{1}{\beta} \|\mathbf{M}^\top Y\|_2 \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2, \\
&= (\|X_0\|_2 + \frac{2t-1}{\beta} \|\mathbf{M}^\top Y\|_2) \|\mathcal{D}(P_t^{k+1}) - \mathcal{D}(P_t^k)\|_2 + \|X_{t-1}^{k+1} - X_{t-1}^k\|_2, \\
&\stackrel{\textcircled{4}}{\leq} (\|X_0\|_2 + \frac{2t-1}{\beta} \|\mathbf{M}^\top Y\|_2) \\
&\quad \left(\frac{1}{2}(1+\beta) \|X_{t-1}^{k+1} - X_{t-1}^k\|_2 \Theta_L \right. \\
&\quad \left. + \frac{1}{2} (\|X_{t-1}^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)\|_2) \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \\
&\quad + \|X_{t-1}^{k+1} - X_{t-1}^k\|_2, \\
&= \left(1 + (\|X_0\|_2 + \frac{2t-1}{\beta} \|\mathbf{M}^\top Y\|_2) \frac{1+\beta}{2} \Theta_L \right) \|X_{t-1}^{k+1} - X_{t-1}^k\|_2, \\
&\quad + \frac{1}{2} (\|X_0\|_2 + \frac{2t-1}{\beta} \|\mathbf{M}^\top Y\|_2) \\
&\quad (\|X_{t-1}^k\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{t-1}^k - Y)\|_2) \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2, \\
&\stackrel{\textcircled{5}}{\leq} \frac{1}{2} \sum_{s=1}^t \left(\prod_{j=s+1}^t (1 + (\|X_0\|_2 + \frac{2j-1}{\beta} \|\mathbf{M}^\top Y\|_2) \frac{1+\beta}{2} \Theta_L) \right) \\
&\quad \underbrace{(\|X_0\|_2 + \frac{2s-1}{\beta} \|\mathbf{M}^\top Y\|_2) ((1+\beta)\|X_0\|_2 + (2s-1 + \frac{2s-2}{\beta}) \|\mathbf{M}^\top Y\|_2)}_{\Lambda_s} \\
&\quad \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2,
\end{aligned}$$

where $\textcircled{1}$ is from triangle inequality. $\textcircled{2}$ is from Lemma A.6. $\textcircled{3}$ is due to inductive summation to $t = 1$. $\textcircled{4}$ is due to the semi-smoothness of NN's output in Lemma A.3. $\textcircled{5}$ is from induction.

Remark 6. We note that the above upper bound relaxation is non-loose. Current existing approaches derive semi-smoothness in terms of NN functions, where parameters matrices are linearly applied and activation functions are Lipschitz continuous. However, in our setting under [23], the sigmoid activation is not Lipschitz continuous. Moreover, the input that is utilized to generate X_t^{k+1} is from X_{t-1}^{k+1} , which is not identical to the X_{t-1}^k for generating X_{t-1}^k .

□

A.5.4 Gradients are Bounded

In this section, we derive bound for the gradient of each layer's parameter at the given iteration k .

Proof for Lemma 4.1 We demonstrate that the gradients of Math-L2O's each layer are bounded.

Proof. For $\ell = L$, we calculate the gradient on W_L^k (Equation (8)):

$$\begin{aligned}
& \left\| \frac{\partial F}{\partial W_L^k} \right\|_2 \\
&= \frac{1}{\beta} \left\| \sum_{t=1}^T (\mathbf{M}^\top (\mathbf{M} X_T^k - Y))^\top \right. \\
&\quad \left. \left(\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M} \right) \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y)) \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) G_{L-1,t}^k \right\|_2^\top, \\
&\stackrel{\textcircled{1}}{\leq} \frac{1}{\beta} \sum_{t=1}^T \left\| \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \right\|_2 \prod_{j=T}^{t+1} \left\| \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M} \right) \right\|_2 \\
&\quad \left\| \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y)) \right\|_2 \left\| \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) \right\|_2 \left\| G_{L-1,t}^k \right\|_2, \\
&\stackrel{\textcircled{2}}{\leq} \frac{1}{2\sqrt{\beta}} \left\| \mathbf{M} X_T^k - Y \right\|_2 \sum_{t=1}^T (\left\| \mathbf{M}^\top \mathbf{M} X_{t-1}^k \right\|_2 + \left\| \mathbf{M}^\top Y \right\|_2) \left\| G_{L-1,t}^k \right\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \frac{\sqrt{\beta}}{2} \left\| \mathbf{M} X_T^k - Y \right\|_2 \prod_{\ell=1}^{L-1} \bar{\lambda}_\ell \sum_{t=1}^T ((1 + \beta) \|X_0\|_2 + (2t - 1 + \frac{2t-2}{\beta}) \left\| \mathbf{M}^\top Y \right\|_2) \\
&\quad (\|X_0\|_2 + \frac{2t-1}{\beta} \left\| \mathbf{M}^\top Y \right\|_2), \\
&= \frac{\sqrt{\beta}}{2} \left\| \mathbf{M} X_T^k - Y \right\|_2 \prod_{\ell=1}^{L-1} \bar{\lambda}_\ell \sum_{t=1}^T \underbrace{(1 + \beta) \|X_0\|_2^2 + ((4t - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \left\| \mathbf{M}^\top Y \right\|_2 + \frac{(2T-1)(\beta(2T-1)+(2T-2))}{\beta^2} \left\| \mathbf{M}^\top Y \right\|_2^2}_{\Lambda_t}, \\
&= \frac{\sqrt{\beta} \Theta_L S_{\Lambda, T}}{2 \lambda_L} \left\| \mathbf{M} X_T^k - Y \right\|_2,
\end{aligned}$$

where ① is from triangle and Cauchy-Schwarz inequalities. ② is from the bound of “ p ” in Lemma A.1. ③ is from the bound of L2O model’s output in Lemma A.6 and inner outputs in Lemma A.5.

For $\ell \in [L - 1]$, we calculate gradient on W_ℓ^k (Equation (7)) at iteration k by:

$$\begin{aligned}
& \left\| \frac{\partial F}{\partial W_\ell^k} \right\|_2 \\
&= \left\| -\frac{1}{\beta} \sum_{t=1}^T (\mathbf{M}^\top (\mathbf{M} X_T^k - Y))^\top \left(\prod_{j=T}^{t+1} \mathbf{I}_d - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M} \mathcal{D}(P_j^k) \right) \right. \\
&\quad \left. \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y)) \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) (\mathbf{I}_d \otimes W_L^k) \right. \\
&\quad \left. \prod_{j=L-1}^{\ell+1} \mathbf{D}_{j,t}^k \mathbf{I}_d \otimes W_j^k \mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^k \right\|_2^\top, \\
&\stackrel{\textcircled{1}}{\leq} \frac{1}{\beta} \sum_{t=1}^T \left\| \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \right\|_2 \prod_{j=T}^{t+1} \left\| \mathbf{I}_d - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M} \mathcal{D}(P_j^k) \right\|_2 \left\| \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y)) \right\|_2 \\
&\quad \left\| \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) (\mathbf{I}_d \otimes W_L^k) \right\|_2 \left\| \prod_{j=L-1}^{\ell+1} \mathbf{D}_{j,t}^k \mathbf{I}_d \otimes W_j^k \mathbf{I}_{n_\ell} \otimes G_{\ell-1,t}^k \right\|_2^\top, \\
&\stackrel{\textcircled{2}}{\leq} \frac{\sqrt{\beta}}{2} \left\| \mathbf{M} X_T^k - Y \right\|_2 \prod_{j=\ell+1}^L \|W_j^k\|_2 \sum_{t=1}^T (\left\| \mathbf{M}^\top \mathbf{M} X_{t-1}^k \right\|_2 + \left\| \mathbf{M}^\top Y \right\|_2) \left\| G_{\ell-1,t}^k \right\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \frac{\sqrt{\beta}}{2} \left\| \mathbf{M} X_T^k - Y \right\|_2 \prod_{j=1, j \neq \ell}^L \bar{\lambda}_j \sum_{t=1}^T \underbrace{(1 + \beta) \|X_0\|_2^2 + ((4t - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \left\| \mathbf{M}^\top Y \right\|_2 + \frac{(2T-1)(\beta(2T-1)+(2T-2))}{\beta^2} \left\| \mathbf{M}^\top Y \right\|_2^2}_{\Lambda_t}, \\
&= \frac{\sqrt{\beta} \Theta_L}{2 \lambda_\ell} S_{\Lambda, T} \left\| \mathbf{M} X_T^k - Y \right\|_2,
\end{aligned}$$

① is from triangle and Cauchy-Schwarz inequalities. Inequality ② is from bounds of “ p ” in Lemma A.1 and we make a rearrangement in it. In inequality ②, we use norm’s triangle inequality of dot product and Kronecker product, bounds of NN’s inner output in Lemma A.5, and we calculate $\prod_{j=1, j \neq \ell}^L \|W_j^k\|_2 = \prod_{j=\ell+1}^L \|W_j^k\|_2 * \prod_{s=1}^{\ell-1} \|W_s^k\|_2$. We reuse the result in the proof for the last layer’s gradient upper bound for case $\ell = L$ in equality ③ to get the final result. \square

A.6 Bound Linear Convergence Rate

Now we are able to substitute the above formulation into three bounding targets in Equation (44) and bound them one-by-one by the NTK theorem. We summarize the main idea of NTK theory before

the proof. The main technique of NTK theory is the establishment of non-singularity of the kernel matrix by a wide-NN layer, where kernel matrix is for the gradient of loss to learnable parameters. This invokes the Polyak-Lojasiewicz condition (a more relaxed condition than strongly convex) for linear convergence. Due to the page limit, we eliminate the explicit formulation of kernel matrix in main page. Following the methodology in [29], the non-singularity of kernel matrix is established by $\sigma_{\min}(G_{L-1,T}^0) > 0$. It is guaranteed by the conditions in Theorem 4.3 and implemented by the initialization strategy in Section 5.

Proof. We start to prove the Theorem 4.3 by proving the following lemma.

Lemma A.7.

$$\begin{cases} \|W_\ell^r\|_2 \leq \bar{\lambda}_\ell, & \ell \in [L], \quad r \in [0, k], \\ \sigma_{\min}(G_{L-1,T}^r) \geq \frac{1}{2}\alpha_0, & r \in [0, k], \\ F([W]^r) \leq (1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4)^r F([W]^0), & r \in [0, k]. \end{cases} \quad (35)$$

Remark 7. The first inequality means that there exists a scalar $\bar{\lambda}_\ell$ that bounds each layer's learnable parameter. The second inequality means that the last inner output is lower bounded. The last inequality is the linear rate of training.

A.6.1 Induction Part 1: NN's Parameter and the Last Inner Output are Bounded

For $k = 0$, Equation (35) degenerates and holds trivially. Assume Equation (35) holds up to iteration k , we aim to prove it still holds for iteration $k + 1$. First, we calculate the following term:

$$\begin{aligned} \|W_\ell^{k+1} - W_\ell^0\|_2 &\stackrel{\textcircled{1}}{\leq} \sum_{s=0}^k \|W_\ell^{s+1} - W_\ell^s\|_2 \\ &\stackrel{\textcircled{2}}{=} \eta \sum_{s=0}^k \left\| \frac{\partial F}{\partial W_\ell^s} \right\|_2 \\ &\stackrel{\textcircled{3}}{\leq} \eta \sum_{s=0}^k \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} \|\mathbf{M}X_T^s - Y\|_2, \\ &\stackrel{\textcircled{4}}{\leq} \eta \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} \sum_{s=0}^k (1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4)^{s/2} \|\mathbf{M}X_T^0 - Y\|_2, \end{aligned}$$

where $\textcircled{1}$ is due to triangle inequality. $\textcircled{2}$ is due to the definition of gradient descent. $\textcircled{3}$ is due the gradient is being upper-bounded in Lemma 4.1 and our assumption that $\|W_\ell^r\|_2 \leq \bar{\lambda}_\ell$, $\ell \in [L], \forall r \in [0, k]$. $\textcircled{4}$ is due to the linear rate in our induction assumption.

Define $u := \sqrt{1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4}$, we calculate the summation of geometric sequence by:

$$\begin{aligned} \eta \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} \sum_{s=0}^k u^s \|\mathbf{M}X_T^0 - Y\|_2 &= \eta \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} \frac{1-u^{k+1}}{1-u} \|\mathbf{M}X_T^0 - Y\|_2, \\ &\stackrel{\textcircled{1}}{=} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} (1-u^2) \frac{1-u^{k+1}}{1-u} \|\mathbf{M}X_T^0 - Y\|_2, \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} \|\mathbf{M}X_T^0 - Y\|_2, \\ &\stackrel{\textcircled{3}}{\leq} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} \frac{\sqrt{\beta} \Theta_L}{2\lambda_\ell} S_{\Lambda,T} (\sqrt{\beta} \|X_0\|_2 + (2T+1) \|Y\|_2), \\ &\stackrel{\textcircled{4}}{\leq} C_\ell, \end{aligned}$$

where $\textcircled{1}$ is due to $1 - u^2 = \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4$. $\textcircled{2}$ is due to $0 \leq u \leq 1$. $\textcircled{3}$ is due to NN's output's bound in Lemma A.6. $\textcircled{4}$ is due to the lower bound on the singular value of last inner output layer in Equation (13c).

Thus, we have:

$$\|W_\ell^{k+1} - W_\ell^0\|_2 \leq C_\ell. \quad (36)$$

Denote $\sigma_1(\cdot)$ as calculating the smallest singular value of any matrices, due to Weyl's inequality [28], we have:

$$\begin{aligned} |\|W_\ell^{k+1}\|_2 - \|W_\ell^0\|_2| &\leq \sigma_1(W_\ell^{k+1} - W_\ell^0), \\ &\leq \|W_\ell^{k+1} - W_\ell^0\|_2, \\ &\leq C_\ell. \end{aligned}$$

where the first inequality is from Weyl's inequality and the last inequality is due to Equation (36). Then, we directly have $\|W_\ell^{k+1}\|_2 - \|W_\ell^0\|_2 \leq C_\ell$ and $\|W_\ell^{k+1}\|_2 \leq \|W_\ell^0\|_2 + C_\ell = \bar{\lambda}_\ell$.

