Learning quadratic neural networks in high dimensions: SGD dynamics and scaling laws

Gérard Ben Arous¹, Murat A. Erdogdu^{2,3}, Nuri Mert Vural^{2,3}, Denny Wu^{1,4}

¹New York University, ²University of Toronto, ³Vector Institute, ⁴Flatiron Institute gba@cims.nyu.edu, {erdogdu,vural}@cs.toronto.edu, dennywu@nyu.edu

Abstract

We study the optimization and sample complexity of gradient-based training of a two-layer neural network with quadratic activation function in the high-dimensional regime, where the data is generated as $y \propto \sum_{j=1}^r \lambda_j \sigma\left(\langle \pmb{\theta}_j, \pmb{x} \rangle\right)$, $\pmb{x} \sim \mathcal{N}(0, \pmb{I}_d)$, σ is the 2nd Hermite polynomial, and $\{\pmb{\theta}_j\}_{j=1}^r \subset \mathbb{R}^d$ are orthonormal signal directions. We consider the extensive-width regime $r \approx d^\beta$ for $\beta \in [0,1)$, and assume a power-law decay on the (non-negative) second-layer coefficients $\lambda_j \approx j^{-\alpha}$ for $\alpha \geq 0$. We present a sharp analysis of the SGD dynamics in the feature learning regime, for both the population limit and the finite-sample (online) discretization, and derive scaling laws for the prediction risk that highlight the power-law dependencies on the optimization time, sample size, and model width. Our analysis combines a precise characterization of the associated matrix Riccati differential equation with novel matrix monotonicity arguments to establish convergence guarantees for the infinite-dimensional effective dynamics.

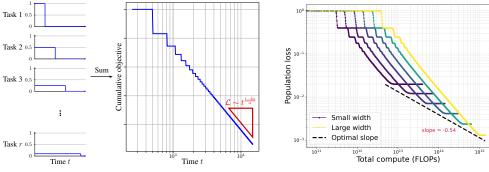
1 Introduction

We study the problem of learning a two-layer neural network (NN) with quadratic activation on isotropic Gaussian data. The target function (or the "teacher" model) is defined as

$$y = f_*(\boldsymbol{x}) + \epsilon \text{ with } f_*(\boldsymbol{x}) = \frac{1}{\|\boldsymbol{\Lambda}\|_F} \sum_{j=1}^r \lambda_j \sigma\left(\langle \boldsymbol{\theta}_j, \boldsymbol{x} \rangle\right) \text{ and } \boldsymbol{x} \sim \mathcal{N}(0, \boldsymbol{I}_d),$$
 (1.1)

where $\sigma(z)=z^2-1$ is the 2nd Hermite polynomial; ϵ is zero-mean, independent sub-Gaussian noise; $\{\boldsymbol{\theta}_j\}_{j=1}^r\subset\mathbb{R}^d$ are unknown signal directions (index features) which we assume to be orthonormal; $\lambda_1>\lambda_2>\cdots>\lambda_r>0$ are their respective contributions; and $\boldsymbol{\Lambda}=\mathrm{diag}(\lambda_1,\cdots,\lambda_r)$ collects the second-layer coefficients. Our goal is to learn this target network using a "student" two-layer neural network with quadratic activation and r_s neurons, trained via a gradient-based optimization algorithm. This setting encompasses several well-known problems:

- Phase retrieval (r=1). The learning of one quadratic neuron has been studied extensively [Fie82, CC15, TV23]. The quadratic σ has information exponent k=2 (defined as the index of the lowest non-zero Hermite coefficient [DH18, BAGJ21]). This entails that randomly initialized parameters are close to a saddle point in high dimensions; hence the SGD dynamics exhibit a plateau ("search" phase) of length $\log d$ before the loss decreases sharply ("descent" phase).
- Multi-spike PCA $(r = \Theta_d(1))$. The target (1.1) is a subclass of Gaussian multi-index models, for which various algorithms have been proposed for the finite-rank case $r_s = \Theta_d(1)$ [CM20, DLS22, BBPV23]. The setting also closely relates to the multi-spike PCA problem, for which online SGD [AGP24] and other streaming algorithms has been studied [OK85, JJK⁺16, AZL17].
- Linear-width quadratic NN $(r \approx d)$. The regime where the teacher width r grows proportionally with dimensionality d has also been studied, typically in the well-conditioned setting (e.g., identical λ_j 's). Recent works characterized the landscape [SJL18, DL18, VBB19, GKZ19, GMMM19], optimization dynamics [MVEZ20, MBB23], and statistical efficiency [MTM $^+$ 24, ETZK25, BCN $^+$ 25].



- (a) Additive model hypothesis for scaling laws.
- (b) SGD risk curves for quadratic NN.

Figure 1: (a) Illustration of the additive model hypothesis, i.e., sum of emergent learning curves at different timescales yields a power law cumulative loss. (b) Population loss vs. compute for two-layer quadratic NNs trained with online SGD with batch size d on MSE loss. We set d = 3200, and for the teacher r = 2400, $\alpha = 1$.

In this work we focus on the "extensive-rank" regime where $r \times d^{\beta}$ for $\beta \in (0,1)$ and $r_s \times d^{\gamma}$ for $\gamma \in [0,1)$, and place a power-law assumption on the second-layer coefficients: $\lambda_j \times j^{-\alpha}$ for $\alpha \geq 0$. Our setting is motivated by the following lines of research.

Neural scaling laws & emergence. Recent empirical studies on large language models (LLMs) reveal that increasing the model or training data size often results in a predictable, power-law decrease in the loss known as *neural scaling laws* [HNA⁺17, KMH⁺20, HBM⁺22]. While such scaling of generalization error has been derived for sketched linear models [MRS22, BAP24, PPXP24, LWK⁺24, DLM24], these analyses assume random projection with no *feature learning*, and hence cannot capture the NN's ability to learn useful features [GDDM14, DCLT18] that adapt to the underlying data structure. We aim to investigate a setting where the training of a nonlinear NN beyond the "lazy" regime exhibits a nontrivial scaling law.

Feature learning in neural networks is often studied theoretically through the learning of *multi-index models*, where the target function depends on a small number of latent directions (see [BH25] and references therein). For these low-dimensional targets, it is known that the training dynamics typically exhibit *emergent* (or staircase-like) behavior — long plateaus followed by sharp drops in loss [BAGJ22, AAM23]. To reconcile this emergent loss curve with smooth power-law decay, recent works hypothesized that the pretraining objective can be decomposed into a sum of losses on individual tasks [MLGT24, NFLL24], the learning of each exhibits a sharp transition, and the superposition of numerous emergent risk curves at different timescales yields a power-law scaling of the cumulative loss (see Figure 1(a)). In this context, the two-layer network (1.1) can be viewed as a sum of single-index phase retrieval tasks, where the length of each $\sim \log d$ plateau in the risk trajectory can be modulated by the second-layer coefficient λ_j . This motivates the following question:

Q1: Does gradient-based training of a two-layer quadratic network yield power-law loss scaling, when the target function is an additive model with varying second-layer coefficients $\{\lambda_i\}_{i=1}^r$?

In Figure 1(b) we empirically observe the affirmative: when the target function has smoothly decaying second-layer weights, online SGD training yields a power-law risk curve that resembles the scaling laws in [KMH⁺20, HBM⁺22]. The goal of this work is to rigorously establish such scaling laws.

Learning extensive-width neural networks. Prior works on multi-index models have shown that when $r = \Theta_d(1)$, gradient-based training succeeds with polynomial sample complexity depending on properties of the link function [AAM22, DLS22, BBSS22]. The "extensive-rank" regime where $r \approx d^{\beta}$ for $\beta > 0$ is relatively under-explored (except for the linear width regime $r \approx d$ [MBB23, MTM+24]); this setting is arguably closer to the practical neural network training (compared to the narrow-width setting), and also bears connections to several observations in the LLM literature such as *superposition* [EHO+22] and *skill localization* [DDH+21, WWZ+22, PSZA23], where the model simultaneously acquires a large number of "skills" during pretraining (see e.g., [OSSW24]).

The learning dynamics of (1.1) with divergingly many neurons is challenging to analyze primarily due to the fact that the effective dynamics may not be captured by a finite set of *summary statistics* [BAGJ22] (as in the finite-r case). Recent works [OSSW24, SBH24] addressed this challenge by assuming that the activation σ has information exponent $k \ge 3$, which allows the learning dynamics

Algorithm	Decay rate (λ_j)	Risk scaling law	Result
Gradient flow	$\alpha > 0.5$	$\bar{t}^{-\frac{2\alpha-1}{\alpha}} + r_s^{-(2\alpha-1)}$	- Theorem 1
Gradient now	$\alpha < 0.5$	$(1 - \bar{t}^{\frac{1-2\alpha}{\alpha}})_+ + (1 - (r_s/r)^{1-2\alpha})_+$	
Online SGD (Stiefel)	$\alpha > 0.5$	$(\eta \bar{t})^{-\frac{2\alpha-1}{\alpha}} + r_s^{-(2\alpha-1)}$	Theorem 2
Omnie GGD (Guerer)	$\alpha < 0.5$	$(1 - (\eta \bar{t})^{\frac{1-2\alpha}{\alpha}})_+ + (1 - (r_s/r)^{1-2\alpha})_+$	

Table 1: Scaling laws for learning quadratic neural network (1.1) using population gradient flow and its online SGD discretization. We omit constant factors in the risk scaling for ease of presentation.

- In $\alpha > 0.5$, for population gradient flow, $\bar{t} \sim t \cdot \log d$ is the rescaled time; for online SGD, $\bar{t} \sim t \cdot \log d$ where t is the number of gradient steps, which is equal to the sample size, and $\eta \sim 1/(d \text{ polylog}(d))$ is the step size.
- In $\alpha < 0.5$, for population gradient flow, $\bar{t} \sim t \cdot r \log d$ is the rescaled time; for online SGD, $\bar{t} \sim t \cdot r \log d$ where t is the number of gradient steps and $\eta \sim 1/(dr^{\alpha} \operatorname{polylog}(d))$ is the step size.

to decouple across feature directions. However, the case $k \leq 2$, which includes the quadratic activation studied in this work, remained open: existing analyses either assumed "isotropic" feature contributions $(\lambda_1 = \lambda_r)$ [RL24, SBH24], or established a computational complexity for SGD that scales with $d^{\Theta(\lambda_1/\lambda_r)}$ [LMZ20], which leads to pessimistic *exponential* dimension dependency in the power-law setting we consider. We therefore ask the following question.

Q2: Can we establish optimization and sample complexity of learning an extensive-width quadratic neural network (1.1) with anisotropic, power-decaying feature contributions?

Finally, although our problem setup does not directly encompass commonly used activation functions such as ReLU, for SGD on multi-index models it is known that the Hermite-2 component is the first harmonic capturing multi-dimensional structure in the target function [DLS22, DKL⁺23b]. Consequently, our quadratic setting represents the first nontrivial mode of feature learning in SGD dynamics. Consistent with this view, Figure 3 shows that the risk trajectory of ReLU networks closely follows our theoretical characterization of quadratic networks over a substantial portion of training.

1.1 Our Contributions

We analyze the risk trajectory of learning (1.1) with both gradient flow on the mean squared error (MSE) loss and its online SGD discretization on Stiefel manifold, covering the extensive-width and power-law settings. We derive scaling laws for feature recovery and population risk as a function of teacher and student network widths r_s , r, the decay exponent α , the optimization time, and the sample size (for the discretized dynamics). Our contributions are summarized as follow (see Table 1).

- 1. In Section 3, we analyze the population gradient flow and tightly characterize the loss decay with respect to time and the student width r_s . We show that signal directions are recovered sequentially, and the population MSE follows a smooth power law specified by the decay rate $\alpha > 0$.
- 2. In Section 4, we consider the online stochastic gradient descent (SGD) dynamics on the Stiefel manifold and derive scaling laws with respect to sample size. When specializing to the isotropic setting $\alpha = 0$, our sample complexity improves upon [RL24] in the extensive-width setting and matches the information theoretic limit (in terms of d, r dependence) up to polylogarithmic factors.

The following technical challenges in the extensive-width regime are central to our analysis:

- Coupled population dynamics. As $r, r_s \to \infty$, we must track infinitely many overlapping student and teacher neurons. [OSSW24, SBH24] assumed high information exponent k > 2, to decouple the dynamics into r independent single-index models, but such property does not hold in our quadratic case (k = 2). We address this by leveraging the closed-form solution of the quadratic problem [MBB23], which satisfies a Matrix Riccati ODE. A key ingredient in our analysis is its monotonicity with respect to its initialization, illustrated in Figures 4(a), which enables sharp risk bounds via comparisons to decoupled models.
- Operator norm discretization error. Prior works [BAGJ21, BBPV23, AGP24] focused on finite-r settings, where Frobenius norm control of the SGD noise was sufficient and natural: it allows bounding error direction-wise without incurring additional dimension dependence. However, in the extensive-width regime, such bounds become pessimistic and lead to suboptimal r-dependent rates. Hence we need to establish operator norm concentration around the population dynamics.

• Matrix-monotone comparison framework. To control discretization error in operator norm, we extend the monotonicity-based argument to discrete time and introduce a novel comparison-based discretization technique. Our approach constructs matrix-valued reference sequences corresponding to decoupled dynamics that tightly bound the discrete evolution from above and below. This yields sharp operator norm control even when the true trajectories are non-monotone (see Figure 4), as the analysis avoids relying on the trajectory itself by comparing against simpler bounding sequences.

1.2 Additional Related Works

Learning multi-index models with SGD. When r=1, the target is a *single-index model* with quadratic link function. The SGD learning of single-index models has been extensively studied in the feature learning literature [BAGJ21, BES+22, MHPG+22, BES+23, MHWSE23, MLHD23, MZD+23, BMZ24, DNGL24, GWB25]; while this model has d parameters to be estimated, the quadratic link (with information exponent k=2) incurs an additional $\log d$ factor in the complexity of online SGD. More generally, the setting where $r=\Theta_d(1)$ is covered by recent analyses of *multi-index models* [AAM22, AAM23, BBPV23, DKL+23a, CWPPS23, AGP24, VE24, MHWE24]; however, these learning guarantees for multi-index models typically yield superpolynomial complexity when the target function is rank-extensive. The sample complexity of gradient-based learning is also connected to statistical query lower bounds [DPVLB24, DTA+24, LOSW24, ADK+24].

Quadratic NNs and additive models. Prior theoretical works on learning two-layer neural network with quadratic activation function have studied the loss landscape [SJL18, DL18, VBB19, GKZ19, GMMM19] and the optimization dynamics [MVEZ20, AKLS23, MBB23, RL24]. While existing optimization and statistical guarantees may cover the extensive-width regime (see e.g., [DL18, MBB23, RL24]), to our knowledge, precise scaling laws have not been established in our extensive-rank and power-law setting. (1.1) is also an instance of the *additive model* [Sto85, HT87, Bac17] where the individual functions are given as (orthogonal) single-index models with *unknown* index features. For this model, [OSSW24, SBH24] established learning guarantees in the well-conditioned regime, under the assumption that the link function σ has information exponent k > 2.

2 Background and Problem Setting

2.1 Student-teacher Setting

Teacher Network. We consider the task of learning a teacher network with a quadratic (second-order Hermite) activation function written as

$$y = f_*(\boldsymbol{x}) + \epsilon \text{ with } f_*(\boldsymbol{x}) \coloneqq \frac{1}{\|\boldsymbol{\Lambda}\|_{\text{F}}} \sum_{j=1}^r \lambda_j (\langle \boldsymbol{\theta}_j, \boldsymbol{x} \rangle^2 - 1) \text{ and } \boldsymbol{x} \sim \mathcal{N}(0, \boldsymbol{I}_d),$$
 (2.1)

where $x \in \mathbb{R}^d$ is the input; ϵ is zero-mean, bounded-variance sub-Gaussian noise; r is the teacher network width; and $\{\theta_j\}_{j=1}^r \subset \mathbb{R}^d$ is an orthonormal set of unknown signal vectors. We collect these as columns of the matrix $\Theta \in \mathbb{R}^{d \times r}$. The contributions of these vectors are determined by the unknown second-layer coefficients $\lambda_1 > \lambda_2 > \dots > \lambda_r > 0$ with a power-law decay $\lambda_j \asymp j^{-\alpha}$ for $\alpha \geq 0$, and Λ is a diagonal matrix whose j-th diagonal entry is λ_j . The normalization in front of summation ensures $\mathbb{E}[y^2]$ is constant. We focus on the regime where $r \asymp d^\beta$ for $\beta \in (0,1)$.

Remark 1. The orthogonality of $\{\theta_j\}_{j=1}^r$ can be assumed without loss of generality. Specifically, consider (2.1) with arbitrary first-layer weights Θ and normalization $\mathbb{E}[y] = 0$, the output can be written as $y \propto \text{Tr}(\boldsymbol{x}\boldsymbol{x}^\top\boldsymbol{\Theta}\boldsymbol{\Lambda}\boldsymbol{\Theta}^\top) + cst$; hence we may redefine (λ_j, θ_j) via the spectral decomposition.

Student Network. We learn the target model with a quadratic student network defined as

$$\hat{y}(x, \mathbf{W}) = \frac{1}{\sqrt{r_s}} \sum_{i=1}^{r_s} \langle \mathbf{w}_i, \mathbf{x} \rangle^2 - \|\mathbf{w}_i\|_2^2,$$
 (2.2)

where r_s is the width of the student network, and $\{\boldsymbol{w}_j\}_{j=1}^{r_s} \subset \mathbb{R}^d$ denotes the set of trainable weights. We collect these weights as the columns of the matrix $\boldsymbol{W} \in \mathbb{R}^{d \times r_s}$, and omit the dependence on \boldsymbol{x} in $\hat{y}(\boldsymbol{x}, \boldsymbol{W})$ when clear from the context. Note that the norm subtraction ensures $\mathbb{E}_{\boldsymbol{x}}[\hat{y}(\boldsymbol{x}, \boldsymbol{W})] = 0$.

We may equivalently write the student network as $\hat{y}(\boldsymbol{x}, \boldsymbol{W}) = \frac{1}{\sqrt{r_s}} \sum_{j=1}^{r_s} \|\boldsymbol{w}_j\|_2^2 \cdot (\langle \bar{\boldsymbol{w}}_j, \boldsymbol{x} \rangle^2 - 1)$ where $\bar{\boldsymbol{w}}_j$ is unit-norm; since our student does not have trainable second-layer, the norm component $\|\boldsymbol{w}_j\|_2^2$ allows the model to adapt to the target second-layer λ_j ; this homogeneous parameterization has been studied in prior works [CB20, GRWZ21].