Next, we bound $G_{L-1,T}^{k+1}$ by calculating:

$$\begin{aligned} &\|G_{L-1,T}^{k+1} - G_{L-1,T}^0\|_2 \\ &\stackrel{\textcircled{1}}{\leq} (1 + \beta) \|X_{T-1}^{k+1} - X_{T-1}^0\|_2 \prod_{j=1}^{L-1} \bar{\lambda}_j \\ &\quad + (\|X_{T-1}^0\|_2 + \|\mathbf{M}^\top (\mathbf{M}X_{T-1}^0 - Y)\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^0\|_2, \\ &\stackrel{\textcircled{2}}{\leq} (1 + \beta) 2(\|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \\ &\quad + (\|X_{T-1}^0\|_2 + \|\mathbf{M}^\top (\mathbf{M}X_{T-1}^0 - Y)\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^0\|_2, \\ &\stackrel{\textcircled{3}}{\leq} (1 + \beta) \underbrace{\sum_{i=0}^k \frac{1}{2} \Theta_L \sum_{s=1}^{T-1} \left(\prod_{j=s+1}^{T-1} (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right)}_{\delta_1^{T-1}} \Lambda_s \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{i+1} - W_\ell^i\|_2 \prod_{j=1}^{L-1} \bar{\lambda}_j \\ &\quad + (\|X_{T-1}^0\|_2 + \|\mathbf{M}^\top (\mathbf{M}X_{T-1}^0 - Y)\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^0\|_2, \end{aligned} \tag{37}$$

where $\textcircled{1}$ is due to the semi-smoothness of NN's inner output in Lemma A.4. $\textcircled{2}$ is due to the triangle inequality. $\textcircled{3}$ is due to semi-smoothness of L2O in Lemma 4.2.

Further, based on the inner results in the former demonstration for $\|W_\ell^{k+1} - W_\ell^0\|_2$, we have:

$$\sum_{i=0}^k \|W_\ell^{i+1} - W_\ell^i\|_2 \leq \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} \frac{\sqrt{\beta} \Theta_L}{2\bar{\lambda}_\ell} S_{\Lambda,T} \|\mathbf{M}X_T^0 - Y\|_2.$$

Substituting above result back into Equation (37) yields:

$$\begin{aligned} &\|G_{L-1,T}^{k+1} - G_{L-1,T}^0\|_2 \\ &\leq (1 + \beta) 2(\|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \\ &\quad + (\|X_{T-1}^0\|_2 + \|\mathbf{M}^\top (\mathbf{M}X_{T-1}^0 - Y)\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} \frac{\sqrt{\beta} \Theta_L}{2\bar{\lambda}_\ell} S_{\Lambda,T} \|\mathbf{M}X_T^0 - Y\|_2, \\ &\stackrel{\textcircled{1}}{\leq} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} (1 + \beta) \zeta_2 (\sqrt{\beta} \|X_0\|_2 + (2T + 1) \|Y\|_2) S_{\Lambda,T} \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \frac{\sqrt{\beta} \Theta_L}{2\bar{\lambda}_\ell} \\ &\quad + 2(1 + \beta) (\|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j, \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} (1 + \beta) \zeta_2 (\sqrt{\beta} \|X_0\|_2 + (2T + 1) \|Y\|_2) S_{\Lambda,T} \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \frac{\sqrt{\beta} \Theta_L}{2\bar{\lambda}_\ell} \\ &\quad + \frac{1}{4} \alpha_0, \\ &\stackrel{\textcircled{3}}{\leq} \frac{1}{2} \alpha_0, \end{aligned} \tag{38}$$

where $\textcircled{1}$ is due to NN's output's bound in Lemma A.6 and $\textcircled{2}$ and $\textcircled{3}$ are due to the other lower bound for minimal singular value of NN's inner output in Equation (13a) and Equation (13d). The inequality in Equation (38) implies $\sigma_{\min}(G_{L-1}^{k+1}) \geq \frac{1}{2} \alpha_0$ since $\sigma_{\min}(G_{L-1}^0) = \alpha_0$.

Based on the above two inequalities, we prove the linear rate in Theorem 4.3 step-by-step in the following sub-section.

A.6.2 Induction Part 2: Linear Convergence

In this section, we aim to prove that $F([W]^{k+1}) \leq (1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4)^{k+1} F([W]^0)$.

Step 1: Split Perfect Square By leveraging term $\mathbf{M}X_T^k$, we can split the perfect square in objective $F([W]^{k+1})$ as:

$$F([W]^{k+1}) = F([W]^k) + \frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2 + (\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k)^\top (\mathbf{M}X_T^k - Y). \quad (39)$$

Based on [29], we aim to demonstrate that $F([W]^{k+1})$ can be upper-bounded by $c_k F([W]^k)$, where $c_k < 1$ is a coefficient related to training iteration k .

Step 2: Bound Term-by-Term We aim to upper bound all terms in Equation (39) by $F([W]^k)$.

Bound the first term $\frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2$. First, based on the β -smoothness of objective F , we calculate

$$\begin{aligned} \frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2 &= \frac{1}{2} (X_T^{k+1} - X_T^k)^\top \mathbf{M}^\top \mathbf{M} (X_T^{k+1} - X_T^k), \\ &\leq \frac{1}{2} \|X_T^{k+1} - X_T^k\|_2^2 \|\mathbf{M}^\top \mathbf{M}\|_2, \\ &\leq \frac{\beta}{2} \|X_T^{k+1} - X_T^k\|_2^2. \end{aligned}$$

The above inequality shows that we need to bound the distance between outputs of two iterations. Moreover, since our target is to construct linear convergence rate, we need to find the upper bound of above inequality w.r.t. the objective $F([W]^k)$, i.e., $\frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2$. We apply Lemma 4.2 to derive the following lemma.

Lemma A.8. Denote $\ell \in [L]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$, we assume $\max(\|W_\ell^{k+1}\|_2, \|W_\ell^k\|_2) \leq \bar{\lambda}_\ell, \forall k$. Using quantities from Equation (12), we further define the following quantities with $i, j \in [T]$:

$$\begin{aligned} \Lambda_i &= (1 + \beta) \|X_0\|_2^2 + ((4i - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 \\ &\quad + \frac{(2i-1)(\beta(2i-1)+(2i-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2, \\ \Phi_j &= \|X_0\|_2 + \frac{2j-1}{\beta} \|\mathbf{M}^\top Y\|_2, \\ \delta_1^T &= \left(\sum_{s=1}^T \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \left(\sum_{j=1}^s \Lambda_j \right) \right). \end{aligned}$$

We have the following upper bounding property:

$$\frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2 \leq \frac{\beta^2 \eta^2}{16} (\delta_1^T)^2 \left(S_{\Lambda, T} \right)^2 \left(\Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)^2 \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2. \quad (40)$$

Proof. We calculate:

$$\begin{aligned} \frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2 &\leq \frac{\beta}{2} \|X_T^{k+1} - X_T^k\|_2^2, \\ &\stackrel{\textcircled{1}}{\leq} \frac{\beta}{2} \left(\sum_{s=1}^T \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \frac{1}{2} \Lambda_s \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right)^2, \\ &\stackrel{\textcircled{2}}{=} \frac{\beta \eta^2}{2} \left(\sum_{s=1}^T \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \frac{1}{2} \Lambda_s \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \left\| \frac{\partial F}{\partial W_\ell^k} \right\|_2 \right)^2, \\ &\stackrel{\textcircled{3}}{\leq} \frac{\beta \eta^2}{2} \left(\sum_{s=1}^T \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \frac{1}{2} \Lambda_s \Theta_L \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \frac{\sqrt{\beta} \Theta_L}{2 \bar{\lambda}_\ell} \left(S_{\Lambda, T} \right) \|\mathbf{M}X_T^k - Y\|_2 \right)^2, \\ &= \frac{\beta^2 \eta^2}{32} \underbrace{\left(\left(\sum_{s=1}^T \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \Lambda_s \right) \left(S_{\Lambda, T} \right) \Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)}_{\delta_1^T} \|\mathbf{M}X_T^k - Y\|_2^2, \\ &= \frac{\beta^2 \eta^2}{16} (\delta_1^T)^2 \left(S_{\Lambda, T} \right)^2 \left(\Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)^2 \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2, \end{aligned} \quad (41)$$

① is from semi-smoothness of L2O's output in Lemma 4.2, Appendix A.5.3. ② is due to gradient descent with learning rate η . ③ is from gradient bounds in Lemma 4.1. \square

Bound the second term $(\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k)^\top (\mathbf{M}X_T^k - Y)$. We calculate:

$$\begin{aligned} & (\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k)^\top (\mathbf{M}X_T^k - Y) \\ &= (X_T^{k+1} - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \\ &= (X_T^{k+1} - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y). \end{aligned} \quad (42)$$

Following the methodology in [29], we hold all other learnable parameters fixed and focus the analysis on the gradient with respect to the last layer, W_L . This approach facilitates the construction of a non-singular NTK, which in turn establishes the PL condition, thereby guaranteeing a linear convergence rate.

Given last NN layer's learnable parameter W_L^{k+1} at iteration $k+1$, due to the GD formulation in Equation (21), we define the following quantity:

$$Z = X_{T-1}^k - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1} G_{L-1,T}^k)^\top) \mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y), \quad (43)$$

where $G_{L-1,T}^k$ represents inner output of layer $L-1$ at training iteration k .

With Z , we reformulate Equation (42) as:

$$\begin{aligned} & (X_T^{k+1} - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \\ &= (X_T^{k+1} - Z + Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \\ &= (X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) + (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \end{aligned} \quad (44)$$

where X_T^{k+1} at training iteration $k+1$ with W_L^{k+1} and solution X_T^k at training iteration k with W_L^k are defined as:

$$X_T^{k+1} = X_{T-1}^{k+1} - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1} G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top (\mathbf{M}X_{T-1}^{k+1} - Y).$$

$$X_T^k = X_{T-1}^k - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^k G_{L-1,T}^k)^\top) \mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y).$$

Then, we have the following lemmas to bound the two terms, respectively:

Lemma A.9. Denote $\ell \in [L]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$ with $j \in [T]$, we assume $\max(\|W_\ell^{k+1}\|_2, \|W_\ell^k\|_2) \leq \bar{\lambda}_\ell$. Define the following quantities with $t \in [T]$:

$$\begin{aligned} \Lambda_t &= (1 + \beta) \|X_0\|_2^2 + ((4t - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 \\ &\quad + \frac{(2T-1)(\beta(2T-1) + (2T-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2, \\ \Phi_j &= \|X_0\|_2 + \frac{2j-1}{\beta} \|\mathbf{M}^\top Y\|_2, \\ \Theta_L &= \Theta_L, \\ \delta_2 &= \sum_{s=1}^{T-1} \left(\prod_{j=s+1}^T (1 + \frac{1+\beta}{2} \Theta_L \Phi_j) \right) \Lambda_s. \end{aligned}$$

We have the following upper bounding property:

$$(X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \leq \frac{\beta\eta}{2} (\Lambda_T + \delta_2) \Theta_L^2 S_{\bar{\lambda},L} S_{\Lambda,T}^{\frac{1}{2}} \|\mathbf{M}X_T^k - Y\|_2^2.$$

Proof. We straightforwardly apply upper bound relaxation in this part, where we reuse the results of the first term $\frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2$'s upper bound in Lemma A.8.

To reuse the results, we would like to construct the $X_{T-1}^{k+1} - X_{T-1}^k$ term. We substitute Equation (46) into above equation and use the Cauchy-Schwarz inequality for vectors to split our bounding targets

into two parts and relax the L_2 -norm of vector summations into each element by triangle inequalities:

$$\begin{aligned}
& (X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
&= \left(X_{T-1}^{k+1} - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top (\mathbf{M}X_{T-1}^{k+1} - Y) \right. \\
&\quad \left. - \left(X_{T-1}^k - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top) \mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y) \right) \right)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \\
&\stackrel{\textcircled{1}}{\leq} \left(\left\| \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top \mathbf{M} \right) X_{T-1}^{k+1} \right. \right. \\
&\quad \left. \left. - \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top) \mathbf{M}^\top \mathbf{M} \right) X_{T-1}^k \right\|_2 \right. \\
&\quad \left. + \frac{1}{\beta} \left\| \underbrace{\left(\mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) - \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top) \right)}_{C_{k+1}} \mathbf{M}^\top Y \right\|_2 \right) \\
&\quad \left\| \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \right\|_2, \\
&\stackrel{\textcircled{2}}{\leq} \left(\left\| \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top \mathbf{M} \right) (X_{T-1}^{k+1} - X_{T-1}^k) \right\|_2 \right. \\
&\quad \left. + \left\| \left(\left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top \mathbf{M} \right) \right. \right. \right. \\
&\quad \left. \left. - \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top) \mathbf{M}^\top \mathbf{M} \right) \right) X_{T-1}^k \right\|_2 \\
&\quad \left. + \frac{1}{\beta} \|C_{k+1} \mathbf{M}^\top Y\|_2 \right) \left\| \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \right\|_2, \\
&\stackrel{\textcircled{3}}{\leq} \left(\left\| \left(\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top \mathbf{M} \right) \right\|_2 \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \right. \\
&\quad \left. + \frac{1}{\beta} \|C_{k+1} \mathbf{M}^\top \mathbf{M}\|_2 \|X_{T-1}^k\|_2 + \frac{1}{\beta} \|C_{k+1}\|_2 \|\mathbf{M}^\top Y\|_2 \right) \left\| \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \right\|_2, \\
&\stackrel{\textcircled{4}}{\leq} \left(\|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + \|X_{T-1}^k\|_2 \|C_{k+1}\|_2 + \frac{1}{\beta} \|\mathbf{M}^\top Y\|_2 \|C_{k+1}\|_2 \right) \left\| \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \right\|_2, \\
&\stackrel{\textcircled{5}}{\leq} \left(\|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + (\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2) \|C_{k+1}\|_2 \right) \left\| \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \right\|_2,
\end{aligned} \tag{45}$$

where ① is due to triangle and Cauchy-Schwarz inequalities. ② is due to triangle inequality. ③ is due to Cauchy-Schwarz inequality. ④ is due to β -smooth definition that $\mathbf{M}^\top \mathbf{M} \leq \beta$ and $\|\mathbf{I}_d - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) \mathbf{M}^\top \mathbf{M}\|_2 \leq 1$ in Lemma A.1. ⑤ is due to the upper bound of X_{T-1} in Lemma A.6.

Further, we bound $C_{k+1} := \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) - \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top)$. We apply the Mean Value Theorem and assume a point v_1^k . For v_1^k 's each entry $(v_1^k)_i$, for some $\alpha_{1i}^k \in [0, 1]$, we calculate $(v_1^k)_i$ as:

$$(v_1^k)_i = \alpha_{1i}^k ((W_L^{k+1}G_{L-1,T}^{k+1})^\top)_i + (1 - \alpha_{1i}^k) ((W_L^{k+1}G_{L-1,T}^k)^\top)_i.$$

Then, we can represent quantity $\|C_{k+1}\|_2$ by:

$$\begin{aligned}
& \|\mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^{k+1})^\top) - \mathcal{D}(2\sigma(W_L^{k+1}G_{L-1,T}^k)^\top)\|_2 \\
&\stackrel{\textcircled{1}}{\leq} \left\| \frac{\partial 2\sigma}{\partial v_1^k} \odot (W_L^{k+1}G_{L-1,T}^{k+1} - W_L^{k+1}G_{L-1,T}^k)^\top \right\|_\infty, \\
&\stackrel{\textcircled{2}}{\leq} \frac{1}{2} \left\| (W_L^{k+1}G_{L-1,T}^{k+1} - W_L^{k+1}G_{L-1,T}^k)^\top \right\|_\infty, \\
&\stackrel{\textcircled{3}}{\leq} \frac{1}{2} \|W_L^{k+1}\|_2 \|G_{L-1,T}^{k+1} - G_{L-1,T}^k\|_2 \leq \frac{1}{2} \bar{\lambda}_L \|G_{L-1,T}^{k+1} - G_{L-1,T}^k\|_2,
\end{aligned}$$

where ① is from the Mean Value Theorem. ② is from the gradient upper bound of Sigmoid function. ③ is from triangle inequality and definition of learnable parameter W_L .

We further substitute the upper bound of $\|G_{L-1,T}^{k+1} - G_{L-1,T}^k\|_2$ in Lemma A.4 and calculate:

$$\begin{aligned}
& \frac{1}{2} \bar{\lambda}_L \|G_{L-1,T}^{k+1} - G_{L-1,T}^k\|_2 \\
& \leq \frac{1}{2} \bar{\lambda}_L \left((1 + \beta) \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \prod_{j=1}^{L-1} \bar{\lambda}_j \right. \\
& \quad \left. + (\|X_{T-1}^k\|_2 + \|\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y)\|_2) \prod_{j=1}^{L-1} \bar{\lambda}_j \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \\
& \stackrel{\textcircled{1}}{\leq} \frac{1}{2} (1 + \beta) \Theta_L \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \\
& \quad + \frac{1}{2} \left((1 + \beta) \|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2 \right) \Theta_L \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2.
\end{aligned}$$

where $\textcircled{1}$ is due to upper bound of X_{T-1} in Lemma A.6.

Substituting the above inequality back into Equation (45) yields:

$$\begin{aligned}
& (X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
& \leq \left(\|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + (\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2) \|C_{k+1}\|_2 \right) \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \leq \left(\|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \right. \\
& \quad \left. + (\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2) \right. \\
& \quad \left(\frac{1}{2} (1 + \beta) \Theta_L \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \right. \\
& \quad \left. + \frac{1}{2} ((1 + \beta) \|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2) \Theta_L \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \\
& \quad \left. \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \right. \\
& = \left(\left(1 + \frac{1+\beta}{2} \Theta_L (\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2) \right) \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 \right. \\
& \quad \left. + \left(\frac{1}{2} ((1 + \beta) \|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2) \right. \right. \\
& \quad \left. \left. (\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2) \Theta_L \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \right) \\
& \quad \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& = \left(\left(1 + \frac{1+\beta}{2} \Theta_L \underbrace{(\|X_0\|_2 + \frac{2T-1}{\beta} \|\mathbf{M}^\top Y\|_2)}_{\Phi_T} \right) \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + \right. \\
& \quad \left. \underbrace{\left(\frac{1}{2} (1 + \beta) \|X_0\|_2^2 + ((4T - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 + \frac{(2T-1)(\beta(2T-1) + (2T-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right)}_{\Lambda_T} \right. \\
& \quad \left. \Theta_L \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& = \left(\left(1 + \frac{1+\beta}{2} \Theta_L \Phi_T \right) \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + \frac{1}{2} \Lambda_T \Theta_L \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \\
& \quad \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2,
\end{aligned}$$

Further, we apply semi-smoothness of L2O model in Lemma 4.2 and upper bound of gradient in Lemma 4.1 to derive the upper bound. We calculate:

$$\begin{aligned}
& (X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
& \leq \left(\left(1 + \frac{1+\beta}{2}\Theta_L\Phi_T\right) \|X_{T-1}^{k+1} - X_{T-1}^k\|_2 + \frac{1}{2}\Lambda_T\Theta_L\sum_{\ell=1}^{L-1}\bar{\lambda}_\ell^{-1}\|W_\ell^{k+1} - W_\ell^k\|_2 \right) \\
& \quad \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \stackrel{\textcircled{1}}{\leq} \left(\left(1 + \frac{1+\beta}{2}\Theta_L\Phi_T\right) \frac{1}{2}\Theta_L\sum_{s=1}^{T-1} \left(\prod_{j=s+1}^{T-1} \left(1 + \frac{1+\beta}{2}\Theta_L\Phi_j\right) \right) \Lambda_s \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right. \\
& \quad \left. + \frac{1}{2}\Lambda_T\Theta_L\sum_{\ell=1}^{L-1}\bar{\lambda}_\ell^{-1}\|W_\ell^{k+1} - W_\ell^k\|_2 \right) \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \leq \left(\frac{1}{2}\Theta_L \underbrace{\sum_{s=1}^{T-1} \left(\prod_{j=s+1}^T \left(1 + \frac{1+\beta}{2}\Theta_L\Phi_j\right) \right) \Lambda_s \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2}_{\delta_2} \right. \\
& \quad \left. + \frac{1}{2}\Lambda_T\Theta_L\sum_{\ell=1}^{L-1}\bar{\lambda}_\ell^{-1}\|W_\ell^{k+1} - W_\ell^k\|_2 \right) \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& = \frac{1}{2}\Theta_L \left(\delta_2 \bar{\lambda}_L^{-1} \|W_L^{k+1} - W_L^k\|_2 + (\Lambda_T + \delta_2) \sum_{\ell=1}^{L-1} \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \right) \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \stackrel{\textcircled{2}}{\leq} \frac{1}{2}\Theta_L (\Lambda_T + \delta_2) \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2,
\end{aligned}$$

where $\textcircled{1}$ is due to Lemma 4.2. $\textcircled{2}$ is due to $\Lambda_T \geq 0$.

Further, based on the gradient descent, i.e., $W_\ell^{k+1} = W_\ell^k - \eta \frac{\partial F}{\partial W_\ell^k}$, we substitute the bound of gradient in Lemma 4.1 and calculate:

$$\begin{aligned}
& (X_T^{k+1} - Z)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
& \leq \frac{1}{2}\Theta_L (\Lambda_T + \delta_2) \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \|W_\ell^{k+1} - W_\ell^k\|_2 \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \leq \frac{\eta}{2}\Theta_L (\Lambda_T + \delta_2) \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \left\| \frac{\partial F}{\partial W_\ell^k} \right\|_2 \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \stackrel{\textcircled{1}}{\leq} \frac{\eta}{2}\Theta_L (\Lambda_T + \delta_2) \sum_{\ell=1}^L \bar{\lambda}_\ell^{-1} \frac{\sqrt{\beta}\Theta_L}{2\lambda_\ell} S_{\Lambda,T} \|\mathbf{M}X_T^k - Y\|_2 \|\mathbf{M}^\top (\mathbf{M}X_T^k - Y)\|_2, \\
& \stackrel{\textcircled{2}}{\leq} \frac{\beta\eta}{2} (\Lambda_T + \delta_2) \Theta_L^2 S_{\bar{\lambda},L} S_{\Lambda,T} \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2,
\end{aligned}$$

where $\textcircled{1}$ is due to Lemma 4.1 and $\textcircled{2}$ is due to $\|M\|_2 \leq \sqrt{\beta}$. \square

Lemma A.10. Define the following quantities with $t \in [T]$:

$$\begin{aligned}
\Lambda_t &= (1 + \beta) \|X_0\|_2^2 + ((4t - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 \\
& \quad + \frac{(2T-1)(\beta(2T-1)+(2T-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2, \\
\Phi_j &= \|X_0\|_2 + \frac{2j-1}{\beta} \|\mathbf{M}^\top Y\|_2, \\
\Theta_L &= \Theta_L, \\
\delta_3 &= ((1 + \beta) \|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2).
\end{aligned}$$

We have the following upperly bounding property:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
& \leq \left(-\eta\sigma(\delta_3\Theta_L)^2(1 - \sigma(\delta_3\Theta_L))^2 \frac{\beta^2}{\beta^2} \alpha_0^2 + \frac{\eta\beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t \right) \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2.
\end{aligned}$$

Proof. In our above demonstrations, we have constructed a non-negative coefficient of the upper bound w.r.t. the objective $\frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2$. To achieve the requirement of the linear convergence

rate, we would like a negative one from our remaining bounding target. We calculate:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \\
&= \left(X_{T-1}^k - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^{k+1} G_{L-1,T}^k)^\top) (\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \right. \\
&\quad \left. - \left(X_{T-1}^k - \frac{1}{\beta} \mathcal{D}(2\sigma(W_L^k G_{L-1,T}^k)^\top) (\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \right) \right)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y), \quad (46) \\
&= -\frac{1}{\beta} (\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y))^\top \mathcal{D}(2\sigma(W_L^{k+1} G_{L-1,T}^k)^\top - 2\sigma(W_L^k G_{L-1,T}^k)^\top) \\
&\quad (\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)).
\end{aligned}$$

Similarly, due to Mean Value Theorem, suppose $v_{2,i}^k = \alpha_i(W_L^{k+1} G_{L-1,T}^k)_i + (1 - \alpha_i)(W_L^k G_{L-1,T}^k)_i$, $v_{2,i}^k \in [0, 1]$, based on Mean Value Theorem, we calculate:

$$2\sigma(W_L^{k+1} G_{L-1,T}^k)_i^\top - 2\sigma(W_L^k G_{L-1,T}^k)_i^\top = \frac{\partial(2\sigma(v_{2,i}^k))}{\partial(v_{2,i}^k)_i} (W_L^{k+1} G_{L-1,T}^k)_i - (W_L^k G_{L-1,T}^k)_i.$$

Denote $v_{2,i}^k := \lceil \frac{\partial(2\sigma(v_{2,i}^k))}{\partial(v_{2,i}^k)_i} \rceil$, we calculate:

$$\begin{aligned}
& \mathcal{D}(2\sigma(W_L^{k+1} G_{L-1,T}^k)^\top - 2\sigma(W_L^k G_{L-1,T}^k)^\top) \\
&= \mathcal{D}\left(\left[\frac{\partial 2\sigma(v_{2,i}^k)}{\partial v_{2,i}^k} ((W_L^{k+1} G_{L-1,T}^k)_i - (W_L^k G_{L-1,T}^k)_i)\right]^\top\right), \\
&= \mathcal{D}\left([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top ((W_L^{k+1} - W_L^k) G_{L-1,T}^k)_i\right)^\top, \\
&= \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(((W_L^{k+1} - W_L^k) G_{L-1,T}^k)^\top), \\
&\stackrel{\textcircled{1}}{=} -\eta \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}\left(\frac{\partial F}{\partial W_L^k} G_{L-1,T}^k\right)^\top,
\end{aligned}$$

where $v_{2,i}^k := \alpha_i(W_L^{k+1} G_{L-1,T}^k)_i + (1 - \alpha_i)(W_L^k G_{L-1,T}^k)_i$ is an interior point between the corresponding entries of $W_L^{k+1} G_{L-1,T}^k$ and $W_L^k G_{L-1,T}^k$. $\textcircled{1}$ is from gradient descent formulation of W_L^k in Equation (8).

Substituting above into Equation (46) yields:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \\
&= \frac{\eta}{\beta} (\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y))^\top \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}\left(\frac{\partial F}{\partial W_L^k} G_{L-1,T}^k\right)^\top (\mathbf{M}^\top (\mathbf{M} X_T^k - Y)), \\
&= \frac{\eta}{\beta} \frac{\partial F}{\partial W_L^k} G_{L-1,T}^k \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) (\mathbf{M}^\top (\mathbf{M} X_T^k - Y)),
\end{aligned}$$

Further, we substitute the gradient formulation in Equation (8) and calculate:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \\
&= -\frac{\eta}{\beta^2} \sum_{t=1}^T (\mathbf{M}^\top (\mathbf{M} X_t^k - Y))^\top \left(\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j) \mathbf{M}^\top \mathbf{M} \right) \\
&\quad \mathcal{D}((\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y))) \mathcal{D}(P_t \odot (1 - P_t/2)) G_{L-1,t}^k G_{L-1,T}^k \quad (47) \\
&\quad \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) (\mathbf{M}^\top (\mathbf{M} X_T^k - Y)), \\
&= -\frac{\eta}{\beta^2} (\mathbf{M} X_T^k - Y)^\top \mathbf{M} \mathbf{B}_T^k \mathbf{M}^\top (\mathbf{M} X_T^k - Y),
\end{aligned}$$

where \mathbf{B}_T^k is defined by:

$$\begin{aligned}
& \mathbf{B}_T^k \\
&= \sum_{t=1}^T \left(\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M} \right) \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{t-1}^k - Y)) \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) G_{L-1,t}^k \\
&\quad G_{L-1,T}^k \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)).
\end{aligned}$$

We discuss the definite property of \mathbf{B}_T^k case-by-case.