2.2 Training Objective

Training constitutes to minimizing the squared loss; we define the instantaneous loss on (x, y) as

$$\mathcal{L}(\boldsymbol{W}; (\boldsymbol{x}, y)) \coloneqq \frac{1}{16} (y - \hat{y}(\boldsymbol{x}, \boldsymbol{W}))^2,$$

where the prefactor is included for notational convenience in the gradient computation. We omit the dependence on (x, y) when clear from context. The population risk can be written as

$$R(\boldsymbol{W}) := \mathbb{E}_{(\boldsymbol{x},y)}[\mathcal{L}(\boldsymbol{W})] = \frac{1}{8} \left\| \frac{1}{\sqrt{r_s}} \boldsymbol{W} \boldsymbol{W}^{\top} - \frac{1}{\|\boldsymbol{\Lambda}\|_{F}} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \right\|_{F}^{2} + \frac{1}{16} \mathbb{E}[\epsilon^{2}].$$
 (2.3)

Alignment. Observe that the student network is invariant to right-multiplication of its weight matrix by an orthonormal matrix, i.e., $\hat{y}(x, W) = \hat{y}(x, WO)$ for any $O \in \mathbb{R}^{r_s \times r_s}$ with $O^{\top}O = I$. Consequently, any notion of alignment that depends on individual directions in W may not be informative. To capture directional learning in a way that respects this symmetry, we define alignment in terms of the subspace spanned by student weights. We formalize this using the polar decomposition:

$$W := UQ^{1/2}$$
, where $Q := W^{\top}W$ and $U^{\top}U = I_{r_s}$. (2.4)

Here, Q denotes the radial component of the student weights, while U is an orthonormal matrix that encodes their directional component. We quantify the alignment between the student network and the jth teacher feature by the squared norm of the projection of θ_i onto the column space of W:

$$Alignment(\boldsymbol{W}, \boldsymbol{\theta}_j) := \|\boldsymbol{U}^{\top} \boldsymbol{\theta}_j\|_2^2. \tag{2.5}$$

Alignment(W, θ_j) takes values in the interval [0, 1]; it is 0 if θ_j is orthogonal to W (no alignment), while it is 1 if θ_j is in the column space of W (perfect alignment)¹.

3 Continuous Dynamics: Population Gradient Flow

We first analyze the continuous-time population gradient flow dynamics for (2.3), given as

$$\partial_t \mathbf{W}_t = -\nabla R(\mathbf{W}_t), \text{ where } \mathbf{W}_0 \in \mathbb{R}^{d \times r_s}, \ \mathbf{W}_{0,ij} \sim_{iid} \mathcal{N}(0, 1/d),$$
 (GF)

and the population gradient reads

$$\nabla R(\boldsymbol{W}_t) = -\frac{1}{2\sqrt{r_s}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}} \left(\boldsymbol{\Theta}\boldsymbol{\Lambda}\boldsymbol{\Theta}^{\top} - \frac{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} \boldsymbol{W}_t \boldsymbol{W}_t^{\top}\right) \boldsymbol{W}_t.$$

For notational convenience, we write $\mathcal{R}(t) := R(\boldsymbol{W}_t)$ and $\mathcal{A}(t, \boldsymbol{\theta}_j) := \operatorname{Alignment}(\boldsymbol{W}_t, \boldsymbol{\theta}_j)$. The following theorem sharply characterizes the timescale for alignment and the limiting risk curve. For ease of exposition, we drop the prefactor $\frac{1}{8}$ in the population risk so that it starts at 1.

Theorem 1. Let $\lambda_j = j^{-\alpha}$ and $r \approx d^{\beta}$ for some $\alpha \geq 0$ and $\beta \in (0,1)$. Consider the regime

$$\begin{cases} \frac{r_s}{r} \to \varphi \in (0, \infty) & \text{and } d \ge \Omega_{\alpha, \beta, \varphi}(1), & \text{if } \alpha \in [0, 0.5), \\ r_s \times 1, & \text{and } d \ge \Omega_{\alpha, r_s}(1), & \text{if } \alpha > 0.5. \end{cases}$$
(3.1)

Define the effective student width and effective timescale as

$$r_{\mathrm{eff}} \coloneqq \begin{cases} \lfloor r_s (1 - \log^{-1/8} d) \wedge r \rfloor, & \textit{if } \alpha \in [0, 0.5) \\ r_s, & \textit{if } \alpha > 0.5. \end{cases} \quad \textit{and} \ \, \mathsf{T}_{\mathrm{eff}} \coloneqq \sqrt{r_s} \| \mathbf{\Lambda} \|_{\mathrm{F}} \log d / r_s.$$

Then, the population (GF) dynamics satisfy the following with probability $1 - o(1/d^2) - \Omega(1/r_s^2)$:

¹The definition in (2.5) may fail to converge to 1 when $\alpha = 0$ and $r_s < r$, due to rotational symmetry in the teacher network. In this case, a more suitable notion of alignment can be defined using the principal angles between the subspaces spanned by W and Θ , which provides a rotation-invariant characterization of directional overlap. Specifically, for $\alpha = 0$, we define Alignment(W, θ_j) as the jth largest eigenvalue of $\Theta^T UU^T \Theta$.

1. Alignment: For $j \leq r_{\text{eff}}$ and t > 0 satisfying $t \asymp r^{\alpha}$ when $\alpha \in [0, 0.5)$ and $t \asymp 1$ when $\alpha > 0.5$, $\mathcal{A}(t\mathsf{T}_{\text{eff}}, \boldsymbol{\theta}_j) = \mathbb{1}\{t \geq \frac{1}{\lambda_i}\} + o_d(1)$. (3.2)

2. **Risk curve:** Under the same time scaling,

$$\mathcal{R}(t\mathsf{T}_{\text{eff}}) = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\text{F}}^2} \sum_{i=1}^{r_{\text{eff}}} \lambda_j^2 \mathbb{1}\{t \ge \frac{1}{\lambda_j}\} + o_d(1). \tag{3.3}$$

Remark 2. We make the following remarks about our result in Theorem 1:

- The spectral decay rate α determines both the choice of student width r_s and the timescale needed for learning in Theorem 1. Specifically, when $\alpha > 1/2$ (i.e., light-tailed regime), the target coefficients $\{\lambda_j\}_{j=1}^r$ are square-summable, making the teacher model effectively finite-dimensional. Therefore, a finite-width student suffices, and only finitely many directions need to be learned to achieve small loss, which results in a timescale of order $\log d$. In contrast, for the heavy-tailed regime $\alpha < 1/2$, we need to recover linear-in-r directions to achieve small population loss, which require both proportional width $r_s/r \to \varphi$ and a longer timescale $r \log d$. This difference in timescale will be made explicit in Corollary 1.
- Theorem 1 verifies the additive model hypothesis [MLGT24] for quadratic NNs in the feature learning regime; specifically, (3.2) identifies sharp transition time in alignment between student weights and the j-th teacher direction, and (3.3) suggests that the cumulative loss can be decomposed into individual emergent risk curves where the timescale is decided by signal strength λ_j .

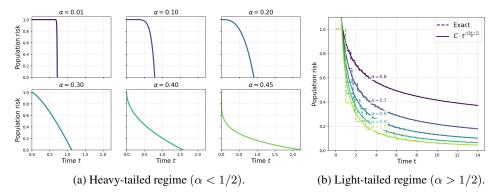


Figure 2: Illustration of the limiting risk trajectories and scaling behavior given in Corollary 1.

Neural scaling laws. As a corollary of Theorem 1, we obtain the following risk characterization. **Corollary 1.** *By Theorem 1, the asymptotic risk of* (GF) *is given as follows:*

• Heavy-tailed regime ($\alpha \in [0, 0.5)$): Almost surely, for all t > 0

$$\mathcal{R}(tr\log d) \xrightarrow{d\to\infty} \left(1 - Ct^{\frac{1-2\alpha}{\alpha}}\right)_+ \vee \left(1 - \varphi^{1-2\alpha}\right)_+.$$

• Light-tailed regime ($\alpha > 0.5$): With probability $1 - \Omega(1/r_s)$, for all t > 0, the risk $\mathcal{R}(t \log d)$ converges as $d \to \infty$ to a deterministic limit satisfying

$$\mathcal{R}(t\log d) \xrightarrow{d \to \infty} \Theta\left(t^{-\frac{2\alpha-1}{\alpha}} + r_s^{-(2\alpha-1)}\right).$$

Corollary 1 shows that, over appropriate timescales, the cumulative effect of these emergent transitions yields a smoothly decaying risk curve. Intuitively speaking, the power-law exponent arises from the Riemann integral approximation of the infinite sum (3.3) – see Appendix D.5 for details.

The asymptotic risk behavior in Corollary 1 is visualized in Figure 2 (see also Figure 1(b) for empirical simulation). The figure illustrates how the sharp, step-like emergent curve at $\alpha=0$ (as observed in earlier works on multi-index learning [BAGJ21, AAM23]) gradually transitions into a smooth curve as α increases. Notably, in the light-tailed regime $\alpha>1/2$, our risk curve resembles the neural scaling laws in [KMH+20, HBM+22] which takes the form of $\mathcal{R}\sim 1/(\mathrm{Data\ size})^a+1/(\mathrm{Model\ size})^b$, where the data size can be connected to optimization time under the one-pass discretization, which we analyze in the ensuing section.

Discrete Dynamics: Online Stochastic Gradient Descent

Now we analyze the finite-sample, discrete-time counterpart of the population dynamics (GF) and establish computational and statistical guarantees. We first discretize the directional component of the dynamics via online SGD with Stiefel constraint (see Proposition 2), and then introduce a finetuning step with negligible statistical and computational cost to fit the radial component; this mirrors the layer-wise training paradigm commonly used in theoretical analyses of gradient-based feature learning [AAM22, DLS22, BES+22, BEG+22]. The procedure is summarized in Algorithm 1.

Algorithm 1 Online Stochastic Gradient Descent (Stiefel)

1: **for**
$$t = 1, 2, \dots$$
 do

2:
$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathrm{St}} \mathcal{L}(\mathbf{W}_{t-1})$$

2:
$$\widetilde{\boldsymbol{W}}_t = \boldsymbol{W}_{t-1} - \eta \nabla_{\operatorname{St}} \mathcal{L}(\boldsymbol{W}_{t-1})$$

3: $\boldsymbol{W}_t = \widetilde{\boldsymbol{W}}_t \left(\widetilde{\boldsymbol{W}}_t^{\top} \widetilde{\boldsymbol{W}}_t\right)^{-1/2}$ \triangleright Feature learning

4: end for

5:
$$W_t^{\text{final}} = W_t \Omega_*$$
 where $\Omega_* = \underset{\Omega \in \mathbb{R}^{r_s \times r_s}}{\arg \min} \sum_{j=1}^{N_{\text{Ft}}} \mathcal{L} \left(W_t \Omega; (\boldsymbol{x}_{t+j}, y_{t+j}) \right)$ $ightharpoonup$ Fine-tuning

In the feature learning step, we update the first-layer weights W_t to recover the subspace spanned by the teacher directions. To this end, we use online SGD on Stiefel manifold [AGP24] with polar retraction. The Riemannian gradient on the Stiefel manifold is given by:

$$\nabla_{\operatorname{St}} \mathcal{L}(\boldsymbol{W}_{t-1}) \coloneqq \nabla \mathcal{L}(\boldsymbol{W}_{t-1}) - \frac{1}{2} \boldsymbol{W}_{t-1} \left(\boldsymbol{W}_{t-1}^{\top} \nabla \mathcal{L}(\boldsymbol{W}_{t-1}) + \nabla \mathcal{L}(\boldsymbol{W}_{t-1})^{\top} \boldsymbol{W}_{t-1} \right),$$

where the instantaneous loss is defined for sample (x_t, y_t) . Since the goal is to ensure subspace alignment (2.5), the overlap of individual student-teacher weights is not relevant during this phase.

After the feature learning phase, we perform a fine-tuning step to rotate W_t so that each w_i aligns with the corresponding teacher direction θ_j . This is achieved by solving an empirical risk minimization problem over N_{Ft} fresh samples. The optimal fine-tuning matrix Ω_* admits a closed-form solution that is also numerically easy to compute. Importantly, the computational and statistical complexity of this step scales only quadratically with the student width r_s , which is negligible compared to the cost of feature learning. The complexity analysis for this phase is provided in Appendix E.

Remark 3. Recall that the stage-wise training procedure is not required in our continuous analysis in Section 3. This is because we employ a Stiefel gradient similar to [BBPV23, AGP24] – which cannot fit the radial component - to simplify the discretization analysis. We conjecture that standard Euclidean discretization can achieve the same risk scaling; see Figure 1(b) for empirical evidence.

We define the population risk of the output of Algorithm 1, the alignment with a teacher direction θ_i , and the optimal risk achievable by a student neural network with width r_s respectively as

$$\mathcal{R}(t) \coloneqq R(\boldsymbol{W}_t^{ ext{final}}), \quad \mathcal{A}(t, \boldsymbol{\theta}_j) \coloneqq \operatorname{Alignment}(\boldsymbol{W}_t, \boldsymbol{\theta}_j), \quad \mathcal{R}_{\operatorname{opt}} \coloneqq \frac{1}{\|\mathbf{\Lambda}\|_{\mathbb{R}}^2} \sum_{j=(r_s \wedge r)+1}^r \lambda_j^2.$$

Intuitively, \mathcal{R}_{opt} is the risk achieved by exactly fitting the top $r_s \leq r$ components of the teacher model. Note that the alignment $A(t, \theta_i)$ depends only on the directional component of W_t ; thus, this quantity remains unchanged during fine-tuning. The following theorem characterizes the alignment and risk curve for the discrete-time Algorithm 1.

Theorem 2. Let the parameters $\{\lambda_j\}_{j=1}^r$, r, r_s , r_{eff} and T_{eff} , and the scaling regime (3.1) be as in Theorem 1. Suppose the student weights are initialized uniformly on the Stiefel manifold, and that the step size η and fine-tuning sample size N_{Ft} satisfy

$$\eta \approx \frac{1}{d} \begin{cases} \frac{1}{r^{\alpha} \log^{C_{\alpha}}(1+d/r_{s})}, & \alpha \in [0, 0.5) \\ \frac{1}{\log^{C_{\alpha}}d}, & \alpha > 0.5 \end{cases} \text{ and } N_{Ft} \approx r_{s}^{2} \log^{5} d,$$

for some constant $C_{\alpha} > 0$ depending only on α . Then with probability $1 - o_d(1/d^2) - \Omega(1/r_s^2)$,

1. Runtime and sample complexity: 1

$$T \ge \begin{cases} dr^{1+\alpha} \log^{C_{\alpha}+1} (1 + d/r_s), & \alpha \in [0, 0.5) \\ d \log^{C_{\alpha}+1} d, & \alpha > 0.5. \end{cases}$$
(4.1)

we have $\mathcal{R}(T) = \mathcal{R}_{\text{opt}} + o_d(1)$.

2. Alignment & Risk curve: For t > 0 satisfying $t \approx r^{\alpha}/\eta$ for $\alpha \in [0, 0.5)$ and $t \approx 1/\eta$ for $\alpha > 0.5$,

$$\bullet \ \mathcal{A}\big(t\mathsf{T}_{\mathrm{eff}},\boldsymbol{\theta}_j\big) = \mathbb{1}\{\eta t \geq \frac{1}{\lambda_j}\} + o_d(1). \quad \bullet \ \mathcal{R}\big(t\mathsf{T}_{\mathrm{eff}}\big) = 1 - \frac{1}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_{\mathrm{eff}}} \lambda_j^2 \ \mathbb{1}\{\eta t \geq \frac{1}{\lambda_j}\} + o_d(1).$$

Remark 4. We make the following remarks on the sample complexity.

- The bound in (4.1) implies a complexity of $n \approx T \simeq dr^{1+\alpha} \operatorname{polylog}(1+d/r_s)$ in the heavy-tailed case, and $T \simeq d \operatorname{polylog}(d)$ in the light-tailed case. Note that due to the one-pass nature of the algorithm, the runtime and sample complexity are identical (up to the negligible fine-tuning step).
- In the light-tailed regime ($\alpha > 1/2$), the required sample size $n \simeq d \operatorname{polylog}(d)$ is information theoretically optimal up to logarithmic factors. Note that kernel methods and neural networks in the lazy regime [JGH18, COB19] requires $n \gtrsim d^2$ samples to learn a quadratic target function; thus our sample complexity bound illustrates the benefit of feature learning.
- In the heavy-tailed regime ($\alpha < 1/2$), we obtain (nearly) information theoretically optimal sample complexity when $\alpha = 0$. For the intermediate regime $\alpha \in (0, 1/2)$, we conjecture that the optimal sample complexity is $T \simeq dr$, which implies our current bound is suboptimal by r^{α} .

Isotropic Setting ($\alpha=0$). In the isotropic case, where the goal is to estimate the r-dimensional subspace spanned by the teacher weights, the above theorem yields a sample and runtime complexity $n \asymp T \asymp dr \operatorname{polylog}(1+d/r_s)$. This interpolates between the $n \simeq d \operatorname{polylog}(d)$ rate for phase retrieval r=1 [TV23, BAGJ21], and $n \simeq d^2$ as $r \to d$, which matches the sample complexity in the linear-width regime [MTM+24, ETZK25]. Notably, our r-dependence improves upon the recent work of [RL24], which established a sufficient sample size of $n \gtrsim d \operatorname{poly}(r)$ for a similar quadratic setting. We expect our result to be optimal up to polylogarithmic factors due to the intrinsic dr-dimensional nature of the subspace recovery problem.

Scaling laws in discrete time. As indicated by the alignment and risk expressions in Theorem 2, a sufficiently small learning rate η ensures that running online SGD for t steps closely tracks the population gradient flow trajectory (GF) at time ηt , exhibiting the same scaling behavior. The following corollary formalizes the discrete-time counterpart of Corollary 1.

Corollary 2. We consider $\eta t \xrightarrow{d \to \infty} t_c > 0$. By Theorem 2, we have

• Heavy-tailed case ($\alpha \in [0, 0.5)$): Almost surely,

$$\mathcal{R}(tr\log d) \xrightarrow{d\to\infty} \left(1 - Ct_c^{\frac{1-2\alpha}{\alpha}}\right) \vee \left(1 - \varphi^{1-2\alpha}\right)_+.$$

• Light-tailed case ($\alpha > 0.5$): With probability $1 - \Omega(1/r_s)$, $\mathcal{R}(t \log d)$ has an asymptotic limit

$$\mathcal{R}(t \log d) \xrightarrow{d \to \infty} \Theta(t_c^{-\frac{2\alpha - 1}{\alpha}} + r_s^{-(2\alpha - 1)}).$$

5 Overview of Proof Techniques

To avoid notational confusion between discrete-time and continuous-time dynamics, we adopt the following convention throughout this section. Subscripts (e.g., W_t) denote discrete-time quantities, while parentheses (e.g., W(t)) denote continuous-time trajectories. Specifically, $\{W_t\}_{t\in\mathbb{N}}$ refers to the iterates of online SGD; $\{W(t)\}_{t\geq 0}$ denotes the continuous-time gradient flow governed by (GF). Since our proof strategy heavily relies on the matrix (Loewner) order for symmetric matrices, we introduce the following notations. For symmetric matrices $G_1, G_2 \in \mathbb{R}^{d \times d}$, we write $G_1 \prec G_2$ if $G_2 - G_1$ is positive definite. The reverse relations are denoted by $G_1 \succ G_2$ and $G_1 \succeq G_2$.

5.1 Proof Sketch of Theorem 1

We first observe that both the population risk $R(\boldsymbol{W}(t))$ and the alignment $\operatorname{Alignment}(\boldsymbol{W}(t), \boldsymbol{\theta}_j)$ depend on $\boldsymbol{W}(t)$ through two Gram matrices: the weight Gram matrix $\boldsymbol{G}_W(t) \coloneqq \boldsymbol{W}(t) \boldsymbol{W}(t)^{\top}$, and the alignment Gram matrix $\boldsymbol{G}_U(t) \coloneqq \boldsymbol{\Theta}^{\top} \boldsymbol{U}(t) \boldsymbol{U}(t)^{\top} \boldsymbol{\Theta}$, where $\{\boldsymbol{U}(t)\}_{t\geq 0}$ denotes the directional component of $\boldsymbol{W}(t)$, as defined in (2.4). The proof proceeds by analyzing the evolution of these matrices, each governed by an autonomous ODE; in particular, a matrix Riccati differential equation.