Case 1: $t = T$. $\Pi_{j=T}^{T+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j) \mathbf{M}^\top \mathbf{M}$ degenerates to be 1. The Equation (47) becomes:

$$\begin{aligned}
& [(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y)]_{\text{Part 1}} \\
&= -\frac{\eta}{\beta^2} (\mathbf{M} X_T^k - Y)^\top \mathbf{M} \\
&\quad \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \\
&\quad \mathcal{D}(P_T^k \odot (1 - P_T^k/2)) \\
&\quad G_{L-1,T}^k{}^\top G_{L-1,T}^k \\
&\quad \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \\
&\quad \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \mathbf{M}^\top (\mathbf{M} X_T^k - Y),
\end{aligned} \tag{48}$$

We first present the following corollary to show that there exists a negative upper bound of $[(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y)]_{\text{Part 1}}$:

Corollary A.11. *RHS of Equation (48) < 0 if $\lambda_{\min}(G_{L-1,T}^k{}^\top G_{L-1,T}^k) > 0$.*

Proof. Due to definition of eigenvalue and Cauchy-Schwarz inequality, we calculate:

$$\begin{aligned}
& (\mathbf{M} X_T^k - Y)^\top \mathbf{M} \\
& \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \\
& \mathcal{D}(P_T^k \odot (1 - P_T^k/2)) G_{L-1,T}^k{}^\top G_{L-1,T}^k \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \\
& \mathcal{D}(\mathbf{M}^\top (\mathbf{M} X_{T-1}^k - Y)) \mathbf{M}^\top (\mathbf{M} X_T^k - Y), \\
& \geq (P_T^k \odot (1 - P_T^k/2))_{\min} ([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min} \\
& \quad \lambda_{\min}(G_{L-1,T}^k{}^\top G_{L-1,T}^k) \lambda_{\min}(\mathbf{M} \mathbf{M}^\top) \|\mathbf{M}^\top (\mathbf{M} X_T^k - Y)\|_2^2, \\
& \stackrel{\textcircled{1}}{>} 0,
\end{aligned}$$

where $\textcircled{1}$ is due to Sigmoid function is non-negative, $\lambda_{\min}(G_{L-1,T}^k{}^\top G_{L-1,T}^k) > 0$, and $\lambda_{\min}(\mathbf{M} \mathbf{M}^\top) > 0$ by definition. Thus, $(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y) < 0$ by nature. $()_{\min}$ means the minimal value among all entries. \square

To get an upper bound, we expect $G_{L-1,T}^k{}^\top G_{L-1,T}^k$ to be positive definition, in which we require $n_{L-1} \geq N$. Thus, we can easily get the upper bound from its minimal eigenvalue.

Based on Corollary A.11, we calculate the negative lower bound of Equation (48) by:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y) \\
& \leq -\frac{\eta}{\beta^2} (P_T^k \odot (1 - P_T^k/2))_{\min} ([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min} \\
& \quad \lambda_{\min}(G_{L-1,T}^k{}^\top G_{L-1,T}^k) \lambda_{\min}(\mathbf{M} \mathbf{M}^\top) \|\mathbf{M}^\top (\mathbf{M} X_T^k - Y)\|_2^2,
\end{aligned} \tag{49}$$

The remaining task is to calculate $(P_T^k \odot (1 - P_T^k/2))_{\min}$ and $([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min}$. We achieve that by calculating the values on the boundary of closed sets.

First, denote $v_3^k := W_L^k G_{L-1,T}^k$, we represent $P_T^k \odot (1 - P_T^k/2)$ by:

$$P_T^k \odot (1 - P_T^k/2) = 2\sigma(v_3^k)^\top \odot (1 - \sigma(v_3^k))^\top.$$

Since the Sigmoid function is a coordinate-wise non-decreasing function, we can straightforwardly find $([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min}$ and $(2\sigma(v_3^k)^\top \odot (1 - \sigma(v_3^k))^\top)_{\min}$ by on the closed sets of v_2^k and v_3^k , respectively, which is achieved by the following lemma.

Lemma A.12. For some $b, B \in \mathbb{R}^{k^2}$, $\forall v^k, b \leq v^k \leq B$, we calculate $(2\sigma(v^k)^\top \odot (1 - \sigma(v^k))^\top)_{\min}$ by:

$$(2\sigma(v^k)^\top \odot (1 - \sigma(v^k))^\top)_{\min} = \begin{cases} \min(2\sigma(b)(1 - \sigma(b))^\top, 2\sigma(B)(1 - \sigma(B))^\top) & -b \neq B, \\ 2\sigma(B)(1 - \sigma(B)) & -b = B. \end{cases}$$

Proof. Since $\sigma(x) \in (0, 1) \forall x$, $\mathcal{D}(2\sigma(x) \odot (1 - \sigma(x)))$ is a quadratic function w.r.t. x . Since $\sigma(x) \in (0, 1) \forall x$, $\mathcal{D}(2\sigma(x) \odot (1 - \sigma(x))) > 0$. Since the coefficient before the x^2 term is negative, its lower bound is either the value on the boundary or 0.

Since $\sigma(b), \sigma(B) \in (0, 1)$, if $-b \neq B$, the lower bound is the smaller one, i.e., $\min(2\sigma(b) \odot (1 - \sigma(b)), 2\sigma(B) \odot (1 - \sigma(B)))$. Otherwise, since both $\sigma(x)$ and $\mathcal{D}(2\sigma(x) \odot (1 - \sigma(x)))$ are symmetric around $\frac{1}{2}$, we have $2\sigma(B) \odot (1 - \sigma(B)) = 2\sigma(b) \odot (1 - \sigma(b))$. \square

Further, we calculate the bounds of v_2^k and v_3^k and invoke Lemma A.12 to get $([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min}$ and $(2\sigma(v_3^k)^\top \odot (1 - \sigma(v_3^k))^\top)_{\min}$.

We present the following two lemmas to show the closed sets that v_2^k and v_3^k belong to.

Lemma A.13. Denote $\ell \in [L]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$, we assume $\|W_\ell^k\|_2 \leq \bar{\lambda}_\ell$. We define the following quantity:

$$\delta_3 = ((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2),$$

$$\Theta_L = \prod_{\ell=1}^L \bar{\lambda}_\ell.$$

For $v_{2,i}^k := \alpha_i(W_L^{k+1}G_{L-1,T}^k)_i + (1 - \alpha_i)(W_L^k G_{L-1,T}^k)_i$, $\alpha_i \in [0, 1]$, v_2^k belongs to the following closed set:

$$v_2^k \in [-\delta_3 \Theta_L, \delta_3 \Theta_L].$$

Proof. We calculate v_2^k 's upper bound by:

$$\begin{aligned} \|v_2^k\|_\infty &= \|\alpha \odot (W_L^{k+1}G_{L-1,T}^k) + (1 - \alpha) \odot (W_L^k G_{L-1,T}^k)\|_\infty, \\ &= \max_i \|\alpha_i(W_L^{k+1}G_{L-1,T}^k)_i + (1 - \alpha_i)(W_L^k G_{L-1,T}^k)_i\|_\infty, \\ &\stackrel{\textcircled{1}}{\leq} \max_i \alpha_i \|(W_L^{k+1}G_{L-1,T}^k)_i\|_\infty + (1 - \alpha_i)\|(W_L^k G_{L-1,T}^k)_i\|_\infty, \\ &\stackrel{\textcircled{2}}{\leq} \max_i \max(\|(W_L^{k+1}G_{L-1,T}^k)_i\|_\infty, \|(W_L^k G_{L-1,T}^k)_i\|_\infty), \\ &= \max(\max_i \|(W_L^{k+1}G_{L-1,T}^k)_i\|_\infty, \max_i \|(W_L^k G_{L-1,T}^k)_i\|_\infty), \\ &\leq \max(\|W_L^{k+1}G_{L-1,T}^k\|_\infty, \|W_L^k G_{L-1,T}^k\|_\infty), \end{aligned} \tag{50}$$

where $\textcircled{1}$ is due to triangle inequality and $\textcircled{2}$ is due to $\alpha_i \in [0, 1]$ and upper bound of NN's inner output in Lemma A.5.

We calculate the bound of $\|W_L^{k+1}G_{L-1,T}^k\|_2$ by:

$$\begin{aligned} \|W_L^{k+1}G_{L-1,T}^k\|_\infty &\stackrel{\textcircled{1}}{\leq} \|W_L^{k+1}\|_2 \|G_{L-1,T}^k\|_2, \\ &\stackrel{\textcircled{2}}{\leq} \bar{\lambda}_L ((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2) \prod_{\ell=1}^{L-1} \bar{\lambda}_\ell, \\ &= \underbrace{((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2)}_{\delta_3} \underbrace{\prod_{\ell=1}^L \bar{\lambda}_\ell}_{\Theta_L}, \end{aligned}$$

where $\textcircled{1}$ is due to Cauchy-Schwarz inequality and $\textcircled{2}$ is due to definition and upper bound of NN's inner output in Lemma A.5. Similarly, we can get $\|W_L^{k+1}G_{L-1,T}^k\|_2 \leq \delta_3 \Theta_L$.

² \mathbb{R}^k means the space at training iteration k .

Substituting back to Equation (50) yields:

$$\|v_2^k\|_\infty \leq \delta_3 \Theta_L.$$

Thus, we have the following bound for vector v_2^k by nature:

$$-\delta_3 \Theta_L \leq v_2^k \leq \delta_3 \Theta_L.$$

It is worth noting that the above lower bound is non-trivial since we cannot have $v_2^k \geq 0$, which can be easily violated by a little perturbation from gradient descent. \square

Lemma A.14. Denote $\ell \in [L]$, for some $\bar{\lambda}_\ell \in \mathbb{R}$, we assume $\|W_\ell^k\|_2 \leq \bar{\lambda}_\ell$. We define the following quantity:

$$\begin{aligned} \delta_3 &= ((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2), \\ \Theta_L &= \prod_{\ell=1}^L \bar{\lambda}_\ell. \end{aligned}$$

For $v_3^k := W_L^k G_{L-1,T}^k, \forall k$, v_3^k belongs to the following closed set:

$$v_3^k \in [-\delta_3 \Theta_L, \delta_3 \Theta_L].$$

Proof. We prove the lemma by a similar method. We calculate the bound of $\|W_L^k G_{L-1,T}^k\|_2$ by:

$$\begin{aligned} \|v_3^k\|_\infty &= \|W_L^k G_{L-1,T}^k\|_\infty \\ &\stackrel{\textcircled{1}}{\leq} \|W_L^k\|_2 \|G_{L-1,T}^k\|_2, \\ &\stackrel{\textcircled{2}}{\leq} \bar{\lambda}_L ((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2) \prod_{\ell=1}^{L-1} \bar{\lambda}_\ell, \\ &= \underbrace{((1 + \beta)\|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2)}_{\delta_3} \underbrace{\prod_{\ell=1}^L \bar{\lambda}_\ell}_{\Theta_L}, \end{aligned}$$

where $\textcircled{1}$ is due to Cauchy-Schwarz inequality and $\textcircled{2}$ is due to definition and upper bound of NN's inner output in Lemma A.5.

We have the following bound for v_3^k by nature:

$$-\delta_3 \Theta_L \leq v_3^k \leq \delta_3 \Theta_L. \quad \square$$

We calculate $([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min}$ by substituting Lemma A.13 into Lemma A.12:

$$([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min} = 2\sigma(\delta_3 \Theta_L)(1 - \sigma(\delta_3 \Theta_L)).$$

Similarly, we get $(P_T^k \odot (1 - P_T^k/2))$ by substituting Lemma A.14 into Lemma A.12:

$$(P_T^k \odot (1 - P_T^k/2))_{\min} = 2\sigma(\delta_3 \Theta_L)(1 - \sigma(\delta_3 \Theta_L)).$$

Substituting the above results into Equation (49) and Equation (48) yields:

$$\begin{aligned} &[(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M} X_T^k - Y)]_{\text{Part 1}} \\ &\leq -\frac{\eta}{\beta^2} (P_T^k \odot (1 - P_T^k/2))_{\min} ([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)_{\min} \\ &\quad \lambda_{\min}(G_{L-1,T}^k)^\top G_{L-1,T}^k \lambda_{\min}(\mathbf{M} \mathbf{M}^\top) \|\mathbf{M}^\top (\mathbf{M} X_T^k - Y)\|_2^2, \\ &\leq -\frac{\eta}{\beta^2} 4\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \lambda_{\min}(G_{L-1,T}^k)^\top G_{L-1,T}^k \lambda_{\min}(\mathbf{M} \mathbf{M}^\top) \|\mathbf{M}^\top (\mathbf{M} X_T^k - Y)\|_2^2, \\ &\stackrel{\textcircled{1}}{\leq} -\eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta^2}{\beta^2} \alpha_0^2 \frac{1}{2} \|\mathbf{M} X_T^k - Y\|_2^2, \end{aligned} \tag{51}$$

where $\textcircled{1}$ is from definition.

Case 2: $t < T$. We derive the upper bound of above term by Cauchy-Schwarz inequality:

$$\begin{aligned}
& [(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y)]_{\text{Part 2}} \\
&= -\frac{\eta}{\beta^2} (\mathbf{M}X_T^k - Y)^\top \mathbf{M} \left(\sum_{t=1}^{T-1} (\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M}) \right. \\
&\quad \mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)) \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) G_{L-1,t}^k{}^\top G_{L-1,T}^k \\
&\quad \left. \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y)) \right) \mathbf{M}^\top (\mathbf{M}X_T^k - Y), \\
&\stackrel{\textcircled{1}}{\leq} \frac{\eta}{\beta^2} \left\| \sum_{t=1}^{T-1} (\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M}) \right\|_2 \\
&\quad \mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y)) \mathcal{D}(P_t^k \odot (1 - P_t^k/2)) G_{L-1,t}^k{}^\top G_{L-1,T}^k \\
&\quad \mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top) \mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y)) \Big\|_2 \|\mathbf{M}\mathbf{M}^\top\|_2 \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\leq \frac{\eta}{\beta^2} \sum_{t=1}^{T-1} \left\| (\prod_{j=T}^{t+1} \mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_j^k) \mathbf{M}^\top \mathbf{M}) \right\|_2 \\
&\quad \|\mathcal{D}(P_t^k \odot (1 - P_t^k/2))\|_2 \|G_{L-1,t}^k\|_2 \|G_{L-1,T}^k\|_2 \|\mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)\|_2 \\
&\quad \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y))\|_2 \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y))\|_2 \|\mathbf{M}\mathbf{M}^\top\|_2 \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\stackrel{\textcircled{2}}{\leq} \frac{\eta}{\beta} \sum_{t=1}^{T-1} \|\mathcal{D}(P_t^k \odot (1 - P_t^k/2))\|_2 \|G_{L-1,t}^k\|_2 \|G_{L-1,T}^k\|_2 \|\mathcal{D}([2\sigma(v_{2,i}^k)(1 - \sigma(v_{2,i}^k))]^\top)\|_2 \\
&\quad \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y))\|_2 \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y))\|_2 \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\stackrel{\textcircled{3}}{\leq} \frac{\eta}{4\beta} \sum_{t=1}^{T-1} \|G_{L-1,t}^k\|_2 \|G_{L-1,T}^k\|_2 \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{t-1}^k - Y))\|_2 \|\mathcal{D}(\mathbf{M}^\top (\mathbf{M}X_{T-1}^k - Y))\|_2 \\
&\quad \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\leq \frac{\eta}{4\beta} (\beta \|X_0\|_2 + \frac{2T}{\beta} \|\mathbf{M}^\top Y\|_2) + \|\mathbf{M}^\top Y\|_2 \|G_{L-1,T}^k\|_2 \\
&\quad \sum_{t=1}^{T-1} \|G_{L-1,t}^k\|_2 (\beta \|X_0\|_2 + \frac{2t}{\beta} \|\mathbf{M}^\top Y\|_2) + \|\mathbf{M}^\top Y\|_2 \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\leq \frac{\eta}{4\beta} (\beta \|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2) + \|\mathbf{M}^\top Y\|_2 ((1 + \beta) \|X_0\|_2 + (2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2) \\
&\quad \prod_{s=1}^{L-1} \bar{\lambda}_s \sum_{t=1}^{T-1} ((1 + \beta) \|X_0\|_2 + (2t - 1 + \frac{2t-2}{\beta}) \|\mathbf{M}^\top Y\|_2) \\
&\quad \prod_{s=1}^{L-1} \bar{\lambda}_s (\beta \|X_0\|_2 + \frac{2t-2}{\beta} \|\mathbf{M}^\top Y\|_2) + \|\mathbf{M}^\top Y\|_2 \|\mathbf{M}X_T^k - Y\|_2^2,
\end{aligned}$$

where ① is due to Cauchy-Schwarz inequality. It is worth noting that ① is non-trivial since \mathbf{B}_{T-1}^k is non-necessarily to be positive definite. ② is due to upper bound of NN's output in Lemma A.1. ③ is based on the Sigmoid function is bounded.

Further, due to the definition of the quantities, we calculate:

$$\begin{aligned}
& [(Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y)]_{\text{Part 2}} \\
&\leq \frac{\eta\beta}{4} \\
&\quad \underbrace{((1 + \beta) \|X_0\|_2^2 + ((4T - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 + \frac{(2T-1)(\beta(2T-1) + (2T-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2)}_{\Lambda_T} \\
&\quad \sum_{t=1}^{T-1} \\
&\quad \underbrace{((1 + \beta) \|X_0\|_2^2 + ((4t - 3)(1 + \frac{1}{\beta}) + 1) \|X_0\|_2 \|\mathbf{M}^\top Y\|_2 + \frac{(2T-1)(\beta(2T-1) + (2T-2))}{\beta^2} \|\mathbf{M}^\top Y\|_2^2)}_{\Lambda_t} \\
&\quad \Theta_{L-1}^2 \|\mathbf{M}X_T^k - Y\|_2^2, \\
&= \frac{\eta\beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2.
\end{aligned} \tag{52}$$

Combining the two parts in Equation (51) and Equation (52) yields:

$$\begin{aligned}
& (Z - X_T^k)^\top \mathbf{M}^\top (\mathbf{M}X_T^k - Y) \\
&\leq \left(\frac{\eta\beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t - \eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 \right) \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2.
\end{aligned}$$

□

Using quantities from Equation (12), substituting the upper bounds in Lemma A.8, Lemma A.9, and Lemma A.10 into Equation (39), we calculate:

$$\begin{aligned}
& F([W]^{k+1}) \\
&= F([W]^k) + \frac{1}{2} \|\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k\|_2^2 + (\mathbf{M}X_T^{k+1} - \mathbf{M}X_T^k)^\top (\mathbf{M}X_T^k - Y), \\
&\leq F([W]^k) + \frac{\beta^2 \eta^2}{16} (\delta_1^T)^2 \left(S_{\Lambda, T} \right)^2 \left(\Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)^2 \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2 \\
&\quad + \frac{\beta \eta}{2} (\Lambda_T + \delta_2) \Theta_L^2 S_{\bar{\lambda}, L} S_{\Lambda, T} \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2 \\
&\quad + \left(-\eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 + \frac{\eta \beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t \right) \frac{1}{2} \|\mathbf{M}X_T^k - Y\|_2^2, \\
&\stackrel{\textcircled{1}}{=} F([W]^k) + \frac{\beta^2 \eta^2}{16} (\delta_1^T)^2 \left(S_{\Lambda, T} \right)^2 \left(\Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)^2 F([W]^k) \\
&\quad + \frac{\beta \eta}{2} (\Lambda_T + \delta_2) \Theta_L^2 S_{\bar{\lambda}, L} S_{\Lambda, T} F([W]^k) \\
&\quad + \left(-\eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 + \frac{\eta \beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t \right) F([W]^k), \\
&= \left(1 + \frac{\eta^2 \beta^2}{16} (\delta_1^T)^2 \left(S_{\Lambda, T} \right)^2 \left(\Theta_L^2 \sum_{\ell=1}^L \bar{\lambda}_\ell^{-2} \right)^2 + \frac{\eta \beta}{2} (\Lambda_T + \delta_2) S_{\Lambda, T} \Theta_L^2 S_{\bar{\lambda}, L} \right. \\
&\quad \left. + \frac{\eta \beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t - \eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 \right) F([W]^k), \\
&\stackrel{\textcircled{2}}{\leq} \left(1 + \eta \beta (\Lambda_T + \delta_2) S_{\Lambda, T} \Theta_L^2 S_{\bar{\lambda}, L} + \frac{\eta \beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t - \eta 8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 \right) \\
&\quad F([W]^k), \\
&= \left(1 - \eta (8\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 - \beta (\Lambda_T + \delta_2) S_{\Lambda, T} \Theta_L^2 S_{\bar{\lambda}, L} - \frac{\beta}{2} \Theta_{L-1}^2 \Lambda_T \sum_{t=1}^{T-1} \Lambda_t) \right) \\
&\quad F([W]^k), \\
&\stackrel{\textcircled{3}}{\leq} \underbrace{\left(1 - \eta 4\sigma(\delta_3 \Theta_L)^2 (1 - \sigma(\delta_3 \Theta_L))^2 \frac{\beta_0^2}{\beta^2} \alpha_0^2 \right)}_{4\eta \frac{\beta_0^2}{\beta^2} \delta_4} F([W]^k),
\end{aligned}$$

where ① is due to the definition of objective. ② is due to upper bound of learning rate in Equation (14a) and $\delta_1^T = \delta_2 + \sum_{j=1}^T \Lambda_j$ in definition. ③ is due to the lower bound of the least eigenvalue α_0 in Equation (13b).

Due to learning rate's upper bound in Equation (14b), we know $0 < 1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4 < 1$, which yields the following linear rate by nature:

$$F([W]^k) \leq (1 - \eta 4\eta \frac{\beta_0^2}{\beta^2} \delta_4)^k F([W]^0).$$

□

B Details for Initialization

B.1 Preliminary

To begin with, we define the following quantities:

$$\begin{aligned}
\delta_5 &= \sigma \left((2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2 \Theta_L \right)^{-2} \left(1 - \sigma \left((2T - 1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2 \Theta_L \right) \right)^{-2}, \\
\delta_6 &= \sigma_{\min} \left(\left(\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M})^{T-t} \mathbf{M}^\top Y \right) \mathbf{M}^\top (\mathbf{M} (\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M})^{T-t} \mathbf{M}^\top Y) - Y) \right), \\
\delta_7 &= \sigma_{\min} \left(\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M})^{T-t} \right).
\end{aligned}$$

Analysis for the numerical stability of δ_5 . δ_5 is a function with Λ_t , which is also enlarged w.r.t. e^L . In general, it is possible to push $\sigma(1 - \sigma((2T - 1 + \frac{2T-2}{\beta})\|\mathbf{M}^\top Y\|_2 \Theta_L))$ to zero and let RHS of above inequality to be ∞ when $e^L \rightarrow \infty$. As presented in the lemma, we claim that the required e is not necessarily to be ∞ . Thus, δ_5 can be regarded as a $\mathcal{O}(e^{L-1}) \ll \infty$ constant. In the following proofs, we demonstrate that it holds since e is finite.

We calculate the following exact formulations of the quantities defined in Theorem 4.3:

$$\begin{aligned}
\Lambda_T &= (1 + \beta)\|X_0\|_2^2 + ((4T - 3)(1 + \frac{1}{\beta}) + 1)\|X_0\|_2\|\mathbf{M}^\top Y\|_2 \\
&\quad + \frac{(2T-1)(\beta(2T-1)+(2T-2))}{\beta^2}\|\mathbf{M}^\top Y\|_2^2, \\
&= \frac{4(\beta+1)}{\beta^2}\|\mathbf{M}^\top Y\|_2^2 T^2 + \left(\frac{4(1+\beta)}{\beta}\|X_0\|_2\|\mathbf{M}^\top Y\|_2 - \frac{4\beta+6}{\beta^2}\|\mathbf{M}^\top Y\|_2^2\right)T \\
&\quad + (1 + \beta)\|X_0\|_2^2 - (2 + \frac{3}{\beta})\|X_0\|_2\|\mathbf{M}^\top Y\|_2 + \frac{\beta+2}{\beta^2}\|\mathbf{M}^\top Y\|_2^2, \\
&\stackrel{\textcircled{1}}{=} \frac{4(\beta+1)}{\beta^2}\|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{4\beta+6}{\beta^2}\|\mathbf{M}^\top Y\|_2^2 T + \frac{\beta+2}{\beta^2}\|\mathbf{M}^\top Y\|_2^2,
\end{aligned} \tag{53}$$

where $\textcircled{1}$ is due to $X_0 = 0$ and

$$\begin{aligned}
\sum_{i=1}^T \Lambda_i &= \sum_{i=1}^T (1 + \beta)\|X_0\|_2^2 + ((4i - 3)(1 + \frac{1}{\beta}) + 1)\|X_0\|_2\|\mathbf{M}^\top Y\|_2 \\
&\quad + \frac{(2i-1)(\beta(2i-1)+(2i-2))}{\beta^2}\|\mathbf{M}^\top Y\|_2^2 \\
&= \frac{4(\beta+1)}{3\beta^2}\|\mathbf{M}^\top Y\|_2^2 T^3 + \left(\frac{2(1+\beta)}{\beta}\|X_0\|_2\|\mathbf{M}^\top Y\|_2 - \frac{1}{\beta^2}\|\mathbf{M}^\top Y\|_2^2\right)T^2 \\
&\quad + \left((1 + \beta)\|X_0\|_2^2 - \frac{1}{\beta}\|X_0\|_2\|\mathbf{M}^\top Y\|_2 - \frac{\beta+1}{3\beta^2}\|\mathbf{M}^\top Y\|_2^2\right)T, \\
&\stackrel{\textcircled{1}}{=} \frac{4(\beta+1)}{3\beta^2}\|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2}\|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2}\|\mathbf{M}^\top Y\|_2^2 T,
\end{aligned} \tag{54}$$

where $\textcircled{1}$ is due to $X_0 = 0$.

Then, we analyze the expansion of $\sigma_{\min}(G_{L-1,T}^0)$ w.r.t. $[W]_L = e[W]_L$. Due to the one line form equation of L2O model in Equation (22), we have $\sigma_{\min}(G_{L-1,T}^0)$ is calculated by:

$$\sigma_{\min}(G_{L-1,T}^0) = \sigma_{\min}(\text{ReLU}(\text{ReLU}([X_{T-1}^0, \mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)]W_1^{0\top}) \cdots W_{L-1}^{0\top})),$$

where due to Equation (22), X_{T-1}^0 is given by:

$$\begin{aligned}
X_{T-1}^0 &= \prod_{t=T-1}^1 (\mathbf{I} - \frac{1}{\beta}\mathcal{D}(P_t^0)\mathbf{M}^\top\mathbf{M})X_0 + \frac{1}{\beta}\sum_{t=1}^{T-1} \prod_{s=T-1}^{t+1} (\mathbf{I} - \frac{1}{\beta}\mathcal{D}(P_s^0)\mathbf{M}^\top\mathbf{M})\mathcal{D}(P_t^0)\mathbf{M}^\top Y, \\
&\stackrel{\textcircled{1}}{=} (\mathbf{I} - \frac{1}{\beta}\mathbf{M}^\top\mathbf{M})^{T-1}X_0 + \frac{1}{\beta}\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta}\mathbf{M}^\top\mathbf{M})^{T-t}\mathbf{M}^\top Y, \\
&\stackrel{\textcircled{2}}{=} \frac{1}{\beta}\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta}\mathbf{M}^\top\mathbf{M})^{T-t}\mathbf{M}^\top Y,
\end{aligned} \tag{55}$$

where $\textcircled{1}$ is due to $P_t = \sigma(\mathbf{0}) = \mathbf{I}$ since $W_L = 0$. The result shows that X_{T-1}^0 is unrelated to $[W]_L$ with $W_L = 0$. $\textcircled{2}$ is due to $X_0 = 0$.