Proposition 1. The Gram matrices defined above satisfy the following matrix Riccati ODEs:

• Weight Gram matrix:
$$\partial_t G_W(t) = \frac{0.5}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}} \Big(\mathbf{\Theta} \mathbf{\Lambda} \mathbf{\Theta}^{\top} G_W(t) + G_W(t) \mathbf{\Theta} \mathbf{\Lambda} \mathbf{\Theta}^{\top} - \frac{2\|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} G_W^2(t) \Big).$$

• Alignment Gram matrix:
$$\partial_t G_U(t) = \frac{0.5}{\|\mathbf{\Lambda}\|_{\mathbb{F}}\sqrt{r_s}} (\mathbf{\Lambda} G_U(t) + G_U(t)\mathbf{\Lambda} - 2G_U(t)\mathbf{\Lambda} G_U(t)).$$

Both equations in Proposition 1 take the form of matrix Riccati ODEs [BLW91], whose structural properties play a central role in the proof. To illustrate the core idea, we focus on the alignment dynamics. For simplicity, we write $\vec{G}(t) \coloneqq G_U(t)$ and consider $\partial_t G(t) = \frac{0.5}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}} \left(\mathbf{\Lambda} G(t) + G(t)\mathbf{\Lambda} - 2G(t)\mathbf{\Lambda} G(t)\right)$.

$$\partial_t \boldsymbol{G}(t) = \frac{0.5}{\|\boldsymbol{\Lambda}\|_{\text{F}}\sqrt{r_s}} \left(\boldsymbol{\Lambda} \boldsymbol{G}(t) + \boldsymbol{G}(t)\boldsymbol{\Lambda} - 2\boldsymbol{G}(t)\boldsymbol{\Lambda} \boldsymbol{G}(t)\right).$$
 (5.1)

Note that Alignment $(W(t), \theta_j)$ corresponds to the j^{th} diagonal entry of G(t). To characterize its trajectory, we leverage the monotonicity of the matrix Riccati flow with respect to its initialization, i.e., if $G_0^+ \succeq G_0^-$, the corresponding solutions satisfy $G(t, G_0^+) \succeq G(t, G_0^-)$ for all $t \geq 0$, where $G(t, G_0)$ denotes the solution to (5.1) with initial condition G_0 . Our proof strategy builds on this monotonicity and proceeds as follows:

- 1. Diagonalization & decoupling. If G_0 is diagonal, the solution $\{G(t)\}_{t\geq 0}$ remains diagonal under (5.1), reducing the dynamics to independent scalar ODEs that govern each diagonal entry. Moreover, each scalar ODE admits a closed-form solution.
- 2. Asymptotic characterization. For general G_0 , we construct diagonal matrices $G_0^+ \succeq G_0 \succeq G_0^-$. By monotonicity, the corresponding trajectories upper and lower bound $\{G(t)\}_{t\geq 0}$. These bounding systems are diagonal and decoupled, and as $d \to \infty$, they converge to the same limit.

We apply this strategy in Appendix D.3 to derive the exact asymptotics stated in Theorem 1.

Remark 5. We note that while the Riccati flow is monotone with respect to its initialization, this does not imply that its solution is monotone in time. That is, the trajectory G(t) may not evolve monotonically in matrix order, even though a larger initialization yields a trajectory that remains above that of a smaller one for all $t \geq 0$. This distinction is illustrated in Figure 4.

5.2 Proof Sketch of Theorem 2

Extending Monotonicity Arguments to Discrete Dynamics

We begin by observing that online SGD on the Stiefel manifold approximates the directional component of the continuous-time gradient flow, with stochastic gradients arising from online sampling. The proposition below formalizes the idea that online SGD approximates the directional dynamics of the continuous gradient flow at the population level. For the statement, recall that U(t) denotes the directional component of the gradient flow solution W(t) from (GF), as defined in (2.4).

Proposition 2. Let $\widehat{W}(t)$ be the solution to the continuous-time gradient flow on the Stiefel manifold, initialized with U(0). Then for all t > 0, the column spaces of $\widetilde{W}(t)$ and U(t) coincide.

This result justifies studying the online SGD on Stiefel manifold via the directional dynamics of (GF). To this end, we introduce the discrete analog of G(t) above as $G_t = \Theta^T W_t W_t^T \Theta$. Extending the analysis to discrete time is non-trivial due to the loss of monotonicity in the Euler discretization of the

Riccati dynamics (5.1). In particular, the update
$$\underbrace{G_t = G_{t-1} + \frac{0.5\eta}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}} \left(\mathbf{\Lambda}G_{t-1} + G_{t-1}\mathbf{\Lambda} - 2G_{t-1}\mathbf{\Lambda}G_{t-1}\right)}_{\text{non-monotone dynamics}} + \text{(2nd-order terms and noise) (5.2)}$$

no longer preserves the matrix order structure crucial to the continuous-time argument.

To overcome this, we construct an auxiliary discrete system that approximates (5.2) up to second-order terms while preserving monotonicity. Specifically, we define the map

$$G(G_t, \eta) := G_t - \frac{\eta}{2} (2G_t - I_r) \Lambda (2G_t - I_r) (I_r + \eta \Lambda (2G_t - I_r))^{-1} + \eta \Lambda$$
 (5.3)

which matches (5.2) up to second-order terms. Indeed, expanding the inverse term gives

$$\boldsymbol{G}(\boldsymbol{G}_t, \eta) = \underbrace{\boldsymbol{G}_t - \frac{\eta}{2}(2\boldsymbol{G}_t - \boldsymbol{I}_r)\boldsymbol{\Lambda}(2\boldsymbol{G}_t - \boldsymbol{I}_r) + \eta\boldsymbol{\Lambda}}_{=\boldsymbol{G}_t + \eta(\boldsymbol{\Lambda}\boldsymbol{G}_t + \boldsymbol{G}_t\boldsymbol{\Lambda} - 2\boldsymbol{G}_t\boldsymbol{\Lambda}\boldsymbol{G}_t)} + \text{2nd-order terms.}$$

The key advantage of the iteration (5.3) is that it preserves matrix order:

Proposition 3. For
$$\eta > 0$$
, if $G_t^+ \succeq G_t^- \succeq 0$, we have $G(G_t^+, \eta) \succeq G(G_t^-, \eta)$.

We use this to bound the non-monotone dynamics (5.2) via monotone iterates. Roughly, we show that for small enough step size η , the following holds:

$$G(G_{t-1}, (1+\varepsilon)\eta) + \text{Noise} \succeq G_t \succeq G(G_{t-1}, (1-\varepsilon)\eta) + \text{Noise}$$

for some $\varepsilon = o_d(1)$, where we denote the effective learning rate in (5.2) with $\eta = \frac{\eta}{\sqrt{r_s} \|\mathbf{A}\|_F}$. We then follow the same bounding argument used in the continuous case by defining

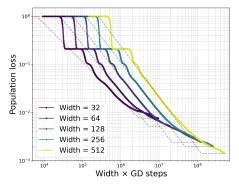
$$G_t^{\pm} = G(G_{t-1}^{\pm}, (1 \pm \varepsilon)\eta) + \text{Noise}, \text{ where } G_0^{+} \succeq G_0 \succeq G_0^{-} \succeq 0,$$

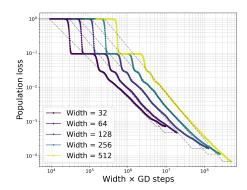
and show that $G_t^+ \succeq G_t \succeq G_t^-$ for all $t \in \mathbb{N}$. Finally, by choosing G_0^\pm to be diagonal, the bounding dynamics reduce to decoupled scalar recursions, which can be analyzed explicitly. This allows us to establish concentration of the original iterates $\{G_t\}_{t \in \mathbb{N}}$ around the bounding sequences, leading to operator-norm convergence of the discrete-time dynamics to their continuous-time counterparts.

6 Conclusion

In this work, we presented a comprehensive theoretical analysis of gradient-based learning in high-dimensional, extensive-width two-layer neural networks with quadratic activation. We established precise scaling laws that characterize both the population gradient flow and its empirical, discrete-time approximation. These results demonstrate how anisotropic signal strengths in the target function fundamentally shapes the convergence behavior and sample efficiency of gradient-based learning.

Beyond quadratic activations. An immediate direction for future research is to extend our analysis to more general activation functions. Link functions with higher information exponent is studied in a companion work [RNWL25], where the precise risk scaling is established by exploiting a decoupling structure that is unique to the information exponent k>2 setting. Importantly, many commonly-used activation functions (ReLU, GeLU, etc.) have information exponent k=1 and also contain a nonzero He_2 component. For such nonlinearities, we conjecture that SGD dynamics exhibits a multi-phase risk curve (analogous to the incremental learning phenomenon in [AAM23, BBPV23]), where the higher Hermite modes affects the learning dynamics after the low-order terms are learned. In Figure 3 we report the SGD risk curves for ReLU networks, in which we observe (i) an initial loss drop driven by the He_1 component (which finds a degenerate rank-1 subspace), followed by (ii) a power-law decay phase driven by the quadratic He_2 component where the empirical scaling exponent align closely with our theoretical predictions, and finally (iii) a slope change late in training likely due to higher Hermite terms (in Figure 5 we confirm that this "late" phase is absent if we remove these higher-order components). Understanding this complex multi-phase learning dynamics remains an interesting challenge for future work.





(a) $\alpha = 1$. Ideal exponent: -1; empirical: -1.01. (b) $\alpha = \frac{3}{2}$. Ideal exponent; $-\frac{4}{3}$, empirical:: -1.26.

Figure 3: Population loss vs. compute for two-layer ReLU network (power-law second-layer with exponent α) trained with population gradient descent. The student network adopts the 2-homogeneous parameterization as in (2.2). Observe that after the initial loss drop due to the He₁ component, the risk curves follow a power-law scaling where the exponent (dashed) nearly matches our theoretical prediction for the quadratic setting $\frac{1-2\alpha}{\alpha}$.

Acknowledgment

The authors thank Florent Krzakala, Jason D. Lee, and Lenka Zdeborová for discussion and feedback. The research of GBA was supported in part by NSF grant 2134216. MAE was partially supported by the NSERC Grant [2019-06167], the CIFAR AI Chairs program, and the CIFAR Catalyst grant. Part of this work was completed when NMV interned at the Flatiron Institute.

References

- [AAM22] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, pages 4782–4887. PMLR, 2022.
- [AAM23] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. SGD learning on neural networks: leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 2552–2623. PMLR, 2023.
- [ADK⁺24] Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan. Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index functions. *arXiv preprint arXiv:2405.15459*, 2024.
 - [AGP24] Gérard Ben Arous, Cédric Gerbelot, and Vanessa Piccolo. High-dimensional optimization for multi-spiked tensor pca. *arXiv preprint arXiv:2408.06401*, 2024.
- [AKLS23] Luca Arnaboldi, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Escaping mediocrity: how two-layer networks learn hard generalized linear models with sgd. arXiv preprint arXiv:2305.18502, 2023.
- [AZL17] Zeyuan Allen-Zhu and Yuanzhi Li. First efficient convergence for streaming k-pca: a global, gap-free, and near-optimal rate. In 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pages 487–492. IEEE, 2017.
- [Bac17] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- [BAGJ21] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *The Journal of Machine Learning Research*, 22(1):4788–4838, 2021.
- [BAGJ22] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. High-dimensional limit theorems for sgd: Effective dynamics and critical scaling. *Advances in Neural Information Processing Systems*, 35:25349–25362, 2022.
- [BAP24] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling laws. *arXiv preprint arXiv:2402.01092*, 2024.
- [BBPV23] Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning Gaussian multiindex models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- [BBSS22] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks. *Advances in Neural Information Processing Systems*, 35:9768–9783, 2022.
- [BCN⁺25] Jean Barbier, Francesco Camilli, Minh-Toan Nguyen, Mauro Pastore, and Rudy Skerk. Statistical mechanics of extensive-width bayesian neural networks near interpolation. *arXiv preprint arXiv:2505.24849*, 2025.
- [BEG⁺22] Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neural Information Processing Systems*, 35:21750–21764, 2022.

- [BES⁺22] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. *Advances in Neural Information Processing Systems*, 35:37932–37946, 2022.
- [BES⁺23] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
 - [BH25] Joan Bruna and Daniel Hsu. Survey on algorithms for multi-index models. *arXiv* preprint arXiv:2504.05426, 2025.
- [BLW91] Sergio Bittanti, Alan J. Laub, and Jan C. Willems. The riccati equation. 1991.
- [BMZ24] Raphaël Berthier, Andrea Montanari, and Kangjie Zhou. Learning time-scales in two-layers neural networks. Foundations of Computational Mathematics, pages 1–84, 2024.
 - [BR14] Davide Barilari and Luca Rizzi. Comparison theorems for conjugate points in subriemannian geometry. *ESAIM: Control, Optimisation and Calculus of Variations*, 22:439–472, 2014.
- [Bub14] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8:231–357, 2014.
- [CB20] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on learning theory*, pages 1305–1338. PMLR, 2020.
- [CC15] Yuxin Chen and Emmanuel Candes. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Advances in Neural Information Processing Systems*, 28, 2015.
- [CM20] Sitan Chen and Raghu Meka. Learning polynomials in few relevant dimensions. In *Conference on Learning Theory*, pages 1161–1227. PMLR, 2020.
- [COB19] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- [CWPPS23] Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi. Hitting the high-dimensional notes: An ode for sgd learning dynamics on glms and multi-index models. *arXiv* preprint arXiv:2308.08977, 2023.
 - [DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pretraining of deep bidirectional transformers for language understanding. *arXiv preprint* arXiv:1810.04805, 2018.
- [DDH⁺21] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*, 2021.
 - [DH18] Rishabh Dudeja and Daniel Hsu. Learning single-index models in gaussian space. In *Conference On Learning Theory*, pages 1887–1930. PMLR, 2018.
- [DKL⁺23a] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer neural networks learn, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023.
- [DKL⁺23b] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023.
 - [DL18] Simon Du and Jason Lee. On the power of over-parametrization in neural networks with quadratic activation. In *International conference on machine learning*, pages 1329–1338. PMLR, 2018.

- [DLM24] Leonardo Defilippis, Bruno Loureiro, and Theodor Misiakiewicz. Dimension-free deterministic equivalents for random feature regression. *arXiv preprint* arXiv:2405.15699, 2024.
- [DLS22] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR, 2022.
- [DNGL24] Alex Damian, Eshaan Nichani, Rong Ge, and Jason D Lee. Smoothing the landscape boosts the signal for sgd: Optimal sample complexity for learning single index models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [DPVLB24] Alex Damian, Loucas Pillaud-Vivien, Jason D Lee, and Joan Bruna. The computational complexity of learning gaussian single-index models. *arXiv preprint* arXiv:2403.05529, 2024.
 - [DTA⁺24] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer networks: Breaking the curse of information and leap exponents. *arXiv preprint arXiv*:2402.03220, 2024.
 - [Dur93] Richard Durrett. Probability: Theory and examples. 1993.
 - [EHO⁺22] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
 - [ETZK25] Vittorio Erba, Emanuele Troiani, Lenka Zdeborová, and Florent Krzakala. The nuclear route: Sharp asymptotics of erm in overparameterized quadratic networks. *arXiv* preprint arXiv:2505.17958, 2025.
 - [Fie82] James R Fienup. Phase retrieval algorithms: a comparison. *Applied optics*, 21(15):2758–2769, 1982.
- [GDDM14] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
 - [GKZ19] David Gamarnik, Eren C Kızıldağ, and Ilias Zadik. Stationary points of shallow neural networks with quadratic activation function. *arXiv preprint arXiv:1912.01599*, 2019.
- [GMMM19] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural network. *Advances in Neural Information Processing Systems*, 32, 2019.
- [GRWZ21] Rong Ge, Yunwei Ren, Xiang Wang, and Mo Zhou. Understanding deflation process in over-parametrized tensor decomposition. *Advances in Neural Information Processing Systems*, 34:1299–1311, 2021.
- [GWB25] Margalit Glasgow, Denny Wu, and Joan Bruna. Propagation of chaos in one-hidden-layer neural networks beyond logarithmic time. *arXiv preprint arXiv:2504.13110*, 2025.
- [HBM⁺22] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [HNA⁺17] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Patwary, Mostofa Ali, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.
 - [HT87] Trevor Hastie and Robert Tibshirani. Generalized additive models: some applications. *Journal of the American Statistical Association*, 82(398):371–386, 1987.

- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [JJK+16] Prateek Jain, Chi Jin, Sham M Kakade, Praneeth Netrapalli, and Aaron Sidford. Streaming pca: Matching matrix bernstein and near-optimal finite sample guarantees for oja's algorithm. In *Conference on learning theory*, pages 1147–1164. PMLR, 2016.
- [KMH⁺20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
 - [LM00] B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000.
 - [LMZ20] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer neural networks beyond ntk. In Conference on learning theory, pages 2613–2682. PMLR, 2020.
- [LOSW24] Jason D Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional polynomials with sgd near the information-theoretic limit. *arXiv* preprint arXiv:2406.01581, 2024.
- [LWK⁺24] Licong Lin, Jingfeng Wu, Sham M Kakade, Peter L Bartlett, and Jason D Lee. Scaling laws in linear regression: Compute, parameters, and data. *arXiv preprint* arXiv:2406.08466, 2024.
- [MBB23] Simon Martin, Francis Bach, and Giulio Biroli. On the impact of overparameterization on the training of a shallow neural network in high dimensions. *arXiv preprint arXiv:2311.03794*, 2023.
- [MHPG⁺22] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu. Neural networks efficiently learn low-dimensional representations with sgd. In *The Eleventh International Conference on Learning Representations*, 2022.
- [MHWE24] Alireza Mousavi-Hosseini, Denny Wu, and Murat A Erdogdu. Learning multi-index models with neural networks via mean-field langevin dynamics. *arXiv preprint arXiv:2408.07254*, 2024.
- [MHWSE23] Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A. Erdogdu. Gradient-based feature learning under structured data. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS 2023)*, 2023.
 - [MLGT24] Eric Michaud, Ziming Liu, Uzay Girit, and Max Tegmark. The quantization model of neural scaling. *Advances in Neural Information Processing Systems*, 36, 2024.
 - [MLHD23] Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks. *arXiv* preprint arXiv:2310.07891, 2023.
 - [MRS22] Alexander Maloney, Daniel A Roberts, and James Sully. A solvable model of neural scaling laws. *arXiv preprint arXiv:2210.16859*, 2022.
- [MTM⁺24] Antoine Maillard, Emanuele Troiani, Simon Martin, Florent Krzakala, and Lenka Zdeborová. Bayes-optimal learning of an extensive-width neural network from quadratically many samples. *arXiv preprint arXiv:2408.03733*, 2024.
- [MVEZ20] Sarao Stefano Mannelli, Eric Vanden-Eijnden, and Lenka Zdeborová. Optimization and generalization of shallow neural networks with quadratic activation functions. *Advances in Neural Information Processing Systems*, 33:13445–13455, 2020.