Further, for $t \in [T]$, denote the angle between X_{t-1}^0 and $\mathbf{M}^\top(\mathbf{M}X_{t-1}^0 - Y)$ as θ_{t-1} , we have $\sin(\theta_{t-1}) \in (0, 1)$, setting $[W]_L = e[W]_L$, we calculate $\sigma_{\min}(\tilde{G}_{L-1,T}^0)$ by:

$$\begin{aligned}
\sigma_{\min}(\tilde{G}_{L-1,T}^0) &= \sigma_{\min}(\text{ReLU}(\text{ReLU}([X_{T-1}^0, \mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)]eW_1^{0\top}) \cdots eW_{L-1}^{0\top})), \\
&\geq \sigma_{\min}([X_{T-1}^0 | \mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)]) \prod_{\ell=1}^{L-1} \sigma_{\min}(eW_\ell^0), \\
&\geq \frac{\|X_{T-1}^0\|_2 \|\mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)\|_2 \sin(\theta_{T-1})}{\|X_{T-1}^0\|_2 + \|\mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)\|_2} \prod_{\ell=1}^{L-1} \sigma_{\min}(eW_\ell^0), \\
&= \frac{\sin(\theta_{T-1})}{\frac{1}{\|X_{T-1}^0\|_2} + \frac{1}{\|\mathbf{M}^\top(\mathbf{M}X_{T-1}^0 - Y)\|_2}} \prod_{\ell=1}^{L-1} \sigma_{\min}(eW_\ell^0), \\
&\geq \sin(\theta_{T-1}) \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0) \Theta_L \|X_{T-1}^0\|_2.
\end{aligned} \tag{56}$$

Based on the definition of X_{T-1}^0 in Equation (55), we calculate following bound:

$$\begin{aligned}\sigma_{\min}(\tilde{G}_{L-1,T}^0) &\geq \frac{\sin(\theta_{T-1})}{\beta} \left\| \sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M})^{T-t} \mathbf{M}^\top Y \right\|_2 \prod_{\ell=1}^{L-1} \sigma_{\min}(eW_\ell^0), \\ &\geq \frac{\sin(\theta_{T-1})}{\beta} \underbrace{\sigma_{\min}(\sum_{t=1}^{T-1} (\mathbf{I} - \frac{1}{\beta} \mathbf{M}^\top \mathbf{M})^{T-t})}_{\delta_7} \left\| \mathbf{M}^\top Y \right\|_2 e^{L-1} \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0),\end{aligned}\tag{57}$$

where X_{T-1}^0 is a constant related to problem definition.

Substituting Equation (55), we calculate a tighter lower bound of $\|X_{T-1}^0\|_2$ by:

$$\begin{aligned}\|X_{T-1}^0\|_2 &= \left\| \frac{1}{\beta} \sum_{t=1}^{T-1} \prod_{s=T-1}^{t+1} (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M}) \mathcal{D}(P_t^0) \mathbf{M}^\top Y \right\|_2, \\ &\geq \frac{1}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \sigma_{\min} \left(\sum_{t=1}^{T-1} \prod_{s=T-1}^{t+1} (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M}) \mathcal{D}(P_t^0) \right), \\ &\stackrel{\textcircled{1}}{\geq} \frac{1}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \sum_{t=1}^{T-1} \sigma_{\min} \left(\prod_{s=T-1}^{t+1} (\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M}) \right) \sigma_{\min}(\mathcal{D}(P_t^0)), \\ &\geq \frac{1}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \sum_{t=1}^{T-1} \left(\prod_{s=T-1}^{t+1} \sigma_{\min}(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M}) \right) \sigma_{\min}(\mathcal{D}(P_t^0)),\end{aligned}\tag{58}$$

where $\textcircled{1}$ is due to all matrices in the summation are positive semi-definite by definition.

We calculate lower bound for $\sigma_{\min}(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M})$ by:

$$\begin{aligned}\sigma_{\min}(\mathbf{I} - \frac{1}{\beta} \mathcal{D}(P_s^0) \mathbf{M}^\top \mathbf{M}) &\geq 1 - \frac{1}{\beta} \sigma_{\max}(\mathcal{D}(2\sigma(eW_L^0 \tilde{G}_{L-1,s}^0)) \mathbf{M}^\top \mathbf{M}) \\ &\geq 1 - 2 \underbrace{\sigma(\delta_3 \Theta_L)(1 - \sigma(\delta_3 \Theta_L))}_{\delta_4} \sigma_{\max}(eW_L^0 \tilde{G}_{L-1,s}^0),\end{aligned}\tag{59}$$

It is easy to verify that the above equation equal to 1 when $e \rightarrow +\infty$ and it decreases with e . Also, a large e ensures the RHS of above inequality to be positive.

Similarly, we calculate lower bound for $\sigma_{\min}(P_t^0)$ by:

$$\begin{aligned}\sigma_{\min}(\mathcal{D}(P_t^0)) &\stackrel{\textcircled{1}}{=} \min(2\sigma(eW_L^0 \tilde{G}_{L-1,t}^0)), \\ &\stackrel{\textcircled{2}}{=} \min\left(\frac{\partial 2\sigma}{\partial v_4}(eW_L^0 \tilde{G}_{L-1,t}^0)\right), \\ &\stackrel{\textcircled{3}}{\geq} 2\delta_4 \sigma_{\min}(eW_L^0 \tilde{G}_{L-1,t}^0), \\ &\stackrel{\textcircled{4}}{\geq} 2\delta_4 e \|W_L^0\|_2 \sigma_{\min}(\tilde{G}_{L-1,t}^0), \\ &\geq 2\Theta_L \delta_4 \prod_{\ell=1}^L \|W_\ell^0\|_2 \sigma_{\min}([X_{t-1}^0 | \mathbf{M}^\top (\mathbf{M}X_{t-1}^0 - Y)]), \\ &\stackrel{\textcircled{5}}{\geq} 2\Theta_L \delta_4 \prod_{\ell=1}^L \|W_\ell^0\|_2 \sin(\theta_{T-1}) \|X_{t-1}^0\|_2\end{aligned}\tag{60}$$

where $\textcircled{1}$ means we apply the expansion here. $\textcircled{2}$ is due to Mean Value Theorem and v_4 denotes a inner point between 0 and $eW_L^0 \tilde{G}_{L-1,T}^0$. $\textcircled{3}$ is due to Lemma A.12 and Lemma A.14. $\textcircled{4}$ is due to W_L^0 is a vector in definition. $\textcircled{5}$ is similar to the workflow in Equation (56).

Substituting Equation (59) and Equation (60) back into Equation (58) yields:

$$\begin{aligned}\|X_{t-1}^0\|_2 &\geq \frac{1}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \sum_{s=1}^{t-1} 2\Theta_L \delta_4 \prod_{\ell=1}^L \|W_\ell^0\|_2 \sigma_{\min}([X_{s-1}^0 | \mathbf{M}^\top (\mathbf{M}X_{s-1}^0 - Y)]), \\ &\geq \frac{2}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \Theta_L \sum_{s=1}^{t-1} \delta_4 \prod_{\ell=1}^L \|W_\ell^0\|_2 \sin(\theta_{s-1}) \|X_{s-1}^0\|_2,\end{aligned}$$

Similarly, we can get the following lower bound of $\|X_{t-1}^0\|_2$:

$$\|X_{t-1}^0\|_2 \geq \frac{2}{\beta} \left\| \mathbf{M}^\top Y \right\|_2 \Theta_L \sum_{s=1}^{t-1} \delta_4 \prod_{\ell=1}^L \|W_\ell^0\|_2 \sin(\theta_{t-1}) \|X_{s-1}^0\|_2,$$

Based on the above results, we calculate the Ω of $\|X_{T-1}^0\|_2$ as in terms of T and Θ_L as:

$$\|X_{T-1}^0\|_2 \geq \underbrace{\Omega(\Theta_L \sum_{t=1}^{T-1} \Theta_L \sum_{s=1}^{t-1} \Theta_L \sum_{j=1}^{s-1} \cdots \sum_{j=1}^2)}_{T-2 \text{ terms}} = \Omega(\Theta_L^{T-2}).$$

Substituting back into Equation (56) yields:

$$\sigma_{\min}(\tilde{G}_{L-1,T}^0) = \Omega(e^{L-1} e^{(T-2)(L-1)}) = \Omega(e^{(T-1)(L-1)}). \quad (61)$$

B.2 Proof of Lemma 5.1

Proof. Making up the lower bounding relationship with Equation (57) and Equation (62) yields:

$$\begin{aligned} e^{L-1} \|\mathbf{M}^\top Y\|_2 \delta_7 \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0) &\geq 8(1+\beta)(\|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2), \\ &= \frac{8(1+\beta)}{\beta} (2T-2) \|\mathbf{M}^\top Y\|_2, \end{aligned}$$

which yields:

$$e \geq \sqrt[L-1]{\frac{8(1+\beta)}{\beta} \delta_7^{-1} \sigma_{\min}(W_\ell^0)^{-1} (2T-2)}.$$

□

B.3 Proof of Lemma 5.4

We apply a similar workflow to prove Lemma 5.4.

Proof. With $X_0 = 0$, we find the upper bound of the RHS of Equation (13d) by substituting the quantity δ_5 :

$$\begin{aligned} &\frac{(1+\beta)\beta^2\sqrt{\beta}}{2\beta_0^2} \delta_5 (\sqrt{\beta} \|X_0\|_2 + (2T+1) \|Y\|_2) \zeta_2 S_{\Lambda,T} \Theta_{L-1} \left(\sum_{\ell=1}^L \frac{\Theta_\ell}{\bar{\lambda}_\ell^2} \right) \\ &\stackrel{\textcircled{1}}{=} \frac{(1+\beta)\beta^2\sqrt{\beta}}{2\beta_0^2} \delta_5 (\sqrt{\beta} \|X_0\|_2 + (2T+1) \|Y\|_2) \zeta_2 \\ &\quad \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right) \Theta_{L-1} \left(\sum_{\ell=1}^L \frac{\Theta_\ell}{\bar{\lambda}_\ell^2} \right), \\ &\stackrel{\textcircled{2}}{=} \frac{(1+\beta)\beta\sqrt{\beta}}{2\beta_0^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2 (2T-2)(2T+1) \\ &\quad \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right) \Theta_{L-1} \left(\sum_{\ell=1}^L \frac{\Theta_\ell}{\bar{\lambda}_\ell^2} \right), \\ &\stackrel{\textcircled{3}}{\leq} \frac{(1+\beta)\sqrt{\beta}}{6\beta_0^2\beta} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^3 \\ &\quad \left(16(\beta+1)T^5 - (8\beta+20)T^4 - 6(2\beta+1)T^3 + 2(\beta+4)T^2 + 2(\beta+1)T \right) L \Theta_{L-1}^2, \end{aligned} \quad (62)$$

where $\textcircled{1}$ is due to Equation (54) and definition of quantity δ_1^{T-1} in Theorem 4.3. $\textcircled{2}$ is due to $X_0 = 0$. $\textcircled{3}$ is due to $\bar{\lambda}_L = 1$ and $\bar{\lambda}_\ell > 1, \ell \in [L-1]$.

Making up the lower bounding relationship with Equation (57) and Equation (62) yields:

$$\begin{aligned} &\left(e^{L-1} \|\mathbf{M}^\top Y\|_2 \delta_7 \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0) \right)^3 \\ &\geq e^{2L-2} \frac{(1+\beta)\sqrt{\beta}}{6\beta_0^2\beta} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^3 L \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1)^2 \\ &\quad \left(16(\beta+1)T^5 - (8\beta+20)T^4 - 6(2\beta+1)T^3 + 2(\beta+4)T^2 + 2(\beta+1)T \right), \end{aligned} \quad (63)$$

which yields:

$$e \geq \sqrt[L-1]{C_{2,\delta_5} \left(16(\beta+1)T^5 - (8\beta+20)T^4 - 6(2\beta+1)T^3 + 2(\beta+4)T^2 + 2(\beta+1)T \right)}.$$

where C_{2,δ_5} denotes the $\frac{(1+\beta)\sqrt{\beta}}{6\beta_0^2\beta\delta_5^3} \delta_5 \|Y\|_2 L \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1)^2 \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0)^{-3}$ term.

Similarly, the finite RHS of above inequality ensures $\delta_5 \ll \infty$. □

B.4 Proof of Lemma 5.2

Proof. Using quantities from Equation (12), with $X_0 = 0$, we find the upper bound of the RHS of Equation (13b) by substituting the quantity δ_5 :

$$\begin{aligned}
& \frac{\beta^3}{4\beta_0^2} \delta_5 \left(-\frac{1}{2} \Theta_{L-1}^2 \Lambda_T \left(\sum_{t=1}^{T-1} \Lambda_t \right) + \Theta_L^2 S_{\bar{\lambda},L} (\Lambda_T + \delta_2) S_{\Lambda,T} \right) \\
& \stackrel{\textcircled{1}}{=} \frac{\beta^3}{4\beta_0^2} \delta_5 \left(-\frac{1}{2} \Theta_{L-1}^2 \left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \right. \\
& \quad \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 (T-1)^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 (T-1)^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 (T-1) \right) \\
& \quad + \Theta_L^2 S_{\bar{\lambda},L} \left(\left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \right. \\
& \quad \left. + \sum_{s=1}^{T-1} \left(\prod_{j=s+1}^T \left(1 + \frac{1+\beta}{2\beta} (2j-1) \Theta_L \|\mathbf{M}^\top Y\|_2 \right) \right. \right. \\
& \quad \left. \left. \left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \right) \right. \\
& \quad \left. \left. \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right) \right) \right) \\
& \leq \mathcal{O}(e^{2L-2} T^5 + e^{2L-4} T^5 + e^{2L-4} T^6 \sum_{s=1}^{T-1} s^2 \prod_{j=s+1}^T j e^{L-1}), \\
& = \mathcal{O}(e^{TL-T+2L-4} T^{3T+6}).
\end{aligned} \tag{64}$$

where ① is due to Equation (54) and definition of quantity δ_1^{T-1} in Theorem 4.3. ② is due to $X_0 = 0$. ③ is due to $\bar{\lambda}_L = 1$ and $\bar{\lambda}_\ell > 1, \ell \in [L-1]$.