- [MZD⁺23] Arvind Mahankali, Haochen Zhang, Kefan Dong, Margalit Glasgow, and Tengyu Ma. Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks with polynomial width, samples, and time. *Advances in Neural Information Processing Systems*, 36, 2023.
- [NFLL24] Yoonsoo Nam, Nayara Fonseca, Seok Hyeong Lee, and Ard Louis. An exactly solvable model for emergence and scaling laws. *arXiv preprint arXiv:2404.17563*, 2024.
 - [OK85] Erkki Oja and Juha Karhunen. On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. *Journal of mathematical analysis and applications*, 106(1):69–84, 1985.
- [OSSW24] Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. Learning sum of diverse features: computational hardness and efficient gradient-based training for ridge combinations. In *Conference on Learning Theory*. PMLR, 2024.
- [PPXP24] Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+ 3 phases of compute-optimal neural scaling laws. *arXiv preprint arXiv:2405.15074*, 2024.
- [PSZA23] Abhishek Panigrahi, Nikunj Saunshi, Haoyu Zhao, and Sanjeev Arora. Task-specific skill localization in fine-tuned language models. *arXiv preprint arXiv:2302.06600*, 2023.
 - [RL24] Yunwei Ren and Jason D Lee. Learning orthogonal multi-index models: A fine-grained information exponent analysis. *arXiv preprint arXiv:2410.09678*, 2024.
- [RNWL25] Yunwei Ren, Eshaan Nichani, Denny Wu, and Jason D Lee. Emergence and scaling laws in sgd learning of shallow neural networks. arXiv preprint arXiv:2504.19983, 2025.
 - [SBH24] Berfin Simsek, Amire Bendjeddou, and Daniel Hsu. Learning gaussian multi-index models with gradient flow: Time complexity and directional convergence. *arXiv* preprint arXiv:2411.08798, 2024.
 - [SJL18] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
 - [Sto85] Charles J Stone. Additive regression and other nonparametric models. *The annals of Statistics*, 13(2):689–705, 1985.
 - [Tro10] Joel A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12:389–434, 2010.
 - [TV23] Yan Shuo Tan and Roman Vershynin. Online stochastic gradient descent with arbitrary initialization solves non-smooth, non-convex phase retrieval. *Journal of Machine Learning Research*, 24(58):1–47, 2023.
 - [VBB19] Luca Venturi, Afonso S Bandeira, and Joan Bruna. Spurious valleys in one-hidden-layer neural network optimization landscapes. *Journal of Machine Learning Research*, 20(133):1–34, 2019.
 - [VE24] Nuri Mert Vural and Murat A Erdogdu. Pruning is optimal for learning sparse features in high-dimensions. *arXiv preprint arXiv:2406.08658*, 2024.
 - [Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [WWZ⁺22] Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. Finding skill neurons in pre-trained transformer-based language models. *arXiv* preprint arXiv:2211.07349, 2022.

Contents

1	Introduction	1				
	1.1 Our Contributions	3				
	1.2 Additional Related Works	4				
2	Background and Problem Setting					
	2.1 Student-teacher Setting	4				
	2.2 Training Objective	5				
3	Continuous Dynamics: Population Gradient Flow	5				
4	Discrete Dynamics: Online Stochastic Gradient Descent	7				
5	Overview of Proof Techniques	8				
	5.1 Proof Sketch of Theorem 1	8				
	5.2 Proof Sketch of Theorem 2	9				
6	Conclusion	10				
A	Additional Figures	18				
В	Preliminaries for Proofs	19				
C	Background: Matrix Riccati Dynamical Systems 22					
	C.1 Continous-time Matrix Riccati ODE	21				
	C.2 Discrete-time Matrix Riccati Difference Equations	21				
D	Proofs for Main Results	24				
	D.1 Proof of Propositions 2 and 3	24				
	D.2 Decomposition of the population risk	24				
	D.3 Proof of Theorem 1	25				
	D.3.1 High-dimensional limit for the alignment	25				
	D.3.2 High-dimensional limit for the risk curve	26				
	D.4 Proof of Theorem 2	30				
	D.5 Proof of Corollary 1 and Corollary 2					
E	Details of the Fine-tuning Step	33				
	E.1 Characterizing the Minimum	33				
	E.2 Computing the Minimum	34				
	E.2.1 Proof of Proposition 11	34				
F	Deferred Proofs for Online SGD	35				
	F.1 Preliminaries	35				
	F.2 Including second-order terms and monotone bounds	37				
	F.2.1 Heavy tailed case - $\alpha \in [0, 0.5)$	37				
	F.2.2 Light tailed case - $\alpha > 0.5$	38				

	F.3	Definitions and bounding systems					
		F.3.1 Proof of Proposition 14					
	F.4	Analysis of the bounding systems					
		F.4.1 Lower bounding system					
		F.4.2 Upper bounding system					
	F.5	Bounds for the second-order terms					
	F.6	Noise characterization					
	F.7	Stability near minima					
G	Auxi	iliary Statements 70					
	G.1	Matrix bounds					
		G.1.1 Additional bounds for continuous-time analysis					
		G.1.2 Additional bounds for discrete-time analysis					
	G.2	Some moment bounds and concentration inequalities					
	G.3	Miscellaneous					

A Additional Figures

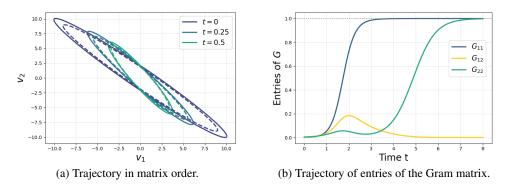


Figure 4: Solutions of the matrix Riccati ODE in (5.1) with $\lambda_1 = 2$, $\lambda_2 = 1$, $r_s = 2$. (a) To visualize the dynamics under matrix order, we plot the level sets of G(t) at times $t \in \{0, 0.25, 0.5\}$ for two initializations: G(0) (solid) and a scaled version 1.25 G(0) (dashed). The dashed ellipses remain enclosed within the solid ones at all times, illustrating monotonicity of the Riccati flow with respect to initialization. However, note that G(t) is not monotone in Loewner order over time, as seen from the lack of nesting among the solid ellipses. (b) Entry-wise evolution of G(t) under a random initialization with d = 1024. The diagonal entry G(t) exhibits non-monotonic behavior, illustrating that the solution trajectory G(t) need not be monotone in time.

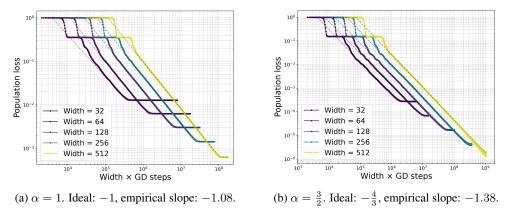


Figure 5: Population loss vs. compute for two-layer neural network with activation function $\sigma \propto {\rm He_1 + He_2}$, trained with population gradient descent. The student network adopts the 2-homogeneous parameterization as in (2.2). Observe that after the initial loss drop due to the ${\rm He_1}$ component, the risk curves exhibit a power-law scaling where the exponent (dashed lines) nearly matches our theoretical prediction for the quadratic setting $\frac{1-2\alpha}{\alpha}$; and unlike the ReLU setting (Figure 3), the loss immediately plateaus after the power-law phase.

Experiment Setting. In Figures 3 and 5, we plot the mean squared error loss for gradient descent with a constant step size on the population loss, using activations $\sigma = \text{ReLU}$ and $\sigma = \text{He}_1 + \text{He}_2$. The teacher model has orthogonal first-layer neurons and power-law decay in the second-layer coefficients with $\alpha \in \{1.0, 1.5\}$. Both teacher and student networks use the same activation function, which we normalize to have zero-mean an unit L^2 norm. The student network uses the 2-homogeneous parameterization:

$$\hat{y}(\boldsymbol{W}) = \frac{1}{\sqrt{r_s}} \sum_{i=1}^{r_s} \|\boldsymbol{w}_i\|_2^2 \cdot \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle) \text{ where } \sigma \in \{\frac{\text{ReLU} - 1/\sqrt{2\pi}}{0.5}, \frac{\text{He}_1 + \text{He}_2}{3}\}.$$

We set dimension d=5000, number of teacher neurons r=2400, student widths $r_s \in \{32,64,128,256,512\}$, and learning rate $\eta=0.5/\sqrt{r}$. To estimate the scaling exponents, we first identify the range of compute exhibiting a linear trend by visual inspection, and then fit the exponent via least squares. The dashed lines in the plots correspond to these fitted lines, and the reported empirical exponents represent the median values across different student widths.

B Preliminaries for Proofs

Proof organization. Section B introduces the notations and definitions used throughout the paper. In Section C, we provide a brief review of matrix Riccati ODEs and difference equations, along with the necessary supporting statements. The main results are proved in Section D. In Section E we discuss the fine-tuning phase for the discretized algorithm. Additional proofs related to online SGD and auxiliary lemmas are deferred to Sections F and G, respectively.

Notation and Definitions. We use $[n] \coloneqq \{1,2,\ldots,n\}$ to denote the first n natural numbers. The Euclidean inner product and norm are denoted by $\langle \cdot, \cdot \rangle$ and $\| \cdot \|_2$, respectively. For matrices, $\| \cdot \|_2$ and $\| \cdot \|_F$ denote the operator norm and Frobenius norm. The positive part is denoted by $(x)_+ \coloneqq \max\{x,0\}$. We write $f_d = o_d(1)$ if $f_d \to 0$ as $d \to \infty$, and $f_d \ll g_d$ if $f_d/g_d \to 0$. We use $O(\cdot)$ or $\Omega(\cdot)$ to suppress constants in upper and lower bounds respectively, and we use subscript to indicate parameter dependence, e.g., $O_{\alpha}(\cdot)$.

The symmetric part of a square matrix $M \in \mathbb{R}^{d \times d}$ is given by $\mathrm{Sym}(M) \coloneqq \frac{1}{2}(M + M^\top)$. For symmetric matrices $A, B \in \mathbb{R}^{d \times d}$, we write $A \prec B$ (or $A \preceq B$) if B - A is positive definite (or positive semidefinite). Moreover, if A and B are mutually diagonalizable, we write $AB^{-1} = \frac{A}{B}$.

We follow the convention that subscripts (e.g., W_t) refer to discrete-time quantities, and parentheses (e.g., W(t)) refer to continuous-time quantities. The overlap matrices of interest are defined as

$$\underbrace{\boldsymbol{G}_{W}(t) \coloneqq \boldsymbol{W}(t) \boldsymbol{W}(t)^{\top}}_{\text{Weight Gram matrix}}, \quad \underbrace{\boldsymbol{G}_{U}(t) \coloneqq \boldsymbol{\Theta}^{\top} \boldsymbol{U}(t) \boldsymbol{U}(t)^{\top} \boldsymbol{\Theta}}_{\text{Alignment Gram matrix}}, \quad \underbrace{\boldsymbol{G}_{t} \coloneqq \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta}}_{\text{Discrete alignment Gram matrix}}$$

Let $Z \in \mathbb{R}^{d \times r_s}$ be a Gaussian matrix with i.i.d. entries distributed as $\mathcal{N}(0, 1/d)$. We define $Z_{1:m} \in \mathbb{R}^{m \times r_s}$ as the submatrix formed by the first m rows:

$$oldsymbol{Z} = egin{bmatrix} oldsymbol{Z}_{1:m} \ oldsymbol{Z}_{ ext{rest}} \end{bmatrix}.$$

Without loss of generality, we assume the teacher directions coincide with the standard basis vectors, i.e., $\theta_j = e_j$. With this, the initialization satisfies:

$$G_W(0) = ZZ^{\top}, \qquad G_U(0) = G_0 = Z_{1:r}(Z^{\top}Z)^{-1}Z_{1:r}^{\top}.$$
 (B.1)

We start with characterizing "good events" for initial matrices given by the following lemma:

Lemma 1.

Both cases ($\alpha \in [0, 0.5) \cup (0.5, \infty)$). For $d \geq \Omega(1)$, the following holds:

(E.1)
$$\frac{1}{1.05} \le \lambda_{\min}(\boldsymbol{Z}^{\top}\boldsymbol{Z}) \le \lambda_{\max}(\boldsymbol{Z}^{\top}\boldsymbol{Z}) \le 1.05.$$

(E.2)
$$1.05 \mathbf{Z}_{1:r} \mathbf{Z}_{1:r}^{\top} \succeq \mathbf{G}_{U}(0) = \mathbf{G}_{0} \succeq \frac{1}{1.05} \mathbf{Z}_{1:r} \mathbf{Z}_{1:r}^{\top}$$

Heavy-tailed case ($\alpha \in [0, 0.5)$). For $d \ge \Omega_{\varphi}(1)$, the following holds:

(H.1) For
$$m \leq r_s (1 - \log^{-1/2} d) \wedge r$$
 uniformly, we consider $\lambda_{\min}(\mathbf{Z}_{1:m} \mathbf{Z}_{1:m}^{\top}) \geq \frac{r_s}{5d} (1 - \frac{m}{r_s})^2$.

(H.2) For all
$$m \leq r_s (1 - \log^{-1/2} d) \wedge r$$
 uniformly, $\lambda_m(\mathbf{Z}_{1:r} \mathbf{Z}_{1:r}^{\top}) \geq \frac{r_s}{5d} (1 - \frac{m}{r_s})^2$.

$$(H.3) \ \lambda_{\max}(\mathbf{Z}_{1:r}\mathbf{Z}_{1:r}^{\top}) \leq \frac{2r_s}{d} (1 + \frac{1}{\sqrt{\varphi}})^2.$$

Light-tailed case ($\alpha \in [0, 0.5)$). For $d \ge \Omega(1)$, the following holds:

(L.1)
$$\frac{1}{r_s^5 d} \leq \lambda_{\min}(\boldsymbol{Z}_{1:r_s} \boldsymbol{Z}_{1:r_s}^{\top})$$

(L.2) For
$$m \in \{1, 2, \cdots, 5r_s, \lceil \log^{2.5} d \rceil, \lceil \log^6 d \rceil, r \}$$
 uniformly, $\lambda_{\max}(\boldsymbol{Z}_{1:m} \boldsymbol{Z}_{1:m}^\top) \leq \frac{5(r_s \vee m)}{d}$

We define

$$\mathcal{G}_{\mathit{init}} \equiv \begin{cases} (E.1) \cap (E.2) \cap (H.1) \cap (H.2) \cap (H.3), & \alpha \in [0, 0.5) \\ (E.1) \cap (E.2) \cap (L.1) \cap (L.2), & \alpha \in > 0.5. \end{cases}$$

We have

$$\mathbb{P}[\mathcal{G}_{init}] \ge \begin{cases} 1 - 3r_s \exp\left(\frac{-r_s}{2\log^2 d}\right), & \alpha \in [0, 0.5) \\ 1 - \Omega(1/r_s^2), & \alpha > 0.5. \end{cases}$$

Proof. We will use the following:

(S.1) By [Ver10, Corollary 5.35], for $m \leq r_s$ and $\sqrt{r_s} - \sqrt{m} \geq t > 0$

$$\mathbb{P}\left[\lambda_{\min}\left(\boldsymbol{Z}_{1:m}\boldsymbol{Z}_{1:m}^{\top}\right) \geq \frac{r_s}{d}\left(1 - \sqrt{\frac{m}{r_s}} - \frac{t}{\sqrt{r_s}}\right)^2\right] \geq 1 - 2e^{-t^2},$$

and for $m \geq r_s$ and $\sqrt{m} - \sqrt{r_s} \geq t > 0$

$$\mathbb{P}\left[\lambda_{\min}\left(\boldsymbol{Z}_{1:m}^{\top}\boldsymbol{Z}_{1:m}\right) \geq \frac{m}{d}\left(1 - \sqrt{\frac{r_s}{m}} - \frac{t}{\sqrt{m}}\right)^2\right] \geq 1 - 2e^{-t^2}.$$

(S.2) By [Ver10, Corollary 5.35], for any fixed m

$$\mathbb{P}\left[\frac{m}{d}\left(1+\sqrt{\frac{r_s}{m}}+\frac{t}{\sqrt{m}}\right)^2 \geq \lambda_{\max}\left(\boldsymbol{Z}_{1:m}^{\top}\boldsymbol{Z}_{1:m}\right)\right] \geq 1-2e^{-t^2}.$$

(S.3) By [Ver10, Theorem 5.38], there exists C, c > 0 such that

$$\mathbb{P}\left[\lambda_{\min}\left(\boldsymbol{Z}_{1:r_s}\boldsymbol{Z}_{1:r_s}^{\top}\right) \geq \frac{\varepsilon^2}{4dr_s}\right] \geq 1 - C\varepsilon - e^{-cr_s}.$$

For the heavy tailed case, we consider d is large enough to guarantee $|\frac{r_s}{r} - \varphi| \leq \frac{\varphi}{2}$. We have

- By using (S.1) and (S.2) with m=d, $t=\sqrt{\frac{d}{\log d}}$, we can show that $\mathbb{P}[(E.1)] \geq 1-e^{\frac{-d}{\log d}}$ for $d \geq \Omega(1)$.
- By (B.1) and (E.1), (E.2) follows.
- For (H.1), by using (S.1) with $t = \frac{\sqrt{r_s} \sqrt{m}}{\sqrt{\log d}} \ge \sqrt{\frac{r_s}{2\log^2 d}}$, we have with probability $1 2r_s \exp\left(\frac{-r_s}{2\log^2 d}\right)$, for $m \le r_s (1 \log^{-1/2} d) \wedge r$ uniformly:

$$\lambda_{\min}\left(\boldsymbol{Z}_{1:m}\boldsymbol{Z}_{1:m}^{\top}\right) > \frac{r_s}{d}\left(1 - \frac{1}{\sqrt{\log d}}\right)^2 \left(1 - \sqrt{\frac{m}{r_s}}\right)^2 \ge \frac{r_s}{5d}\left(1 - \frac{m}{r_s}\right)^2.$$

Therefore, $\mathbb{P}[(\underline{H}.1)] \ge 1 - 2r_s \exp\left(\frac{-r_s}{2\log d}\right)$.

- By Cauchy's eigenvalue interlacing theorem, $\lambda_m(\boldsymbol{Z}_{1:r}\boldsymbol{Z}_{1:r}^\top) \geq \lambda_{\min}\left(\boldsymbol{Z}_{1:m}\boldsymbol{Z}_{1:m}^\top\right)$. Therefore, by (H.1), (H.2) follows.
- For (H.3), by using m = r and $t = 0.4\sqrt{r_s}$ in (S.2), we have $\mathbb{P}[(H.3)] \ge 1 2e^{-0.16r_s}$.
- For (L.1), by using (S.3) with $\varepsilon = \frac{2}{r_s^2}$, we have $\mathbb{P}[(L.1)] \ge 1 \Omega(1/r_s^2)$.
- For (L.2), by using (S.2) with $t = 0.4\sqrt{r_s}$, we have with probability $1 (10r_s + 6)e^{-0.16r_s}$ for $m \in [5r_s] \cup \{\lceil \log^{2.5} d \rceil, \lceil \log^6 d \rceil, r\}$ uniformly:

$$\lambda_{\max}(\boldsymbol{Z}_{1:m}\boldsymbol{Z}_{1:m}^{\top}) \leq \frac{r_s}{d} \left(1.4 + \sqrt{\frac{m}{r_s}}\right)^2 \leq \frac{5(r_s \vee m)}{d}.$$

By union bound, we have the result.