Making up the lower bounding relationship with Equation (61) and Equation (62) yields:

$$(\Omega(e^{(T-1)(L-1)}))^2 \geq \mathcal{O}(e^{TL-T+2L-4} T^{3T+6}),$$

which yields:

$$e = \Omega(T^{\frac{3T+6}{TL-T-4L+6}}).$$

□

B.5 Proof of Lemma 5.3

Proof. Using quantities from Equation (12), with $X_0 = 0$, we find the upper bound of the RHS of Equation (13c) by substituting the quantity δ_5 :

$$\begin{aligned}
& \max_{\ell \in [L]} \frac{\Theta_L}{C_\ell \bar{\lambda}_\ell} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \\
& \underbrace{\sigma \left((2T-1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2 \Theta_L \right)^{-2} \left(1 - \sigma \left((2T-1 + \frac{2T-2}{\beta}) \|\mathbf{M}^\top Y\|_2 \Theta_L \right) \right)^{-2}}_{\delta_5} \\
& S_{\Lambda,T} (2T+1) \|Y\|_2, \\
& \stackrel{\textcircled{1}}{\leq} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 S_{\Lambda,T} (2T+1) \|Y\|_2 \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1), \\
& \stackrel{\textcircled{2}}{=} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right) \\
& (2T+1) \|Y\|_2 \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1), \\
& = \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^2 \left(\frac{8(\beta+1)}{3\beta^2} T^4 + \left(\frac{4(\beta+1)}{3\beta^2} - \frac{2}{\beta^2} \right) T^3 - \left(\frac{1}{\beta^2} + 2 \frac{\beta+1}{3\beta^2} \right) T^2 - \frac{\beta+1}{3\beta^2} T \right) \\
& \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1),
\end{aligned} \tag{65}$$

where ① is due to $\bar{\lambda}_\ell > 1, \ell \in [L-1]$ and $\bar{\lambda}_L = 1$. ② is due to Equation (54).

We analyze the two sides of the above inequality when $[W]_L = e[W]_L$ to demonstrate a sufficient lower bound of e to ensure Equation (65) holds.

If $[W]_L = e[W]_L$, since $e \geq 1$, Equation (65) is upper-bounded by:

$$\begin{aligned}
& \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^2 \left(\frac{8(\beta+1)}{3\beta^2} T^4 + \left(\frac{4(\beta+1)}{3\beta^2} - \frac{2}{\beta^2} \right) T^3 - \left(\frac{1}{\beta^2} + 2\frac{\beta+1}{3\beta^2} \right) T^2 - \frac{\beta+1}{3\beta^2} T \right) \\
& \prod_{\ell=1}^{L-1} (e \|W_\ell^0\|_2 + e) \\
& = e^{L-1} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^2 \\
& \quad \left(\frac{8(\beta+1)}{3\beta^2} T^4 + \left(\frac{4(\beta+1)}{3\beta^2} - \frac{2}{\beta^2} \right) T^3 - \left(\frac{1}{\beta^2} + 2\frac{\beta+1}{3\beta^2} \right) T^2 - \frac{\beta+1}{3\beta^2} T \right) \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1).
\end{aligned} \tag{66}$$

If RHS (lower bound) of Equation (57) greater than the RHS (upper bound) of above result, lower bound condition for minimal singular value in Equation (65) sufficiently holds, which yields:

$$\begin{aligned}
& \left(e^{L-1} \|\mathbf{M}^\top Y\|_2 \delta_7 \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0) \right)^2 \\
& \geq e^{L-1} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2 \delta_6^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^2 \left(\frac{8(\beta+1)}{3\beta^2} T^4 + \left(\frac{4(\beta+1)}{3\beta^2} - \frac{2}{\beta^2} \right) T^3 - \left(\frac{1}{\beta^2} + 2\frac{\beta+1}{3\beta^2} \right) T^2 - \frac{\beta+1}{3\beta^2} T \right) \\
& \quad \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1),
\end{aligned}$$

which yields:

$$e \geq \sqrt[L-1]{C_{1,\delta_5} \left(\frac{8(\beta+1)}{3} T^4 + \left(\frac{4(\beta+1)}{3} - 2 \right) T^3 - \left(1 + 2\frac{\beta+1}{3} \right) T^2 - \frac{\beta+1}{3} T \right)},$$

where C_{1,δ_5} denotes the $\frac{\sqrt{\beta}}{8\beta_0^2 \delta_6^2} \delta_5 \|Y\|_2 \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1) \prod_{\ell=1}^{L-1} \sigma_{\min}(W_\ell^0)^{-2}$ term, which is a ‘‘constant’’ w.r.t. δ_5 .

In the end, it is trivial to evaluate that the RHS of above δ_5 is finite with such e . \square

C Additional Experimental Results

In this section, we present detailed experimental settings and corresponding results. We define problems at three distinct scales, as described in Appendix C.1. The smaller scale is utilized for ablation studies (Section 6.2), whereas the larger scales are adopted for training experiments (Section 6.1 and Appendix C.2) and inference experiments (Appendix C.4).

C.1 Configurations for Different Experiments

Details of the three experimental configurations are presented in Table 1. **Scale 1** involves a DNN trained with input $X \in \mathbb{R}^{32 \times 32}$ and output $Y \in \mathbb{R}^{32 \times 25}$, featuring an $(L-1)$ -th layer dimension of 1024. **Scale 2** utilizes input $X \in \mathbb{R}^{10 \times 512}$ and output $Y \in \mathbb{R}^{10 \times 400}$, with the $(L-1)$ -th layer dimension established at 5120. **Scale 3** employs input $X \in \mathbb{R}^{2048 \times 512}$ and output $Y \in \mathbb{R}^{2048 \times 400}$. This configuration is designed as an under-parameterized system, with an $(L-1)$ -th layer dimension of 5120, specifically to evaluate the robustness of our proposed L2O framework. The third model, although targeting the optimization problem with the same dimension, has a different number of training samples N . We design the scale to align with the training configurations of the baseline model LISTA-CPSS [7]. Moreover, due to the GPU memory limitation, we set a thin NN, whose convergence is not guaranteed by our proposed theorem. The related experimental result is used to further demonstrate our proposed framework in Section 3.

Table 1: Configurations with Different Scales

Index	d	b	Dimension of $L-1$ Layer’s Output	Training Samples
1	32	25	1024	32
2	512	400	5120	10
3	512	400	20	2048

C.2 Additional Training Performance Comparisons verses L2O Baselines

For these experiments, the **Scale 3** configuration is utilized. Both baseline state-of-the-art (SOTA) methods and our proposed L2O framework are trained for 2000 epochs using a learning rate of 0.001. However, the inherent model construction and training scheme of a prominent SOTA method, LISTA-CPSS [7], diverges considerably from the requirements of our problem. Direct application of its original settings to our scenario results in over-fitting and poor training convergence, indicating a lack of robustness for this specific application. The following discussion elaborates on these incompatibilities and the modifications undertaken.

The original LISTA-CPSS framework possesses two key characteristics pertinent to this discussion. First, regarding its model construction, LISTA-CPSS addresses inverse problems by formulating a learnable Least Absolute Shrinkage and Selection Operator (LASSO) problem, wherein it learns a scalar coefficient for the L_1 regularization term [7]. However, our objective in Equation (1) is quadratic. Second, its training protocol is supervised, utilizing an L_2 loss against pre-generated optimal solutions, and employs a layer-wise training scheme. In this scheme, one layer is progressively added to the set of trainable parameters per training iteration, and these parameters are updated using four back-propagation (BP) steps [7]. To adapt LISTA-CPSS for our purposes, we modify both its model architecture and original training scheme to enable unsupervised optimization of our loss function (defined in Equation (2)) and to better align with our established training configuration.

First, to demonstrate the challenges of applying LISTA-CPSS’s original training paradigm to unsupervised quadratic objectives, we evaluate a minimally adapted version. This version is trained unsupervisedly by defining the loss as the objective function value from the final optimization step. Given our quadratic loss in Equation (2), any model components in LISTA-CPSS specifically designed for non-quadratic terms are not directly applicable. Moreover, a critical aspect of the publicly available LISTA-CPSS implementation is its initialization of the neural network (NN) with a fixed matrix \mathbf{M} . This initialization inherently restricts the trained model’s utility to problems featuring this identical, predetermined \mathbf{M} .

We train this minimally adapted LISTA-CPSS variant for 50 epochs (corresponding to 20000 BPs due to its layer-wise updates) using the Adam optimizer³ on a dataset of 2048 randomly generated samples. The loss function defined in Equation (2) is evaluated at an optimization step of $T = 100$. The experimental results, depicted in Figure 6, reveal that this configuration leads to severe over-fitting on the training samples. Specifically, Figure 6a illustrates the convergence of the objective function (at $T = 100$) as a function of the training iteration k . Concurrently, Figure 6b displays the mean objective value across 100 optimization steps during inference. These results indicate that while LISTA-CPSS achieves rapid convergence on the training data (which used a fixed \mathbf{M}), its performance degrades catastrophically (i.e., fails to generalize) when evaluated with a different matrix, \mathbf{M}' , during inference.

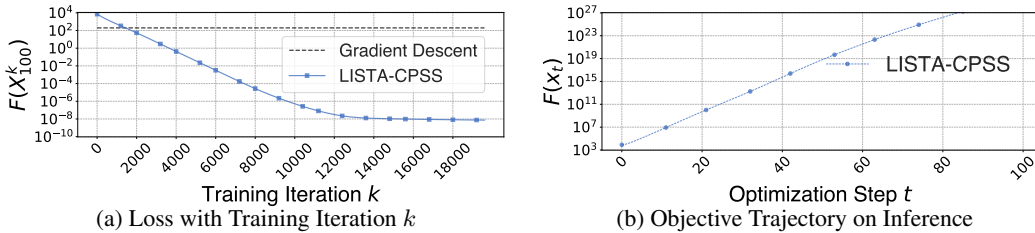


Figure 6: Training Loss and Inference Trajectory of LISTA-CPSS [7] with Fixed \mathbf{M}

Informed by the above observation, a more robust approach is achieved through the random initialization of LISTA-CPSS. Specifically, weights are sampled from a standard Gaussian distribution and subsequently scaled by a factor of $\frac{1}{d \cdot b}$ to mitigate potential numerical overflow in cumulative products. The LISTA-CPSS model is then trained using this initialization strategy.

For our proposed L2O framework, the expansion coefficient e is set to 100. As detailed in **Scale 3** in Table 1, we implement an under-parameterized system wherein the dimension of the $(L - 1)$ -th layer

³Our preliminary experiments indicates that SGD fails to converge with LISTA-CPSS’s original layer-wise training scheme.

is configured to 20. This implementation intentionally deviates from the theoretical requirements stipulated by our proposed theorems, which necessitate that the dimension of the $(L - 1)$ -th layer must be larger than the input dimension. This particular experiment is conducted to demonstrate the robustness of the proposed L2O framework, especially under such conditions that depart from our established theoretical framework.

The training losses of LISTA-CPSS and our proposed L2O framework are depicted in Figure 7, with the performance of non-learnable gradient descent (indicated by a horizontal line in the figure) serving as a baseline. Under scenarios with varied M configurations, LISTA-CPSS exhibits markedly slower convergence compared to both our proposed L2O framework and the gradient descent baseline. Moreover, the fast convergence observed for our L2O framework underscores the robustness and efficacy of its proposed initialization strategy, particularly when applied to under-parameterized models.

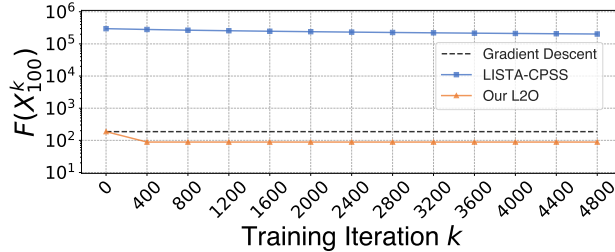


Figure 7: Training Losses with Varied M

C.3 Real-World Training Performance Comparisons

To empirically validate our proposed theorem, we perform an additional experiment comparing the training convergence of our L2O construction against standard Gradient Descent (GD). Utilizing a compact Convolutional Neural Network (CNN) on the MNIST dataset, our method achieved significantly faster convergence, thereby corroborating our theoretical findings.

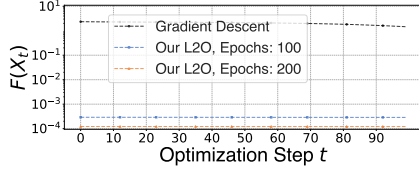
We employ the **Scale 3** configuration (an under-parameterized setting from Table 1). The CNN architecture (Table 2) comprises two convolutional layers, two max-pooling layers, ReLU activation functions, and a final linear layer. The optimization objective is the total cross-entropy loss over 200 randomly selected MNIST samples. The learning rates for training our L2O model and the CNN were set to 10^{-6} and 10^{-2} , respectively.

Table 2: Architecture of a Small CNN Model with MNIST Dataset

Layer	Input Channel	Output Channel	Kernel Size	Input Size	Output Size
Convolution	1	2	3	28×28	28×28
Max Pooling	2	2	2	28×28	14×14
ReLU	2	2	N/A	14×14	14×14
Convolution	2	3	3	14×14	14×14
Max Pooling	3	3	2	14×14	7×7
ReLU	3	3	N/A	7×7	7×7
Linear	147	10	N/A	1	1

To validate our framework, we conducted a comparative analysis of the CNN training loss on the MNIST dataset, contrasting our proposed L2O method with Gradient Descent (GD). The results are depicted in Figure 8a, which plots the training loss over 100 iterations. We evaluate two versions of our L2O optimizer, pre-trained for 100 and 200 epochs, respectively. In both scenarios, our L2O framework yields a substantially lower loss than the GD baseline, which corroborates the effectiveness of our approach for training DNN models.

Additionally, Figure 8b provides a quantitative comparison of the iteration cost for both methods. The proposed L2O framework converges to a more optimal (lower) loss value than GD in substantially fewer iterations, confirming its superior efficiency in training the CNN model.



(a) Training Losses

Method	Loss Value with Iterations
GD	10,000 Iterations: 2.92e-04
Our L2O, Epoch: 100	100 Iterations: 2.91e-04
Our L2O, Epoch: 200	100 Iterations: 1.22e-04

(b) Final Loss Values of CNN on MNIST Dataset

Figure 8: Performance of Training CNN on MNIST Dataset

C.4 Inference Experiment

Beyond analyzing training outcomes, we extend our evaluation to the robustness of the proposed L2O framework by assessing its performance in inference-stage optimization. This involves comparing the convergence characteristics of L2O against the Adam optimizer [10] and standard gradient descent (GD). It should be noted that while our theorems provide convergence guarantees for the training phase, such guarantees do not explicitly extend to this inference optimization context. For this empirical investigation, both our L2O framework and the Adam optimizer are executed across a range of hyperparameter settings for 3000 iterations (longer than 100 iterations in training), and their respective objective function trajectories are plotted as a function of the iteration count.

Adam utilizes momentum to accelerate gradient descent. In addition to the learning rate η , Adam employs two crucial hyperparameters, β_1 and β_2 , which control the exponential moving averages of past gradients and their squared magnitudes, respectively. For the Adam optimizer in our experiments, we set the learning rate $\eta = \frac{1}{\beta}$ (β -smoothness of objective) and explored hyper-parameters $\beta_1 \in \{0.1, 0.3, \dots, 0.9\}$ and $\beta_2 \in \{0.95, 0.955, \dots, 1.0\}$.

Regarding our proposed L2O framework and consistent with the initialization strategy detailed in Section 5, we selected a large expansion coefficient $e = 100$ to enhance training stability. The L2O model is then trained with learning rates η chosen from the set $\{10^{-3}, 10^{-4}, \dots, 10^{-7}\}$.

As illustrated in Figure 9, we present the objective trajectory over 3000 optimization steps, where each point is a mean value of 30 randomly generated problems' objectives. While the objective function initially exhibits rapid decay, the Adam optimizer fails to maintain this convergence, ultimately settling at sub-optimal values and not converging on average. In contrast, our proposed framework demonstrates superior performance compared to the Gradient Descent (GD) algorithm and exhibits robustness across various learning rates.

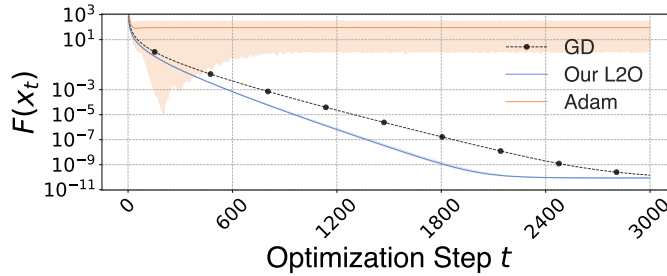


Figure 9: Inference Trajectory of Our Proposed L2O

C.5 Corollary in Ablation Studies

Corollary C.1 (LR's upper bound w.r.t. e).

$$\eta = \mathcal{O}(e^{3-L}T^{-6}) \cap \mathcal{O}(e^{1-L}T^{-4}) \cap \mathcal{O}(e^{\frac{4}{3}(1-L)}T^{-\frac{10}{3}}) \cap \mathcal{O}(e^{-TL-2L+T+4}T^{-3T-6}) \cap \mathcal{O}(T^{-2}).$$

Proof. From Equation (14a), we calculate:

$$\begin{aligned}
& \eta \\
& < \frac{8}{\beta} (\delta_2 + \Lambda_T) \left(\delta_2 + S_{\Lambda, T} \right)^{-1} S_{\Lambda, T}^{-2} \Theta_L^{-1} S_{\bar{\lambda}, L}^{-1}, \\
& < \left(\sum_{s=1}^{T-1} \left(\prod_{j=s+1}^T \left(1 + \frac{1+\beta}{2\beta} (2j-1) \Theta_L \|\mathbf{M}^\top Y\|_2 \right) \right. \right. \\
& \quad \left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \\
& \quad \left. + \left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \right) \\
& \quad \left(\sum_{s=1}^{T-1} \left(\prod_{j=s+1}^T \left(1 + \frac{1+\beta}{2\beta} (2j-1) \Theta_L \|\mathbf{M}^\top Y\|_2 \right) \right. \right. \\
& \quad \left(\frac{4(\beta+1)}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s^2 - \frac{4\beta+6}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 s + \frac{\beta+2}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 \right) \\
& \quad \left. + \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right) \right) \Big)^{-1} \\
& \quad \left(\frac{4(\beta+1)}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^3 - \frac{1}{\beta^2} \|\mathbf{M}^\top Y\|_2^2 T^2 - \frac{\beta+1}{3\beta^2} \|\mathbf{M}^\top Y\|_2^2 T \right)^{-2} \left(e^{L-1} \prod_{\ell=1}^{L-1} \bar{\lambda}_\ell \right)^{-1} S_{\bar{\lambda}, L}^{-1}, \\
& = \mathcal{O}(e^{3-L} T^{-6}).
\end{aligned}$$

From Equation (14b), due to the four lower bounds in Equation (13), we calculate following four upper bounds:

$$\begin{aligned}
& \eta \\
& < \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_4^{-2} \alpha_0^{-2}, \\
& \stackrel{66}{<} \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_5 \left(e^{L-1} \frac{\beta^2 \sqrt{\beta}}{8\beta_0^2} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^2 \right. \\
& \quad \left. \left(\frac{8(\beta+1)}{3\beta^2} T^4 + \left(\frac{4(\beta+1)}{3\beta^2} - \frac{2}{\beta^2} \right) T^3 - \left(\frac{1}{\beta^2} + 2\frac{\beta+1}{3\beta^2} \right) T^2 - \frac{\beta+1}{3\beta^2} T \right) \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1) \right)^{-1}, \\
& = \mathcal{O}(e^{1-L} T^{-4}).
\end{aligned}$$

$$\begin{aligned}
& \eta \\
& < \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_4^{-2} \alpha_0^{-2}, \\
& \stackrel{\text{Equation (63)}}{<} \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_5 \left(e^{2L-2} \frac{(1+\beta)\sqrt{\beta}}{6\beta_0^2\beta} \delta_5 \|Y\|_2 \|\mathbf{M}^\top Y\|_2^3 \right. \\
& \quad \left(16(\beta+1)T^5 - (8\beta+20)T^4 - 6(2\beta+1)T^3 + 2(\beta+4)T^2 + 2(\beta+1)T \right) \\
& \quad \left. L \prod_{\ell=1}^{L-1} (\|W_\ell^0\|_2 + 1)^2 \right)^{-\frac{2}{3}} \\
& = \mathcal{O}(e^{\frac{4}{3}(1-L)} T^{-\frac{10}{3}}).
\end{aligned}$$

$$\eta < \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_4^{-2} \alpha_0^{-2} \stackrel{\text{Equation (64)}}{<} \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_5 \mathcal{O}((e^{TL-T+2L-4} T^{3T+6})^{-1}) = \mathcal{O}(e^{-TL-2L+T+4} T^{-3T-6}).$$

$$\eta < \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_4^{-2} \alpha_0^{-2} \stackrel{\text{Equation (13a)}}{<} \frac{1}{4} \frac{\beta^2}{\beta_0^2} \delta_5 \left(8(1+\beta)(\|X_0\|_2 + \frac{2T-2}{\beta} \|\mathbf{M}^\top Y\|_2) \right)^{-2} = \mathcal{O}(T^{-2}).$$

□

C.6 Additional Ablation Studies for Learning Rates

We present two additional ablation studies with e of 25 and 100. Both use the configuration 1 in Table 1. The results are in Figure 10, which shows a deterministic relationship between LR and expansion coefficient. For $e = 25$ in Figure 10a, the 10^{-7} LR is too small and leads to worse optimality. The large LRs, i.e., 10^{-3} , 10^{-4} , cause unstable convergence. Similarly, for $e = 100$ in Figure 10b, a proper LR is 10^{-4} .

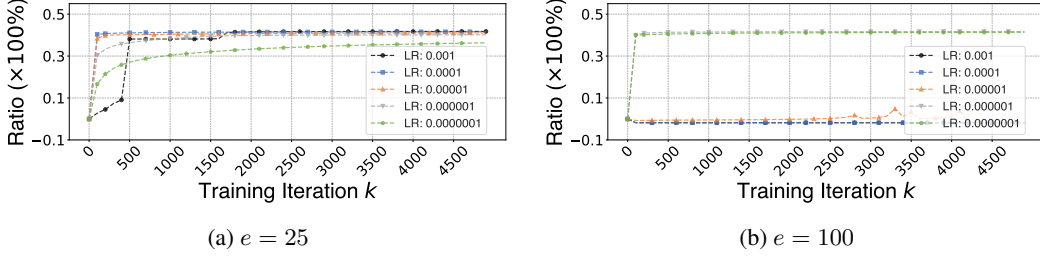


Figure 10: Additional Ablation Studies of Learning Rate with Different e .

C.7 Additional Ablation Studies for Expansion Coefficient e in Initialization

We present two additional ablation studies for e with learning rates of 0.001 and 0.00001. Both use the configuration 1 in Table 1. The results are in Figure 11. For a large LR, a large e may cause poor convergence due to Theorem 4.3. From Figure 11a, $e = 25$ is a proper setting for best convergence with $\eta = 0.001$. Similarly, for $\eta = 0.00001$, $e = 5$ is enough.

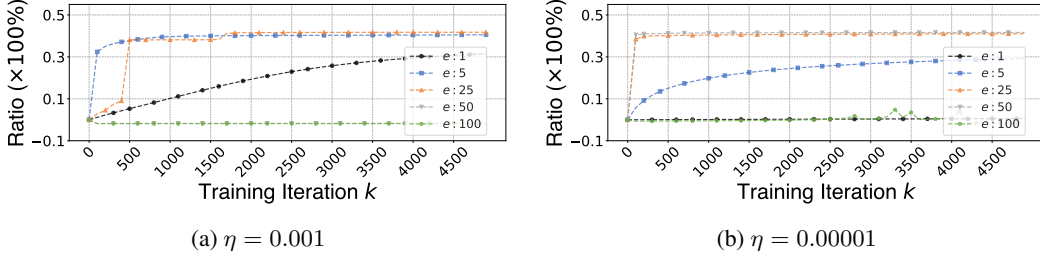


Figure 11: Additional Ablation Studies of e with Different Learning Rates.

D Discussion

Scope of Theoretical Guarantees. Our theoretical analysis establishes convergence guarantees and demonstrates superior convergence rates specifically for *over-parameterized* Math-L2O systems compared to baseline optimization algorithms. While we acknowledge the empirical effectiveness of certain *under-parameterized* Math-L2O systems [23, 34], providing theoretical convergence proofs for them remains challenging due to the inherent non-convexity of the underlying neural network training. Alternative theoretical approaches, such as convex dualization [17, 18, 31], have been explored. However, these methods typically necessitate the inclusion of regularization terms within the loss function, which may deviate from the original optimization objective we aim to solve.

Generalization to Other Objective Functions. The central thesis of Section 3 is that learning can enhance algorithmic convergence. To substantiate this claim, we first require a convergence guarantee for the neural network training process—a well-known complex problem. We leverage Neural Tangent Kernel (NTK) theory, which typically analyzes convergence under an L_2 -norm objective [16]. Despite generalizations of NTK to other loss functions [9, 40], we retain the L_2 -norm for two reasons: (1) it permits the derivation of an explicit convergence rate, rather than a surrogate one [40], and (2) it aids in demonstrating a deterministic initialization strategy, which has practical implications for model and training design.

Choice of Base Algorithm. Our framework utilizes Gradient Descent (GD) as the core algorithm primarily because it admits a direct analytical formulation relating the initial point X_0 to the iterate X_T . This tractability is crucial for our analysis. In contrast, accelerated variants like Nesterov Accelerated Gradient Descent (NAG) [4] generally lack such closed-form expressions for X_T . This absence significantly complicates the derivation of the output bounds required to analyze the L2O system’s dynamics and to prove convergence guarantees. Consequently, rigorously extending our current theoretical framework to momentum-based methods, despite attempts using inductive approaches, remains an open challenge.

We contend that a convergence proof for NAG can be constructed. Our central strategy involves bounding the L2O model’s output to satisfy the convergence conditions of the backbone algorithm. This is analogous to our use of the β -smoothness property to derive the step size in Equation (3) and is a methodology applicable to any provably convergent algorithm. To this end, we aim to bound X_T relative to X_0 . The proof proceeds as follows: First, NAG is formulated as a linear dynamical system where a transition matrix maps X_t to X_{t+1} . Second, we constrain the neural network outputs (i.e., momentum terms and step sizes) to ensure the transition matrix remains bounded over T steps. Finally, by applying the Cauchy-Schwarz and Triangle inequalities to this stable system, a formal bound on X_T is derived.

E Impact Statement

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

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: See Section 1.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: See Appendix D.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: See Appendix A.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: See Section 6 and Appendix A.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: If accepted, we will open source the codes.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 6 and Appendix A.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See Section 6 and Appendix A.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: See Section 6.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: See Appendix E.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See Appendix E.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: See Appendix E.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.