C Background: Matrix Riccati Dynamical Systems

We begin by reviewing Riccati dynamical systems in both continuous and discrete time, establishing the necessary background for the arguments that follow. For the following, we define

$$\mathbf{\Lambda}_{\mathrm{e}} \coloneqq \begin{bmatrix} \mathbf{\Lambda} & 0 \\ 0 & 0 \end{bmatrix}$$
 and $\widetilde{\mathbf{\Lambda}} \coloneqq \frac{\sqrt{r_s}}{\|\mathbf{\Lambda}\|_{\mathrm{F}}} \mathbf{\Lambda}_{\mathrm{e}}$.

For notational convenience, we adapt the abuse of notation:

$$\frac{\boldsymbol{\Lambda}_{\mathrm{e}}}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} = \lim_{\varepsilon \to 0} \frac{(\boldsymbol{\Lambda}_{\mathrm{e}} + \varepsilon \boldsymbol{I}_{d})}{\boldsymbol{I}_{d} - \exp(-t(\boldsymbol{\Lambda}_{\mathrm{e}} + \varepsilon \boldsymbol{I}_{d}))} = \begin{bmatrix} \frac{\boldsymbol{\Lambda}}{\boldsymbol{I}_{r} - \exp(-t\boldsymbol{\Lambda})} & 0\\ 0 & \frac{1}{t}\boldsymbol{I}_{d-r} \end{bmatrix}.$$

C.1 Continous-time Matrix Riccati ODE

In this paper, we study continuous-time matrix Riccati differential equations of the following form:

• Weight Gram matrix: For $T_W=r_s$

$$\partial_t \mathbf{G}_W(t) = \frac{0.5}{T_W} \left(\widetilde{\mathbf{\Lambda}} \mathbf{G}_W(t) + \mathbf{G}_W(t) \widetilde{\mathbf{\Lambda}} - 2\mathbf{G}_W^2(t) \right). \tag{C.1}$$

• Alignment Gram matrix: For $T_U = \|\mathbf{\Lambda}\|_{\mathrm{F}} \sqrt{r_s}$,

$$\partial_t \mathbf{G}_U(t) = \frac{0.5}{T_U} \Big(\mathbf{\Lambda} \mathbf{G}_U(t) + \mathbf{G}_U(t) \mathbf{\Lambda} - 2\mathbf{G}_U(t) \mathbf{\Lambda} \mathbf{G}_U(t) \Big). \tag{C.2}$$

For $\alpha = 0$, we assume that the ODEs are expressed in the eigenbasis of $G_W(0)$ or $G_U(0)$, ensuring that the trajectories remain diagonal. The solutions of these ODEs are characterized in the following statement:

Lemma 2. (C.1) and (C.2) admit the following solutions:

$$G_W(t) = \frac{\widetilde{\Lambda}}{I_d - \exp(-t\widetilde{\Lambda}/T_W)} - \frac{\widetilde{\Lambda} \exp(-0.5t\widetilde{\Lambda}/T_W)}{I_d - \exp(-t\widetilde{\Lambda}/T_W)} \left(G_W(0) + \frac{\widetilde{\Lambda} \exp(-t\widetilde{\Lambda}/T_W)}{I_d - \exp(-t\widetilde{\Lambda}/T_W)}\right)^{-1} \frac{\widetilde{\Lambda} \exp(-0.5t\widetilde{\Lambda}/T_W)}{I_d - \exp(-t\widetilde{\Lambda}/T_W)}$$

$$G_U(t) = \frac{I_r}{I_r - \exp(-t\Lambda/T_U)}$$
$$-\frac{\exp(-0.5t\Lambda/T_U)}{I_r - \exp(-t\Lambda/T_U)} \left(G_U(0) + \frac{\exp(-t\Lambda/T_U)}{I_r - \exp(-t\Lambda/T_U)}\right)^{-1} \frac{\exp(-0.5t\Lambda/T_U)}{I_r - \exp(-t\Lambda/T_U)}$$

Moreover, $(G_W(t))_{t\geq 0}$ and $(G_U(t))_{t\geq 0}$ are monotone with respect to $G_W(0) \succeq 0$ and $G_U(0) \succeq 0$ respectively.

Proof. One can check by direct differentiation that the given closed-form expressions satisfy the ODEs above. The uniqueness of the solutions follow the local Lipschitzness of the drifts. Monotonicity is a consequence of Proposition 25.

C.2 Discrete-time Matrix Riccati Difference Equations

In this section, we will study a particular discretization of Alignment Gram matrix ODE, given as

$$G_{t+1} = G_t - \frac{\eta}{2} (2G_t - I_r) \Lambda (2G_t - I_r) (I_r + \eta \Lambda (2G_t - I_r))^{-1} + \eta \Lambda.$$
 (C.3)

For convenience, we will make a change of variable and define $V_t := 2\Lambda^{\frac{1}{2}}G_t\Lambda^{\frac{1}{2}} - \Lambda$. We write (C.3) in terms of V_t as follows:

$$\boldsymbol{V}_{t+1} = \boldsymbol{V}_t - \eta \boldsymbol{V}_t^2 (\boldsymbol{I}_r + \eta \boldsymbol{V}_t)^{-1} + \eta \boldsymbol{\Lambda}^2.$$

We characterize the dynamics of $(V_t)_{t\in\mathbb{N}}$ as follows:

Lemma 3. We consider

$$\begin{bmatrix} \boldsymbol{X}_{t+1,1} \\ \boldsymbol{X}_{t+1,2} \end{bmatrix} = \begin{bmatrix} \boldsymbol{X}_{t,1} \\ \boldsymbol{X}_{t,2} \end{bmatrix} + \eta \boldsymbol{H} \begin{bmatrix} \boldsymbol{X}_{t,1} \\ \boldsymbol{X}_{t,2} \end{bmatrix} \quad \text{where} \quad \begin{bmatrix} \boldsymbol{X}_{0,1} \\ \boldsymbol{X}_{0,2} \end{bmatrix} = \begin{bmatrix} \boldsymbol{I}_r \\ \boldsymbol{V}_0 \end{bmatrix} \quad \text{and} \quad \boldsymbol{H} \coloneqq \begin{bmatrix} 0 & \boldsymbol{I}_r \\ \boldsymbol{\Lambda}^2 & \eta \boldsymbol{\Lambda}^2 \end{bmatrix}.$$

The following hold for all $n \in \mathbb{N}$:

(R.1) We have

$$\begin{bmatrix} \boldsymbol{A}_{t,11} & \boldsymbol{\Lambda}^{-1} \boldsymbol{A}_{t,12} \\ \boldsymbol{\Lambda} \boldsymbol{A}_{t,12} & \boldsymbol{A}_{t,22} \end{bmatrix} := (\boldsymbol{I}_{2r} + \eta \boldsymbol{H})^t.$$
 (C.4)

where $A_{t,11}$, $A_{t,12}$, $A_{t,22}$ are positive definite diagonal matrices.

(R.2)
$$A_{t,11} + \eta \Lambda A_{t,12} = A_{t,22}$$
 and $A_{t,22}A_{t,11} - A_{t,12}^2 = I_r$.

(R.3) For
$$\eta \leq 1$$
, we have $A_{t,11}A_{t,12}^{-1} \succ \left(I_r + \frac{\eta^2}{4}\Lambda^2\right)^{1/2} - \frac{\eta}{2}\Lambda$ and $A_{t,22}A_{t,12}^{-1} \succ \left(I_r + \frac{\eta^2}{4}\Lambda^2\right)^{1/2} + \frac{\eta}{2}\Lambda$.

(R.4) If $\|\eta \mathbf{\Lambda}\|_2 < 1$,

$$oldsymbol{A}_{t,22}oldsymbol{A}_{t,12}^{-1} \succ rac{(oldsymbol{I}_r + \etaoldsymbol{\Lambda})^t + (oldsymbol{I}_r - \etaoldsymbol{\Lambda})^t}{(oldsymbol{I}_r + \etaoldsymbol{\Lambda})^t - (oldsymbol{I}_r - \etaoldsymbol{\Lambda})^t} \succeq oldsymbol{A}_{t,11}oldsymbol{A}_{t,12}^{-1}.$$

Moreover, if $X_{t,1}$ and $X_{t+1,1}$ are invertible:

(R.7) For
$$V_{t+1} := X_{t+1,2} X_{t+1,1}^{-1}$$
, and $V_t := X_{2,t} X_{t,1}^{-1}$, we have

$$V_{t+1} = V_t - \eta V_t^2 \left(I_r + \eta V_t \right)^{-1} + \eta \Lambda^2.$$

Proof of Lemma 3. We have

$$I_{2r} + \eta H = \begin{bmatrix} I_r & \eta I_r \\ \eta \Lambda^2 & I_r + \eta^2 \Lambda^2 \end{bmatrix}.$$
 (C.5)

Let

$$(\boldsymbol{I}_{2r} + \eta \boldsymbol{H})^t =: egin{bmatrix} \tilde{\boldsymbol{A}}_{t,11} & \tilde{\boldsymbol{A}}_{t,12} \ \tilde{\boldsymbol{A}}_{21,t} & \tilde{\boldsymbol{A}}_{t,12} \end{bmatrix}.$$

Since each submatrix in (C.5) is diagonal positive definite, the matrices in (C.4) are also diagonal positive definite. To prove (R.1) and the first part of (R.2) we use proof by induction. We assume $\tilde{A}_{t,11} + \eta \tilde{A}_{21,t} = \tilde{A}_{t,12}$ and $\tilde{A}_{21,t} \tilde{A}_{12,t}^{-1} = \Lambda^2$. We have

$$\begin{split} \tilde{\boldsymbol{A}}_{12,t+1} = & \tilde{\boldsymbol{A}}_{t,12} + \eta \tilde{\boldsymbol{A}}_{t,12} \stackrel{(a)}{=} \tilde{\boldsymbol{A}}_{t,12} + \eta (\tilde{\boldsymbol{A}}_{t,11} + \eta \tilde{\boldsymbol{A}}_{21,t}) \stackrel{(b)}{=} \eta \tilde{\boldsymbol{A}}_{t,11} + (\boldsymbol{I}_r + \eta^2 \boldsymbol{\Lambda}^2) \tilde{\boldsymbol{A}}_{t,12} \\ &= \eta \tilde{\boldsymbol{A}}_{t,11} + \boldsymbol{\Lambda}^{-2} (\boldsymbol{I}_r + \eta^2 \boldsymbol{\Lambda}^2) \tilde{\boldsymbol{A}}_{21,t} \\ &= \boldsymbol{\Lambda}^{-2} \tilde{\boldsymbol{A}}_{21,t+1}. \end{split}$$

where (a) follows the first assumption, (b) and (c) follow the second assumption. Moreover,

$$\begin{split} \tilde{\pmb{A}}_{11,t+1} + \eta \tilde{\pmb{A}}_{21,t+1} &= \tilde{\pmb{A}}_{t,11} + \eta \tilde{\pmb{A}}_{21,t} + \eta^2 \pmb{\Lambda}^2 \tilde{\pmb{A}}_{t,11} + \eta (\pmb{I}_r + \eta^2 \pmb{\Lambda}^2) \tilde{\pmb{A}}_{21,t} \\ &= (\pmb{I}_r + \eta^2 \pmb{\Lambda}^2) (\tilde{\pmb{A}}_{t,11} + \eta \tilde{\pmb{A}}_{21,t}) + \eta \tilde{\pmb{A}}_{21,t} \\ &\stackrel{(d)}{=} (\pmb{I}_r + \eta^2 \pmb{\Lambda}^2) \tilde{\pmb{A}}_{t,12} + \eta \pmb{\Lambda}^2 \tilde{\pmb{A}}_{t,12} = \tilde{\pmb{A}}_{22,t+1}. \end{split}$$

where (d) follows the first and second assumptions. For the second part of (R.2) we again use proof by induction. We assume $A_{t,22}A_{t,11} - A_{t,12}^2 = I_r$. We have

$$\begin{split} \tilde{\boldsymbol{A}}_{11,t+1}\tilde{\boldsymbol{A}}_{22,t+1} - \tilde{\boldsymbol{A}}_{12,n+1}\tilde{\boldsymbol{A}}_{21,t+1} &= \left(\tilde{\boldsymbol{A}}_{t,11} + \eta \tilde{\boldsymbol{A}}_{21,t}\right) \left(\eta \boldsymbol{\Lambda}^2 \tilde{\boldsymbol{A}}_{t,12} + \left(\boldsymbol{I}_r + \eta^2 \boldsymbol{\Lambda}^2\right) \tilde{\boldsymbol{A}}_{t,12}\right) \\ &- \left(\eta \boldsymbol{\Lambda}^2 \tilde{\boldsymbol{A}}_{t,11} + \left(\boldsymbol{I}_r + \eta^2 \boldsymbol{\Lambda}^2\right) \tilde{\boldsymbol{A}}_{21,t}\right) \left(\tilde{\boldsymbol{A}}_{t,12} + \eta \tilde{\boldsymbol{A}}_{t,12}\right) \end{split}$$

$$= ilde{m{A}}_{t.11} ilde{m{A}}_{t.12}- ilde{m{A}}_{t.12} ilde{m{A}}_{21.t}=m{I}_r.$$

For (R.3), by using (R.2), we have

$$\mathbf{A}_{t,11} \left(\mathbf{A}_{t,11} + \eta \mathbf{\Lambda} \mathbf{A}_{t,12} \right) - \mathbf{A}_{t,12}^2 = \mathbf{I}_r \Rightarrow \left(\mathbf{A}_{t,11} \mathbf{A}_{t,12}^{-1} \right)^2 + \eta \mathbf{\Lambda} \left(\mathbf{A}_{t,11} \mathbf{A}_{t,12}^{-1} \right) - \mathbf{I}_r \succ 0$$

$$\Rightarrow \mathbf{A}_{t,11} \mathbf{A}_{t,12}^{-1} \succ \left(\mathbf{I}_r + \frac{\eta^2}{4} \mathbf{\Lambda} \right)^{1/2} - \frac{\eta}{2} \mathbf{\Lambda}.$$

The second part follows (R.2). For (R.4), we recall that

$$\mathbf{A}_{t+1,12} = \mathbf{A}_{t,12} + \eta \mathbf{\Lambda} \mathbf{A}_{t,22} = (\mathbf{I}_r + \eta^2 \mathbf{\Lambda}^2) \mathbf{A}_{t,12} + \eta \mathbf{\Lambda} \mathbf{A}_{t,11}$$
(C.6)

$$A_{t+1,22} = \eta \Lambda A_{t,12} + (I_r + \eta^2 \Lambda^2) A_{t,22} > \eta \Lambda A_{t,12} + A_{t,22}.$$
 (C.7)

We use proof by induction. Suppose the lower bound for $\frac{A_{t,22}}{A_{t,12}}$ holds. We have

$$\frac{\boldsymbol{A}_{t+1,22}}{\boldsymbol{A}_{t+1,12}} \stackrel{(e)}{\succ} \frac{\eta \boldsymbol{\Lambda} \boldsymbol{A}_{t,12} + \boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12} + \eta \boldsymbol{\Lambda} \boldsymbol{A}_{t,22}} \\
\stackrel{(f)}{\succ} \left(\eta \boldsymbol{\Lambda} + \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t} \right) \left(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda} \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t} \right)^{-1} \\
= \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^{t+1} + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^{t+1}}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^{t+1} - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^{t+1}}.$$

where (e) follows (C.7) and (f) follows the induction hypothesis with that $x \to \frac{x+\eta\lambda}{1+\eta\lambda x}$ is monotonic increasing for $\eta\lambda < 1$. For the upper bound, suppose the lower bound for $\frac{A_{t,11}}{A_{t,12}}$ holds. We have

$$\frac{\boldsymbol{A}_{t+1,11}}{\boldsymbol{A}_{t+1,12}} \stackrel{(g)}{\preceq} \frac{\boldsymbol{A}_{t,11} + \eta \boldsymbol{\Lambda} \boldsymbol{A}_{t,12}}{\boldsymbol{A}_{t,12} + \eta \boldsymbol{\Lambda} \boldsymbol{A}_{t,11}} \\
\stackrel{(h)}{\preceq} \left(\eta \boldsymbol{\Lambda} + \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t} \right) \left(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda} \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t} \right)^{-1} \\
= \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^{t+1} + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^{t+1}}{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^{t+1} - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^{t+1}}.$$

where (g) follows (C.6), and (h) follows the induction hypothesis.

Lastly if $X_{t,1}$ and $X_{t+1,1}$ are invertible,

$$\begin{aligned} \boldsymbol{X}_{t+1,2} \boldsymbol{X}_{t+1,1}^{-1} &= \left(\boldsymbol{X}_{t,2} \boldsymbol{X}_{t,1}^{-1} + \eta^2 \boldsymbol{\Lambda}^2 \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,1}^{-1} + \eta \boldsymbol{\Lambda}^2 \right) \left(\boldsymbol{I}_r + \eta \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,2}^{-1} \right)^{-1} \\ &= \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,2}^{-1} \left(\boldsymbol{I}_r + \eta \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,2}^{-1} \right)^{-1} + \eta \boldsymbol{\Lambda}^2 \\ &= \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,2}^{-1} - \eta \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,1}^{-1} \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,1}^{-1} \left(\boldsymbol{I}_r + \eta \boldsymbol{X}_{t,2} \boldsymbol{X}_{t,1}^{-1} \right)^{-1} + \eta \boldsymbol{\Lambda}^2 .\end{aligned}$$

Corollary 3. For $V_0=2\Lambda_2^{1\over 2}G_0\Lambda_2^{1\over 2}-\Lambda_1$, we define

$$V_{t+1} = V_t - \eta V_t^2 \left(I_r + \eta V_t \right)^{-1} + \eta \hat{\Lambda}^2.$$

If Λ_1 , Λ_2 and $\hat{\Lambda}$ are mutually diagonalizable, and $X_{1,t}$ is invertible for $t \leq t^* \in \mathbb{N}$, we have for $t \leq t^*$

$$G_t = \frac{\frac{\mathbf{\Lambda}_1}{\mathbf{\Lambda}_2} + \frac{\mathbf{A}_{t,22}}{\mathbf{A}_{t,12}} \frac{\hat{\mathbf{\Lambda}}}{\mathbf{\Lambda}_2}}{2} - \frac{1}{4} \frac{\mathbf{A}_{t,12}^{-1} \hat{\mathbf{\Lambda}}}{\mathbf{\Lambda}_2} \left(\frac{\frac{\hat{\mathbf{\Lambda}}_2}{\mathbf{\Lambda}_2} \frac{\mathbf{A}_{t,11}}{\mathbf{A}_{t,12}} - \frac{\mathbf{\Lambda}_1}{\mathbf{\Lambda}_2}}{2} + G_0 \right)^{-1} \frac{\hat{\mathbf{\Lambda}} \mathbf{A}_{t,12}^{-1}}{\mathbf{\Lambda}_2},$$

where A_{11t} , $A_{t,12}$, $A_{t,22}$ are defined with $\hat{\Lambda}$.

Proof. By using (C.4), we can write that

$$\begin{split} \boldsymbol{V}_t &= \left(\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12} + \boldsymbol{A}_{t,22} \boldsymbol{V}_0\right) \left(\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \boldsymbol{A}_{t,11} + \boldsymbol{V}_0\right)^{-1} \hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \\ &= \hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12} \Big(\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \boldsymbol{A}_{t,11} + \boldsymbol{V}_0\Big)^{-1} \hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \\ &+ \boldsymbol{A}_{t,22} \Big(\boldsymbol{I}_r - \boldsymbol{A}_{t,11} \boldsymbol{A}_{t,12}^{-1} \hat{\boldsymbol{\Lambda}} \left(\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \boldsymbol{A}_{t,11} + \boldsymbol{V}_0\right)^{-1} \Big) \hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \\ &= \boldsymbol{A}_{t,22} \boldsymbol{A}_{t,12}^{-1} \hat{\boldsymbol{\Lambda}} - \boldsymbol{A}_{t,12}^{-1} \hat{\boldsymbol{\Lambda}} \left(\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1} \boldsymbol{A}_{t,11} + \boldsymbol{V}_0\right)^{-1} \hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1}. \end{split}$$

Therefore,

$$m{G}_t = rac{rac{m{\Lambda}_1}{m{\Lambda}_2} + rac{m{A}_{t,12}}{m{A}_{t,12}}rac{\hat{m{\Lambda}}}{m{\Lambda}_2}}{2} - rac{1}{4}rac{m{A}_{t,12}^{-1}\hat{m{\Lambda}}}{m{\Lambda}_2} \left(rac{\hat{m{\Lambda}}_2}{m{\Lambda}_2}rac{m{A}_{t,11}}{m{A}_{t,12}} - rac{m{\Lambda}_1}{m{\Lambda}_2}}{2} + m{G}_0
ight)^{-1}rac{\hat{m{\Lambda}}m{A}_{t,12}^{-1}}{m{\Lambda}_2}.$$

Proposition 4. For some symmetric matrix S, we consider

$$V_1 = V_0 + \eta S - \eta V_0^2 (I_r + \eta V_0)^{-1}.$$
 (C.8)

If $V_0^+ \succeq V_0 \succ \frac{-1}{n} I_r$, we have $V_1^+ \succeq V_1$, where V_1^+ is the next iterate if we use V_0^+ in (C.8).

Proof. We have

$$oldsymbol{V}_1 = rac{1}{\eta} \Big(oldsymbol{I}_r - \left(oldsymbol{I}_r + \eta oldsymbol{V}_0
ight)^{-1} \Big) + \eta oldsymbol{S}.$$

The statement follows by Proposition 25.

D Proofs for Main Results

D.1 Proof of Propositions 2 and 3

For Proposition 2, we observe that

$$\widehat{\boldsymbol{G}}(t) \coloneqq \widehat{\boldsymbol{W}}(t) \widehat{\boldsymbol{W}}(t)^{\top} \ \text{ and } \ \widetilde{\boldsymbol{G}}(t) \coloneqq \boldsymbol{U}(t) \boldsymbol{U}(t)^{\top} = \boldsymbol{W}(t) (\boldsymbol{W}(t)^{\top} \boldsymbol{W}(t))^{-1} \boldsymbol{W}(t)^{\top}$$

have the exact same dynamics. Therefore the statement follows. Proposition 3 follows Proposition 25.

D.2 Decomposition of the population risk

The fine-tuning step relies on the following decomposition of the population risk:

Proposition 5. For any $\Omega \in \mathbb{R}^{r_s \times r_s}$, the population risk defined in (2.3) can be written as:

$$R(\boldsymbol{W}_{t}\boldsymbol{\Omega}) = \frac{1}{r_{s}} \left\| \boldsymbol{\Omega} \boldsymbol{\Omega}^{\top} - \frac{\sqrt{r_{s}}}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \right\|_{F}^{2} + \frac{1}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}} \left(\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} - \|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_{t} \boldsymbol{\Lambda}^{\frac{1}{2}} \|_{F}^{2} \right),$$

where $G_t = \Theta^\top W_t W_t^\top \Theta$ is the discrete alignment Gram matrix defined in the previous part.

We observe that both the second term and the matrix $W_t^\top \Theta \Lambda \Theta^\top W_t$ are independent of Ω . Hence, the fine-tuning step reduces to a least squares problem in the matrix $\Omega \Omega^\top$ in population, which is approximated via empirical risk minimization over a fresh batch of samples. By standard concentration arguments, a sample size of $N_{\rm Ft} \geq r_s^2 {\rm polylog} d$ suffices to ensure that the empirical minimizer approximates the population solution with high probability.

Proof. We begin by noting that W_t is an orthonormal matrix. Using this, we can express the population risk as:

$$R(\boldsymbol{W}_{t}\boldsymbol{\Omega}) = \left\| \frac{1}{\sqrt{r_{s}}} \boldsymbol{W}_{t} \boldsymbol{\Omega} \boldsymbol{\Omega}^{\top} \boldsymbol{W}_{t}^{\top} - \frac{1}{\|\boldsymbol{\Lambda}\|_{F}} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \right\|_{F}^{2}$$

$$\begin{split} &= \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^{2}} \left(\|\mathbf{\Lambda}\|_{\mathrm{F}}^{2} + \frac{\|\mathbf{\Lambda}\|_{\mathrm{F}}^{2}}{r_{s}} \|\mathbf{\Omega}\mathbf{\Omega}^{\top}\|_{F}^{2} - \frac{2\|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_{s}}} \mathrm{Tr}(\mathbf{\Omega}\mathbf{\Omega}^{\top} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t}) \pm \|\boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \|_{F}^{2} \right) \\ &= \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^{2}} \left(\|\mathbf{\Lambda}\|_{\mathrm{F}}^{2} - \|\boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \|_{F}^{2} + \left\| \frac{\|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_{s}}} \mathbf{\Omega} \mathbf{\Omega}^{\top} - \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \right\|_{F}^{2} \right) \end{split}$$

By observing that $\| \boldsymbol{W}^{ op} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{ op} \boldsymbol{W} \|_F^2 = \| \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_t \boldsymbol{\Lambda}^{\frac{1}{2}} \|_F^2$, we have the statement. \square

D.3 Proof of Theorem 1

We let

$$t_{\mathrm{sc}} \coloneqq t\sqrt{r_s}\|\mathbf{\Lambda}\|_{\mathrm{F}}, \quad \kappa_{\mathrm{eff}} \coloneqq \begin{cases} r^{\alpha}, & \alpha \in [0, 0.5) \\ 1, & \alpha > 0.5 \end{cases}, \qquad \mathsf{T}_{\mathrm{eff}} \coloneqq \kappa_{\mathrm{eff}}\sqrt{r_s}\|\mathbf{\Lambda}\|_{\mathrm{F}}\log d/r_s.$$

and

$$r_u := \begin{cases} r, & \alpha \in [0, 0.5) \\ \lceil \log^{2.5} d \rceil, & \alpha > 0.5 \end{cases} \qquad r_{u_{\star}} := \begin{cases} \lfloor r_s (1 - \log^{-1/s} d) \wedge r \rfloor, & \alpha \in [0, 0.5) \\ r_s, & \alpha > 0.5. \end{cases}$$

In the following part, we will establish the high-dimensional limit of the risk curve and the alignment.

D.3.1 High-dimensional limit for the alignment

By Lemma 2, we have

$$G_U(t_{\rm sc}) = \frac{I_r}{I_r - \exp(-t\Lambda)} - \frac{\exp(-0.5t\Lambda)}{I_r - \exp(-t\Lambda)} \left(G_U(0) + \frac{\exp(-t\Lambda)}{I_r - \exp(-t\Lambda)}\right)^{-1} \frac{\exp(-0.5t\Lambda)}{I_r - \exp(-t\Lambda)}$$

We define the block matrix forms

$$\boldsymbol{G}_{U}(t) \coloneqq \begin{bmatrix} \boldsymbol{G}_{U,11}(t) & \boldsymbol{G}_{U,12}(t) \\ \boldsymbol{G}_{U,12}^{-1}(t) & \boldsymbol{G}_{U,22}(t) \end{bmatrix}, \ \boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{\mathrm{eff}} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Lambda}_{22} \end{bmatrix}, \ \boldsymbol{\Lambda}_{\mathrm{e},11} \coloneqq \boldsymbol{\Lambda}_{\mathrm{eff}}, \ \boldsymbol{\Lambda}_{\mathrm{e},22} \coloneqq \begin{bmatrix} \boldsymbol{\Lambda}_{22} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix},$$

where $G_{U,11}(t), \Lambda_{\text{eff}} \in \mathbb{R}^{r_{u_{\star}} \times r_{u_{\star}}}$. The following statement characterizes the time-scales for the alignment terms.

Proposition 6. \mathcal{G}_{init} implies that $\mathcal{A}(t\mathsf{T}_{\mathrm{eff}},\boldsymbol{\theta}_j) = \mathbb{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1)$ for $t \neq \lim_{d \to \infty} \frac{1}{\lambda_j \kappa_{\mathrm{eff}}}$ and $j \leq r_{u_{\star}}$.

Proof. For $\alpha = 0$, since the trajectory stays diagonal and the diagonal entries are monotonically increasing, by using the events (E.2) and (H.2) with Lemma 10 we have the result.

In the following, we will prove the result for $\alpha > 0$. By using Proposition 22 with (L.2)

$$G_U(0) \leq \begin{bmatrix} 2.1 \mathbf{Z}_{1:r_{u_{\star}}} \mathbf{Z}_{1:r_{u_{\star}}}^{\top} & 0 \\ 0 & 2.1 \mathbf{Z}_2 \mathbf{Z}_2^{\top} \end{bmatrix},$$

where

$$2.1\lambda_{\max}(\mathbf{Z}_{1:r_{u_{\star}}}\mathbf{Z}_{1:r_{u_{\star}}}^{\top}) \leq \begin{cases} 5(1+\frac{1}{\sqrt{\varphi}})^{2}, & \alpha \in [0,0.5) \\ 15, & \alpha > 0.5. \end{cases}$$

Therefore.

$$\begin{split} \boldsymbol{G}_{U,11}(t_{\text{sc}}) & \leq \frac{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\text{eff}})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\text{eff}})} \\ & - \frac{\exp(-0.5t\boldsymbol{\Lambda}_{\text{eff}})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\text{eff}})} \left(\frac{O(r_s)}{d} \boldsymbol{I}_{r_{u_{\star}}} + \frac{\exp(-t\boldsymbol{\Lambda}_{\text{eff}})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\text{eff}})}\right)^{-1} \frac{\exp(-0.5t\boldsymbol{\Lambda}_{\text{eff}})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\text{eff}})} \end{split}$$

Therefore, by Proposition 36, for $j \leq r_{u_{\star}}$,

$$\mathcal{A}(t\mathsf{T}_{\mathrm{eff}},\pmb{\theta}_j) \leq \frac{1}{1 + \left(\frac{d}{r_s}\frac{1}{\log^3 d} - 1\right)\frac{d}{r_s}^{-t\kappa_{\mathrm{eff}}j^{-\alpha}}} = \mathbb{I}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1).$$

Moreover, for $t \leq (r_{u_{\star}} + 1)^{\alpha} \log \frac{d}{r_{\star}}$, by using the events (H.3) and (L.2), we have

$$\boldsymbol{Z}_2^{\top} \exp(t\boldsymbol{\Lambda}_{22}) \boldsymbol{Z}_2 \preceq \begin{cases} O_{\varphi}(1) \boldsymbol{I}_{r_s}, & \alpha \in (0, 0.5) \\ O(\log^{2.5} d) \boldsymbol{I}_{r_s}, & \alpha > 0.5. \end{cases}$$

Therefore, for $t \leq (r_{u_{\star}} + 1)^{\alpha} \log \frac{d}{r_s}$, we have $G_{U,11}(t_{sc}) \succeq \underline{G}(t)$ where

$$\underline{G}(t) \coloneqq \frac{I_{r_{u_{\star}}}}{I_{r_{u_{\star}}} - \exp(-t\mathbf{\Lambda}_{\text{eff}})} \\
- \frac{\exp(-0.5t\mathbf{\Lambda}_{\text{eff}})}{I_{r_{u_{\star}}} - \exp(-t\mathbf{\Lambda}_{\text{eff}})} \left(\frac{r_s}{d} \frac{O(1)}{\log^4 d} I_{r_{u_{\star}}} + \frac{\exp(-t\mathbf{\Lambda}_{\text{eff}})}{I_{r_{u_{\star}}} - \exp(-t\mathbf{\Lambda}_{\text{eff}})}\right)^{-1} \frac{\exp(-0.5t\mathbf{\Lambda}_{\text{eff}})}{I_{r_{u_{\star}}} - \exp(-t\mathbf{\Lambda}_{\text{eff}})}$$

which implies that for $t < (r_{u_{\star}} + 1)^{\alpha} \log \frac{d}{r_s}$

$$\mathcal{A}(t\mathsf{T}_{\mathrm{eff}},\boldsymbol{\theta}_{j}) \geq \frac{1}{1 + O(\log^{4}d)\frac{d}{r_{-}}}^{1 - t\kappa_{\mathrm{eff}}j^{-\alpha}} = \mathbb{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_{j}}\} + o_{d}(1).$$

To extend the lower bound for $t>(r_{u_\star}+1)^\alpha\log\frac{d}{r_s}$, let us define

$$t_0 \coloneqq (r_{u_\star} + 1)^\alpha \log \frac{d}{r_s}$$
 and $\boldsymbol{\Lambda}_{\text{eff}}^- \coloneqq \boldsymbol{\Lambda}_{\text{eff}} - (r_{u_\star} + 1)^{-\alpha} \boldsymbol{I}_{r_{u_\star}}$.

We have for $t > t_0$,

$$\partial_{t} \boldsymbol{G}_{U,11}(t) = \underbrace{\frac{0.5}{T_{U}} \Big(\boldsymbol{\Lambda}_{\text{eff}}^{-} \boldsymbol{G}_{U,11}(t) + \boldsymbol{G}_{U,11}(t) \boldsymbol{\Lambda}_{\text{eff}}^{-} - 2\boldsymbol{G}_{U,11}(t) \boldsymbol{\Lambda}_{\text{eff}}^{-} \boldsymbol{G}_{U,11}(t) \Big)}_{+ \underbrace{\frac{1}{T_{U}} \Big((r_{u} + 1)^{-\alpha} \boldsymbol{G}_{U,11}(t) (\boldsymbol{I}_{r_{u_{\star}}} - \boldsymbol{G}_{U,11}(t)) - \boldsymbol{G}_{U,12}(t) \boldsymbol{\Lambda}_{22} \boldsymbol{G}_{U,12}^{\top}(t) \Big)}_{\succeq \underbrace{\frac{(r_{u} + 1)^{-\alpha}}{T_{U}} \Big(\boldsymbol{G}_{U,11}(t) (\boldsymbol{I}_{r_{u_{\star}}} - \boldsymbol{G}_{U,11}(t)) - \boldsymbol{G}_{U,12}(t) \boldsymbol{G}_{U,12}^{\top}(t) \Big)}_{\succeq 0}}_{-}$$

Therefore, for $t > t_0$, by monotonicity and [BR14, Theorem 38], $G_{U,11}(t_{sc}) \succeq \underline{G}(t) \succeq \underline{G}(t_0)$, where

$$\begin{split} &\underline{\boldsymbol{G}}(t) = \frac{\boldsymbol{I}_{r_{u_{\star}}}}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})} \\ &- \frac{\exp(-0.5(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})} \left(\underline{\boldsymbol{G}}(t_0) + \frac{\exp(-(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})}\right)^{-1} \underbrace{\exp(-0.5(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})}_{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-(t-t_0)\boldsymbol{\Lambda}_{\text{eff}}^{-})}. \end{split}$$

Therefore, the result extends to $t > t_0$ as well.

D.3.2 High-dimensional limit for the risk curve

For $\mathsf{Err}(t) \coloneqq \|\mathbf{\Lambda}\|_{\mathsf{F}} \big(\frac{\mathbf{\Lambda}_{\mathsf{e}}}{\|\mathbf{\Lambda}\|_{\mathsf{F}}} - \frac{G_W(t)}{\sqrt{r_s}}\big)$, by Lemma 2, we have

$$\begin{split} \mathsf{Err}(t_{\mathrm{sc}}) &= \frac{-\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \\ &+ \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \left(\frac{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_{s}}} \boldsymbol{G}_{W}(0) + \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \right)^{-1} \!\! \boldsymbol{\Lambda}_{\mathrm{e}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})}, (\mathrm{D.1}) \end{split}$$

We define the block matrix forms

$$\boldsymbol{G}_W(t) = \begin{bmatrix} \boldsymbol{G}_{W,11}(t) & \boldsymbol{G}_{W,12}(t) \\ \boldsymbol{G}_{W,12}^\top(t) & \boldsymbol{G}_{W,22}(t) \end{bmatrix}, \ \boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{\text{eff}} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Lambda}_{22} \end{bmatrix}, \ \boldsymbol{\Lambda}_{\text{e},11} \coloneqq \boldsymbol{\Lambda}_{\text{eff}}, \ \boldsymbol{\Lambda}_{\text{e},22} \coloneqq \begin{bmatrix} \boldsymbol{\Lambda}_{22} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix},$$

where $G_{W,11}(t)$, $\Lambda_{\text{eff}} \in \mathbb{R}^{r_u \times r_u}$. Our proof strategy is as follows: In Proposition 7, we show that the off-diagonal and lower-right terms in (D.1) does not contribute to the high-dimensional limit. Then, in Proposition 8, we characterize the limit of the left-top terms. Finally, in Proposition 9, we prove the asymptotic behaviour of the risk curve.

Proposition 7. \mathcal{G}_{init} implies that $\|G_{W,12}(t\mathsf{T}_{\mathrm{eff}})\|_F^2 = o_d(r_s)$ and $\|G_{W,22}(t\mathsf{T}_{\mathrm{eff}})\|_F^2 = o_d(r_s)$.

Proof. We let

$$\boldsymbol{D}_1 \coloneqq \frac{\boldsymbol{\Lambda}_{\mathrm{e},11} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},11})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},11})}, \quad \boldsymbol{D}_2 \coloneqq \frac{\boldsymbol{\Lambda}_{\mathrm{e},22} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},22})}{\boldsymbol{I}_{d-r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},22})}, \quad \boldsymbol{Z} \coloneqq \begin{bmatrix} \boldsymbol{Z}_{1:r_u} \\ \boldsymbol{Z}_2 \end{bmatrix}.$$

and

$$\begin{bmatrix} \boldsymbol{S}_{11} & \boldsymbol{S}_{12} \\ \boldsymbol{S}_{12}^\top & \boldsymbol{S}_{22} \end{bmatrix} \coloneqq \left(\begin{bmatrix} \frac{\|\boldsymbol{\Delta}\|_{\mathbb{F}}}{\sqrt{r_s}} \boldsymbol{Z}_{1:r_u} \boldsymbol{Z}_{1:r_u}^\top + \boldsymbol{D}_1 & \frac{\|\boldsymbol{\Delta}\|_{\mathbb{F}}}{\sqrt{r_s}} \boldsymbol{Z}_{1:r_u} \boldsymbol{Z}_2^\top \\ \frac{\|\boldsymbol{\Delta}\|_{\mathbb{F}}}{\sqrt{r_s}} \boldsymbol{Z}_2 \boldsymbol{Z}_{1:r_u}^\top & \frac{\|\boldsymbol{\Delta}\|_{\mathbb{F}}}{\sqrt{r_s}} \boldsymbol{Z}_2 \boldsymbol{Z}_2^\top + \boldsymbol{D}_2 \end{bmatrix} \right)^{-1}.$$

where

$$egin{aligned} oldsymbol{S}_{11} &= \left(oldsymbol{D}_1 + oldsymbol{Z}_{1:r_u} \left(rac{\sqrt{r_s}}{\|oldsymbol{\Lambda}\|_{ ext{F}}} oldsymbol{I}_{r_s} + oldsymbol{Z}_2^ op oldsymbol{D}_2^{-1} oldsymbol{Z}_2^ op oldsymbol{D}_{1:r_u}
ight)^{-1} oldsymbol{Z}_{1:r_u}, \ oldsymbol{S}_{12} &= -\left(oldsymbol{D}_1 + oldsymbol{Z}_{1:r_u} \left(rac{\sqrt{r_s}}{\|oldsymbol{\Lambda}\|_{ ext{F}}} oldsymbol{I}_{r_s} + oldsymbol{Z}_2^ op oldsymbol{D}_2^{-1} oldsymbol{Z}_2
ight)^{-1} oldsymbol{Z}_{1:r_u} oldsymbol{D}_1^{-1} oldsymbol{Z}_{1:r_u} oldsymbol{Z}_2^ op oldsymbol{Z}_2^ o$$

Off-diagonal terms: By Proposition 26

$$\begin{split} \tilde{\boldsymbol{G}}_{W,12}(t) &\coloneqq \frac{\boldsymbol{\Lambda}_{\mathrm{e},11} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e},11})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},11})} \boldsymbol{S}_{12} \frac{\boldsymbol{\Lambda}_{\mathrm{e},22} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e},11})}{\boldsymbol{I}_{d-r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e},22})} \\ &= \exp(0.5t\boldsymbol{\Lambda}_{\mathrm{e},11}) \boldsymbol{Z}_{1:r_u} \left(\frac{\sqrt{r_s}}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}} \boldsymbol{I}_{r_s} + \boldsymbol{Z}_2^{\top} \boldsymbol{D}_2^{-1} \boldsymbol{Z}_2 + \boldsymbol{Z}_{1:r_u}^{\top} \boldsymbol{D}_1^{-1} \boldsymbol{Z}_{1:r_u} \right)^{-1} \boldsymbol{Z}_2^{\top} \exp(0.5t\boldsymbol{\Lambda}_{\mathrm{e},22}). \end{split}$$

We observe that $\lambda_{\max}(D_2) \leq \frac{1}{t}$ and $\frac{\sqrt{r_s}}{\|\mathbf{\Lambda}\|_F} \approx \kappa_{\text{eff}}$. By using Proposition 27 with $\mathcal{G}_{\text{init}}$ and $\tilde{t} \coloneqq t\kappa_{\text{eff}} \log \frac{d}{r}$, we write

$$\begin{split} &\frac{1}{r_s} \|\boldsymbol{G}_{W,12}(t\mathsf{T}_{\mathrm{eff}})\|_F^2 = \frac{1}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2} \|\tilde{\boldsymbol{G}}_{W,12}(\tilde{t})\|_F^2 \\ &\leq \frac{1}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2} \frac{O(1)}{(\kappa_{\mathrm{eff}} + \tilde{t})^2} \sum_{i=1}^{r_u \wedge r_s} \left(\lambda_{\max}(\boldsymbol{Z}_{1:r_u} \boldsymbol{Z}_{1:r_u}^\top) \exp(\tilde{t}\lambda_i) \wedge (\kappa_{\mathrm{eff}} + \tilde{t}) \frac{\lambda_i \exp(\tilde{t}\lambda_i)}{\exp(\tilde{t}\lambda_i) - 1} \right). \end{split} \tag{D.2}$$

For the heavy-tailed case ($\alpha \in [0, 0.5)$),

$$(\mathbf{D.2}) \leq \frac{O_{\alpha,\varphi,\beta}(1)}{r\log^2 d} \sum_{i < r} \mathbb{1}\{\tilde{t}\lambda_i \leq \log \frac{d}{r_s}\} + \frac{O_{\alpha,\varphi,\beta}(1)}{r^{1-\alpha}\log d} \sum_{i < r} \lambda_i \mathbb{1}\{\tilde{t}\lambda_i > \log \frac{d}{r_s}\} = o_d(1).$$

For the light-tailed case ($\alpha > 0.5$),

$$(D.2) \le \frac{O_{\alpha,r_s,\beta}(1)}{\log^2 d} \sum_{i \le r_s} \mathbb{1}\{\tilde{t}\lambda_i \le \log \frac{d}{r_u r_s}\} + \frac{O_{\alpha,\varphi,\beta}(1)}{\log d} \sum_{i \le r_s} \lambda_i \mathbb{1}\{\tilde{t}\lambda_i > \log \frac{d}{r_s r_u}\} = o_d(1).$$

Lower-right terms: By using Matrix-Inversion lemma, we have

$$-m{D}_2 + m{D}_2 m{S}_{22} m{D}_2 = -m{Z}_2 \left(rac{\sqrt{r_s}}{\|m{\Lambda}\|_{ ext{F}}} m{I}_{r_s} + m{Z}_{1:r_u}^ op m{D}_1^{-1} m{Z}_{1:r_u} + m{Z}_2^ op m{D}_2^{-1} m{Z}_2
ight)^{-1} m{Z}_2^ op.$$

We observe that $\lambda_{\max}(D_2) \leq \frac{1}{t}$ and $\frac{\sqrt{r_s}}{\|\mathbf{\Lambda}\|_{\mathrm{F}}} \approx \kappa_{\mathrm{eff}}$. By using $\tilde{t} \coloneqq t\kappa_{\mathrm{eff}} \log \frac{d}{r_s}$, we have

$$\frac{1}{\sqrt{r_s}}\boldsymbol{G}_{W,22}(t\mathsf{T}_{\mathrm{eff}}) = \frac{1}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}\boldsymbol{Z}_2 \left(\frac{\sqrt{r_s}}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}\boldsymbol{I}_{r_s} + \boldsymbol{Z}_{1:r_u}^{\top}\boldsymbol{D}_1^{-1}\boldsymbol{Z}_{1:r_u} + \boldsymbol{Z}_2^{\top}\boldsymbol{D}_2^{-1}\boldsymbol{Z}_2\right)^{-1}\boldsymbol{Z}_2^{\top}$$

$$\leq \frac{O(1)}{\|\mathbf{\Lambda}\|_{\mathrm{F}}} \mathbf{Z}_2 \left(\kappa_{\mathrm{eff}} \mathbf{I}_{r_s} + \tilde{t} \mathbf{Z}_2^{\top} \mathbf{Z}_2\right)^{-1} \mathbf{Z}_2^{\top}. \tag{D.3}$$

For the heavy-tailed case ($\alpha \in [0, 0.5)$),

$$\|(\mathbf{D.3})\|_F^2 \le \frac{O_{\alpha,\varphi,\beta}(1)}{r^{-2\alpha}\log^2 d} \frac{1}{t^2 r^{2\alpha}\log^2 d} = o_d(1).$$

For the light-tailed case ($\alpha > 0.5$),

$$\|(\mathbf{D.3})\|_F^2 \le O_{\alpha,r_s,\beta}(1) \frac{1}{t^2 \log^2 d} = o_d(1).$$

Proposition 8. For some c > 0, let $G(0) := \frac{c}{t} \mathbf{Z}_{1:r_u} \mathbf{Z}_{1:r_u}^{\top}$ and

$$t \in \begin{cases} \left(0, \frac{(r_{u_{\star}}+1)^{\alpha}}{\kappa_{\text{eff}}}\right), & \alpha \in [0, 0.5) \\ \left(0, \frac{(r_{u_{\star}}+1)^{\alpha}}{\kappa_{\text{eff}}}\right) \setminus \{j^{\alpha}: j \in \mathbb{N}\}, & \alpha > 0.5, \end{cases} \quad d \geq \begin{cases} \Omega_{\varphi,\alpha}(1) & \alpha \in [0, 0.5) \\ \Omega_{r_{s},\alpha}(1), & \alpha > 0.5. \end{cases}$$

We define

$$\begin{split} \mathsf{Err}_{r_u}(t_{\mathrm{sc}}) &\coloneqq \frac{-\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} \\ &+ \frac{\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} \left(\frac{\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} + \boldsymbol{G}(0)\right)^{-1} \!\! \frac{\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}. \end{split}$$

 G_{init} implies that

$$\frac{\|\mathsf{Err}_{r_u}(t\mathsf{T}_{\mathrm{eff}})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_{u_\star}} \lambda_j^2 \mathbb{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1).$$

Proof. We let

$$\boldsymbol{Z}_{1:r_u}\boldsymbol{Z}_{1:r_u}^\top \coloneqq \begin{bmatrix} \boldsymbol{Z}_{1:r_{u_\star}}\boldsymbol{Z}_{1:r_{u_\star}}^\top & \boldsymbol{Z}_1\boldsymbol{Z}_2^\top \\ \boldsymbol{Z}_2\boldsymbol{Z}_1^\top & \boldsymbol{Z}_2\boldsymbol{Z}_2^\top \end{bmatrix}, \qquad \boldsymbol{\Lambda}_{\text{eff}} \coloneqq \begin{bmatrix} \boldsymbol{\Lambda}_{\text{eff},11} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Lambda}_{\text{eff},22} \end{bmatrix},$$

where $\mathbf{\Lambda}_{ ext{eff},11} \in \mathbb{R}^{r_{u_\star} imes r_{u_\star}}$, $\mathbf{Z}_2 \in \mathbb{R}^{(r_u - r_{u_\star}) imes r_s}$. Let

$$\Gamma(t_{\mathrm{sc}}) \coloneqq \left(\frac{\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_{u}} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} + \frac{c}{t}\boldsymbol{Z}_{1:r_{u}}\boldsymbol{Z}_{1:r_{u}}^{\intercal}\right)^{-1} \ \ \text{and} \ \ \mathsf{Err}(t_{\mathrm{sc}},\boldsymbol{\Gamma}) \coloneqq \mathsf{Err}_{r_{u}}(t_{\mathrm{sc}}).$$

By using Proposition 22 and the events (H.3) and (L.2),

$$\boldsymbol{\Gamma}(t_{\mathrm{sc}}) \succeq \underline{\boldsymbol{\Gamma}}(t_{\mathrm{sc}}) \coloneqq \begin{bmatrix} \frac{10c}{t} \frac{r_s}{d} \boldsymbol{I}_{r_{u_{\star}}} + \frac{\boldsymbol{\Lambda}_{\mathrm{eff,11}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff,11}})}{\boldsymbol{I}_{r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff,11}})} & 0 \\ \\ 0 & \frac{2c}{t} \frac{r_s \log^{2.5} d}{d} \boldsymbol{I}_{r_u - r_{u_{\star}}} + \frac{\boldsymbol{\Lambda}_{\mathrm{eff,22}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff,22}})}{\boldsymbol{I}_{r_u - r_{u_{\star}}} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff,22}})} \end{bmatrix}^{-1}.$$

For the upper bound, by using $\varepsilon=\frac{1}{c\log^3 d}$ in Proposition 24, and the events (H.1), (L.1) and (H.3), (L.2), we have for $t\leq (r_{u_\star}+1)^\alpha\log\frac{d}{r_s}$,

$$\boldsymbol{\Gamma}(t_{\rm sc}) \preceq \overline{\boldsymbol{\Gamma}}(t_{\rm sc}) \coloneqq \begin{bmatrix} \frac{0.2/t}{\log^4 d} \frac{r_s}{d} \boldsymbol{I}_{r_{u_\star}} + \frac{\boldsymbol{\Lambda}_{\rm eff,11} \exp(-t\boldsymbol{\Lambda}_{\rm eff,11})}{\boldsymbol{I}_{r_{u_\star}} - \exp(-t\boldsymbol{\Lambda}_{\rm eff,11})} & 0 \\ \\ 0 & \frac{-1.1/t}{\log^{1/2} d} \frac{r_s}{d} \boldsymbol{I}_{r_u - r_{u_\star}} + \frac{\boldsymbol{\Lambda}_{\rm eff,22} \exp(-t\boldsymbol{\Lambda}_{\rm eff,22})}{\boldsymbol{I}_{r_u - r_{u_\star}} - \exp(-t\boldsymbol{\Lambda}_{\rm eff,22})} \end{bmatrix}^{-1}.$$

28

Therefore, for $t < \frac{(r_{u_*}+1)^{\alpha}}{\kappa_{\text{eff}}}$, by Corollary 8, we have

$$\begin{split} \frac{\|\mathsf{Err}(t\mathsf{T}_{\mathrm{eff}},\mathbf{\Gamma})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} &\geq \frac{\|\mathsf{Err}(t\mathsf{T}_{\mathrm{eff}},\underline{\mathbf{\Gamma}})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} = \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_u} \lambda_j^2 \mathbbm{1}\{\frac{1}{\lambda_j} > t\kappa_{\mathrm{eff}}\} + o_d(1) \\ &= 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_{u_\star}} \lambda_j^2 \mathbbm{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1). \end{split}$$

On the other hand, by Corollary 8, we have

$$\begin{split} \frac{\|\mathsf{Err}(t\mathsf{T}_{\mathrm{eff}},\pmb{\Gamma})\|_F^2}{\|\pmb{\Lambda}\|_{\mathrm{F}}^2} &\leq \frac{\|\mathsf{Err}(t\mathsf{T}_{\mathrm{eff}},\overline{\pmb{\Gamma}})\|_F^2}{\|\pmb{\Lambda}\|_{\mathrm{F}}^2} \\ &= \frac{1}{\|\pmb{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_{u_\star}} \lambda_j^2 \mathbbm{1}\{\frac{1}{\lambda_j} > t\kappa_{\mathrm{eff}}\} + \frac{1}{\|\pmb{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=r_{u_\star}+1}^{r_u} \lambda_j^2 + o_d(1) \\ &= 1 - \frac{1}{\|\pmb{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_{u_\star}} \lambda_j^2 \mathbbm{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1). \end{split}$$

Proposition 9. G_{init} implies that

$$\frac{\|\mathsf{Err}\,(t\mathsf{T}_{\mathrm{eff}})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_{u_\star}} \lambda_j^2 \mathbb{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1),$$

for

$$t \in \begin{cases} (0, \infty), & \alpha \in [0, 0.5) \\ (0, \infty) \setminus \{j^{\alpha} : j \in \mathbb{N}\}, & \alpha > 0.5, \end{cases} \quad d \ge \begin{cases} \Omega_{\varphi, \alpha}(1) & \alpha \in [0, 0.5) \\ \Omega_{r_s, \alpha}(1), & \alpha > 0.5. \end{cases}$$

Proof. We recall that

$$m{D}_1 \coloneqq rac{m{\Lambda}_{\mathrm{e},11} \exp(-tm{\Lambda}_{\mathrm{e},11})}{m{I}_{r_u} - \exp(-tm{\Lambda}_{\mathrm{e},11})}, \quad m{D}_2 \coloneqq rac{m{\Lambda}_{\mathrm{e},22} \exp(-tm{\Lambda}_{\mathrm{e},22})}{m{I}_{d-r_u} - \exp(-tm{\Lambda}_{\mathrm{e},22})}, \quad m{Z} \coloneqq egin{bmatrix} m{Z}_{1:r_u} \ m{Z}_2 \end{bmatrix},$$

and

$$\begin{split} \mathsf{Err}(t_{\mathrm{sc}}) &= \frac{-\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \\ &+ \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \left(\frac{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_{s}}} \boldsymbol{Z} \boldsymbol{Z}^{\top} + \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \right)^{-1} \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{e}})}{\boldsymbol{I}_{d} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{e}})} \\ &= \begin{bmatrix} \mathsf{Err}_{r_{u}}(t_{\mathrm{sc}}) & \frac{-1}{\sqrt{r_{s}}} \boldsymbol{G}_{W,12}(t_{\mathrm{sc}}) \\ \frac{-1}{\sqrt{r_{s}}} \boldsymbol{G}_{W,12}(t_{\mathrm{sc}}) & \frac{-1}{\sqrt{r_{s}}} \boldsymbol{G}_{W,22}(t_{\mathrm{sc}}) \end{bmatrix}. \end{split}$$

Note that by Proposition 35, in the time scale we consider we have $\frac{1-o_d(1)}{t\mathsf{T}_{\mathrm{eff}}} \leq \lambda_{\min}(\boldsymbol{D}_2) \leq \lambda_{\max}(\boldsymbol{D}_2) \leq \frac{1}{t\mathsf{T}_{\mathrm{eff}}}$. The by using (E.1),

$$\begin{split} & \mathsf{Err}_{r_u}(t_{\mathrm{sc}}) = \frac{-\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} \\ & + \frac{\boldsymbol{\Lambda}_{\mathrm{eff}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} \bigg(\frac{\Theta(1)}{\frac{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} + t} \boldsymbol{Z}_{1:r_u} \boldsymbol{Z}_{1:r_u}^\top + \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})} \bigg)^{-1} \!\! \frac{\boldsymbol{\Lambda}_{\mathrm{e}} \exp(-0.5t\boldsymbol{\Lambda}_{\mathrm{eff}})}{\boldsymbol{I}_{r_u} - \exp(-t\boldsymbol{\Lambda}_{\mathrm{eff}})}. \end{split}$$

By Propositions 7 and 8, we have

$$\frac{\|\mathsf{Err}\,(t\mathsf{T}_{\mathrm{eff}})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{(r_s \wedge r)} \lambda_j^2 \mathbb{1}\{t\kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1),$$

for

$$t \in \begin{cases} \left(0, \frac{(r_{u_{\star}} + 1)^{\alpha}}{\kappa_{\text{eff}}}\right), & \alpha \in [0, 0.5) \\ \left(0, \frac{(r_{u_{\star}} + 1)^{\alpha}}{\kappa_{\text{eff}}}\right) \setminus \{j^{\alpha} : j \in \mathbb{N}\}, & \alpha > 0.5. \end{cases}$$
(D.4)

To extend the limit for $t>\frac{(r_{u_\star}+1)^\alpha}{\kappa_{\rm eff}},$ we observe that

- $\|\operatorname{Err}(t)\|_F^2$ non increasing since it corresponds to the objective under (GF).
- The global optimum of (GF) and the previous item with (D.4) guarantees that for $t > \frac{(r_{u_*}+1)^{\alpha}}{\kappa_{\text{eff}}}$,

$$1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{(r_s \wedge r)} \!\! \lambda_j^2 \mathbbm{1} \big\{ t \kappa_{\mathrm{eff}} \geq \tfrac{1}{\lambda_j} \big\} \leq \frac{\|\mathsf{Err} \, (t \mathsf{T}_{\mathrm{eff}})\|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \leq 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{(r_s \wedge r)} \!\! \lambda_j^2 \mathbbm{1} \big\{ t \kappa_{\mathrm{eff}} \geq \tfrac{1}{\lambda_j} \big\} + o_d(1).$$

Therefore, the statement extends to $t>\frac{(r_{u_\star}+1)^\alpha}{\kappa_{\rm eff}}.$

D.4 Proof of Theorem 2

We redefine the time-scale and effective-width as:

$$t_{\rm sc} = t\sqrt{r_s} \|\mathbf{\Lambda}\|_{\rm F}, \quad \kappa_{\rm eff} = \begin{cases} r^{\alpha}/\eta, & \alpha \in [0, 0.5) \\ 1/\eta, & \alpha > 0.5 \end{cases}, \qquad \mathsf{T}_{\rm eff} = \kappa_{\rm eff} \sqrt{r_s} \|\mathbf{\Lambda}\|_{\rm F} \log d/r_s.$$

and

$$r_u = \begin{cases} r, & \alpha \in [0, 0.5) \\ \lceil \log^{2.5} d \rceil, & \alpha > 0.5 \end{cases}, \qquad r_{u_\star} := \begin{cases} \lfloor r_s (1 - \log^{-1/s} d) \wedge r \rfloor, & \alpha \in [0, 0.5) \\ r_s, & \alpha > 0.5. \end{cases}$$

We consider the learning rate and fine-tuning sample size given as

$$\eta \asymp \frac{1}{d} \begin{cases} \frac{1}{r^{\alpha} \log^{20}(1+d/r_s)}, & \alpha \in [0,0.5) \\ \frac{1}{r_s^{\alpha\alpha+3} \log^{18} d}, & \alpha > 0.5 \end{cases} \quad \text{and} \quad N_{\mathrm{Ft}} \asymp r_s^2 \log^5 d.$$

We define the effective learning rate η and the hitting time \mathcal{T}_{hit} as follows:

$$\eta \coloneqq \frac{\eta/2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}}, \quad \mathcal{T}_{\mathrm{hit}} \coloneqq \left\{ t \ge 0 \mid 1 - \frac{\|\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{G}_t \mathbf{\Lambda}^{\frac{1}{2}} \|_F^2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \le \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=(r_s \wedge 1)+1}^r \lambda_j^2 + \frac{10}{\log^{\frac{1}{8}} d} \right\}.$$

We note that bounding \mathcal{T}_{hit} suffices to derive sample complexity since by Proposition 11, we have

$$R(\boldsymbol{W}_{t}^{\text{final}}) \leq 1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{G}_{t}\boldsymbol{\Lambda}^{\frac{1}{2}}\|_{F}^{2}}{\|\boldsymbol{\Lambda}\|_{F}^{2}} + \frac{O(1)}{\log d}.$$

The main statement of this part is as follows:

Proposition 10. The intersection of the following events hold with probability $1 - o_d(1/d^2) - \Omega(1/r_s^2)$:

1. We have

$$\mathcal{T}_{\text{hit}} \le \begin{cases} \frac{1}{2\eta} \left(r_s \left(1 - \log^{-1/s} d \right) \wedge r \right)^{\alpha} \log \left(\frac{20d \log^{\frac{3}{4}} (1 + d/r_s)}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2\eta} r_s^{\alpha} \log \left(20 \frac{d \log^{3/4} d}{r_s} \right), & \alpha > 0.5. \end{cases}$$

2. For t > 0,

•
$$\mathcal{A}(t\mathsf{T}_{\mathrm{eff}}, \pmb{\theta}_j) = \mathbb{1}\{\eta t \kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j}\} + o_d(1) \text{ for } t \neq \lim_{d \to \infty} \frac{1}{\eta \kappa_{\mathrm{eff}} \lambda_j} \text{ and } j \leq r_{u_\star}.$$

•
$$\|\mathbf{\Lambda}\|_{\mathrm{F}}^2 - \|\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{G}_{t\mathsf{T}_{\mathrm{eff}}} \mathbf{\Lambda}^{\frac{1}{2}} \|_F^2 = 1 - \sum_{j=1}^{r_{u_{\star}}} \lambda_j^2 \mathbb{1} \{ \eta t \kappa_{\mathrm{eff}} \geq \frac{1}{\lambda_j} \} + o_d(\|\mathbf{\Lambda}\|_{\mathrm{F}}^2).$$

Proof. By using Lemma 1 and Corollary 5, we have with probability $1 - o_d(1/d^2) - \Omega(1/r_s^2)$:

$$\mathcal{T}_{\text{bad}} \geq \begin{cases} \frac{1}{2\eta} \left(r_s \left(1 - \log^{\frac{-1}{2}} d \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2\eta} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha > 0.5, \end{cases}$$

where \mathcal{T}_{bad} is defined in (F.14). Given the lower bound, by Proposition 14, and the third item of Proposition 15, we have the first item.

For the second item, by Proposition 14, and Proposition 15 (for the lower bound) and Proposition 16 (for the upper bound), we have for $r_{u_{\star}} \times r_{u_{\star}}$ dimensional top left submatrices $G_{t,11}$ and Λ_{11} ,

$$\frac{1}{\frac{1.2}{C_{lb}}\frac{d}{r_s}\exp\left(-2t\eta\mathbf{\Lambda}_{11}\right)+1} - o_d(1) \leq \mathbf{G}_{t,11} \leq \left(\frac{C_{ub}r_s}{d}\exp\left(2\eta t\mathbf{\Lambda}_{11}\right) \wedge 1\right) + o_d(1), \text{ (D.5)}$$

and

$$\begin{split} \|\mathbf{\Lambda}\|_F^2 - \|\mathbf{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_t \mathbf{\Lambda}^{\frac{1}{2}} \|_F^2 &\geq \sum_{i=1}^{r_u} \lambda_i^2 \left(1 - \frac{C_{\text{ub}} r_s}{d} \exp\left(2 \eta t \lambda_i \right) \right)_+ - o_d(\|\mathbf{\Lambda}\|_F^2) \\ \|\mathbf{\Lambda}\|_F^2 - \|\mathbf{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_t \mathbf{\Lambda}^{\frac{1}{2}} \|_F^2 &\leq \sum_{i=(r_u, \Lambda^r)+1}^{r} \lambda_i^2 + \sum_{i=1}^{r_{u_\star}} \lambda_i^2 \left(1 - \frac{1}{\frac{1\cdot 2}{C_{\text{lb}}} \frac{d}{r_s} \exp\left(-2 t \eta \lambda_i \right) + 1} \right)^2 + o_d(\|\mathbf{\Lambda}\|_F^2). \end{split}$$

for

$$t \leq \begin{cases} \frac{1}{2\eta} \left(r_s \left(1 - \log^{-1/8} d \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2\eta} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha > 0.5. \end{cases}$$
(D.7)

where

$$C_{\rm ub} = \begin{cases} 2.5 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2, & \alpha \in [0, 0.5) \\ 15, & \alpha > 0.5, \end{cases} \qquad C_{\rm lb} = \frac{1}{15} \begin{cases} \log^{-1/2} d, & \alpha \in [0, 0.5) \\ r_s^{-6}, & \alpha > 0.5. \end{cases}$$

The high-dimensional limits of the alignment and risk up to the time horizon in (D.7) follow from (D.5) (for the alignment), and from (D.6) (for the risk) by Proposition 36. Proposition 21 then allows us to extend these results beyond the time limit in (D.7), yielding the full statement.

D.5 Proof of Corollary 1 and Corollary 2

Finally, we derive the scaling of prediction risk under power-law second-layer coefficients. Since Corollary 2 is a rescaled version of Corollary 1, we will only consider the latter.

Proof of Corollary 1. We will prove heavy and light-tailed cases separately.

Heavy-tailed case ($\alpha \in [0, 0.5)$): We define $C := \left(\frac{(1-\beta)\sqrt{\varphi}}{\sqrt{1-2\alpha}}\right)^{\frac{1}{\alpha}}$. We first fix a $(C\varphi)^{\alpha} > t > 0$. By Proposition 9, for any $d \ge \Omega_{\varphi,\alpha}(1)$, we have with probability at least $1 - o(1/d^2)$

$$\mathcal{R}(tr\log d) = \underbrace{1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_s} \lambda_j^2 \mathbb{1}\left\{\frac{tr^{\alpha}}{C^{\alpha}} \ge \frac{1 \pm o_d(1)}{\lambda_j}\right\}}_{:=\mathcal{R}_d((Ct)^{\frac{1}{\alpha}})} + o_d(1)$$

where we define $\mathcal{R}_d((Ct)^{\frac{1}{\alpha}})$ to isolate the main term and make the dependence on the ambient dimension explicit. By using $\lambda_j=j^{-\alpha}$ in the indicator function, we can rewrite

$$\mathcal{R}_d(t) = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_s} \lambda_j^2 \mathbb{1}\{(1 \pm o_d(1))t \ge \frac{j}{r}\}$$

We define a sequence of measures supported on $\{j/r: j \in [r]\}$, where $\mu_d\{\frac{j}{r}\} \propto j^{-2\alpha}$ for $j=1,\cdots,r$. We observe the following:

• μ_d converges weakly weakly to a limiting probability measure μ supported on [0,1], with cumulative distribution function

$$\mu\{[0,c)\} = \begin{cases} c^{1-2\alpha}, & c < 1\\ 1 & x \ge 1 \end{cases}$$

• Moreover, the risk can be expressed as

$$\mathcal{R}_d(t) = 1 - (1 \pm o_d(1)) \mathbb{E}_{X \sim \mu_d} [\mathbb{1}\{(1 \pm o_d(1))t \land \varphi \ge X\}]$$

By the Portmanteau theorem [Dur93], it follows that for any fixed $t \in (0, \varphi)$,

$$\mathcal{R}_d(t) \to 1 - t^{1-2\alpha}$$
.

almost surely as $d \to \infty$. The almost sure convergence follows from the Borel-Cantelli lemma [Dur93] applied to the failure probabilities.

To extend this result to $t \ge \varphi$, we observe that by (GF), $\mathcal{R}_d(t)$ is non-increasing and $\inf_{t \ge 0} \mathcal{R}_d(t) \ge (1 - \varphi^{1-2\alpha})_+ - o_d(1)$. Hence, for all t > 0, we obtain

$$\mathcal{R}_d(t) = (1 - t^{1-2\alpha})_+ \vee (1 - \varphi^{1-2\alpha})_+$$

The desired result for a fixed t > 0 follows by a change of variable. Finally, since the risk curves are continuous in t, the almost sure convergence extends to all t > 0 pointwise.

Light-tailed case ($\alpha > 0.5$): For this part, we consider the probability space conditioned on \mathcal{G}_{init} which holds with probability at least $1 - o(1/r_s^2)$. We define

$$\mathcal{Z} \coloneqq \sum_{j=1}^{\infty} j^{-2\alpha}, \quad C \coloneqq (r_s \mathcal{Z})^{\frac{1}{2\alpha}}.$$

We first fix a $t \in (0, (Cr_s)^{\alpha}) \setminus \{j^{\alpha} : j \in \mathbb{N}\}$. By Proposition 9, for any $d \geq \Omega_{r_s,\alpha}(1)$, we have,

$$\mathcal{R}(t\log d) = \underbrace{1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_s} \lambda_j^2 \mathbb{I}\left\{\frac{t}{C^{\alpha}} \ge \frac{1 \pm o_d(1)}{\lambda_j}\right\}}_{:=\mathcal{R}_d((Ct)^{\frac{1}{\alpha}})} + o_d(1).$$

By using $\lambda_j = j^{-\alpha}$ in the indicator function, we rewrite

$$\mathcal{R}_d(t) = 1 - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_s} \lambda_j^2 \mathbb{1}\{(1 \pm o_d(1))t \ge j\}$$

We define a sequence of measures supported on \mathbb{N} , where $\mu_d\{j\} \propto j^{-2\alpha}$ for $j=1,\cdots,r_s$. We observe the following:

- μ_d converges weakly weakly to a limiting probability measure μ supported on \mathbb{N} , such that $\mu\{j\} = \frac{j^{-2\alpha}}{\mathcal{Z}}$.
- Moreover, the risk can be expressed as

$$\mathcal{R}_d(t) = \mathbb{E}_{X \sim \mu_d} [\mathbb{1}\{(1 \pm o_d(1)t \vee r_s < X\}]$$

Since $t \notin \mathbb{N}$, we have

$$\mathbb{R}_d(t) \to \mu([t \vee r_s, \infty)).$$

By observing that $\mu([t,\infty)) \in \Theta(t^{1-2\alpha})$, the result follows for a fixed $t \in (0,(Cr_s)^\alpha) \setminus \{j^\alpha: j \in \mathbb{N}\}$. Since the limit is piecewise continuous and non increasing, it is sufficient to take a union over $t \in \{0.5, 1.5, \cdots, r_s + 0.5\}$ to extend the result for all t > 0.

E Details of the Fine-tuning Step

In this part, we describe how to efficiently solve the empirical risk minimization problem used in the fine-tuning step of Algorithm 1. Recall that this step aims to find a rotation matrix $\Omega \in \mathbb{R}^{r_s \times r_s}$ that aligns the learned features with the teacher directions by minimizing the empirical loss over N_{Ft} fresh samples:

$$\Omega_* = \underset{\mathbf{\Omega} \in \mathbb{R}^{r_s \times r_s}}{\operatorname{arg \, min}} \sum_{j=1}^{N_{\mathrm{Ft}}} \mathcal{L}(\mathbf{W}_t \mathbf{\Omega}; (\mathbf{x}_{t+j}, y_{t+j})), \tag{E.1}$$

where each sample loss is given by

$$\mathcal{L}(\boldsymbol{W}_{t}\boldsymbol{\Omega};(\boldsymbol{x}_{t+j},y_{t+j})) = \frac{1}{16} \left(y_{t+j} - \frac{1}{\sqrt{r_{s}}} \text{Tr}(\boldsymbol{\Omega}\boldsymbol{\Omega}^{\top} \boldsymbol{W}_{t}^{\top} (\boldsymbol{x}_{t+j} \boldsymbol{x}_{t+j}^{\top} - \boldsymbol{I}_{d}) \boldsymbol{W}_{t}) \right)^{2}.$$

Let us define $A_j := W_t^\top (x_{t+j} x_{t+j}^\top - I_d) W_t$. We observe that the loss becomes quadratic in the symmetric matrix positive semidefinite matrix $S := \Omega \Omega^\top$. Then, the fine-tuning objective reduces to a standard least squares problem over the cone of symmetric matrix positive semidefinite matrices:

$$S_* := \underset{\boldsymbol{S} = \boldsymbol{S}^{\top}, \boldsymbol{S} \succeq 0}{\arg \min} \underbrace{\frac{1}{2N_{\mathrm{Ft}}} \sum_{j=1}^{N_{\mathrm{Ft}}} \left(\sqrt{r_s} y_{t+j} - \mathrm{Tr}(\boldsymbol{S} \boldsymbol{A}_j) \right)^2}_{:=\mathrm{Ft}(\boldsymbol{S})}. \tag{E.2}$$

For the following, we also define the global minimum of the least square objective in (E.2) as:

$$S_{\text{glob}} := \underset{\boldsymbol{S} \in \mathbb{R}^{r_s \times r_s}}{\min} \operatorname{Ft}(\boldsymbol{S}). \tag{E.3}$$

$$S \in \mathbb{R}^{r_s \times r_s}$$

E.1 Characterizing the Minimum

Since the fine-tuning objective reduces to a least squares regression problem over symmetric matrices, we can write

$$\mathrm{Ft}(oldsymbol{S}) = \mathrm{Ft}(oldsymbol{S}_{\mathrm{glob}}) + \mathrm{Tr}ig((oldsymbol{S} - oldsymbol{S}_{\mathrm{glob}}) \mathsf{L}(oldsymbol{S} - oldsymbol{S}_{\mathrm{glob}})ig)$$

where L is defined as the linear operator acting on symmetric matrices via

$$\mathsf{L}(oldsymbol{S})\coloneqq rac{1}{2N_{\mathrm{Ft}}}\sum_{i=1}^{N_{\mathrm{Ft}}}\mathrm{Tr}(oldsymbol{S}oldsymbol{A}_{j})oldsymbol{A}_{j},$$

which corresponds to the empirical second moment operator associated with the covariates A_j . We note that the operator L is self-adjoint and positive semi-definite on the space of symmetric matrices, and we can write the characterization in (E.2) equivalently

$$egin{aligned} oldsymbol{S}_* \coloneqq & rg \min_{oldsymbol{S} \in \mathbb{R}^{r_s imes r_s} > 0} \mathrm{Tr}ig((oldsymbol{S} - oldsymbol{S}_{\mathrm{glob}}) \mathsf{L} (oldsymbol{S} - oldsymbol{S}_{\mathrm{glob}}) ig). \end{aligned}$$

We define the projection on the cone of symmetric positive semi-definite matrices as:

$$\Pi(\widetilde{\boldsymbol{S}}) \coloneqq \mathop{\arg\min}_{\substack{\boldsymbol{S} \in \mathbb{R}^{r_s \times r_s} \\ \boldsymbol{S} = \boldsymbol{S}^{\top}, \boldsymbol{S} \succeq 0}} \|\boldsymbol{S} - \widetilde{\boldsymbol{S}}\|_F^2.$$

In the following, we will show that the operator L is close to the identity, and thus, S_* is close to $\Pi \circ L(S_{\text{glob}})$. Before proceeding, we make the following observations:

• We observe that by the first-order optimality condition applied in (E.3), we have

$$\mathsf{L}(\boldsymbol{S}_{\mathrm{glob}}) = \frac{\sqrt{r_s}}{2N_{\mathrm{Ft}}} \sum_{i=1}^{N_{\mathrm{Ft}}} y_{t+j} \boldsymbol{A}_j. \tag{E.4}$$

• By the generalized Pythagorean theorem [Bub14, Lemma 3.1], we have

$$\begin{aligned} \|\boldsymbol{S}_{*} - \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}})\|_{F}^{2} &\leq \|\boldsymbol{S}_{*} - \mathsf{L}(\boldsymbol{S}_{\text{glob}})\|_{F}^{2} - \|\boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}) - \mathsf{L}(\boldsymbol{S}_{\text{glob}})\|_{F}^{2} \\ &= \mathrm{Ft}(\boldsymbol{S}_{*}) - \mathrm{Ft}(\boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}})) \\ &- \mathrm{Tr}((\boldsymbol{S}_{*} - \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}))(\mathsf{L} - \mathsf{Id})(\boldsymbol{S}_{*} + \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}))), \end{aligned} \tag{E.5}$$

where we use ld to denote the identity map on symmetric matrices.

E.2 Computing the Minimum

We define the approximate solution for (E.1) as:

$$\hat{\mathbf{\Omega}} := \left(\mathsf{\Pi} \circ \mathsf{L}(\mathbf{S}_{\mathrm{glob}}) \right)^{\frac{1}{2}},\tag{E.6}$$

where $S \to S^{1/2}$ denotes the square root operator on symmetric positive semidefinite matrices. Note that the approximation in (E.6) can be computed by taking the spectral decomposition of $\mathsf{L}(S_{\mathrm{glob}})$ given in (E.4), which requires $\tilde{O}(dr_s^3)$ including the computation of $\mathsf{L}(S_{\mathrm{glob}})$. This is negligible compared to the feature learning phase, whose complexity scales as $O(Tdr_s)$. The following statement shows that $\hat{\Omega}$ is sufficiently close to the fine-tuning solution Ω^* :

Proposition 11. Suppose $N_{\rm Ft} \ge r_s^2 \log^5 d$. Then, with probability at least $1 - 2d^{-3}$, the final risk incurred by $W_t \hat{\Omega}$ is close to that of the optimal fine-tuning solution:

$$R(\boldsymbol{W}_{t}\hat{\boldsymbol{\Omega}}) \leq R(\boldsymbol{W}_{t}\boldsymbol{\Omega}_{*}) + \frac{1}{\log d} \leq 1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{G}_{t}\boldsymbol{\Lambda}^{\frac{1}{2}}\|_{F}^{2}}{\|\boldsymbol{\Lambda}\|_{F}^{2}} + \frac{O(1)}{\log d}.$$

E.2.1 Proof of Proposition 11

We define the operator norm of L as

$$\|\mathsf{L}\|_2 = \sup_{\substack{\boldsymbol{S} \in \mathbb{R}^{r_s \times r_s} \\ \boldsymbol{S} = \boldsymbol{S}^\top}} \|\mathsf{L}(\boldsymbol{S})\|_F.$$

We consider the intersection of the following events

•
$$\|\mathsf{L} - \mathsf{Id}\|_2 \leq \frac{6}{\sqrt{\log d}}$$

•
$$\left\| \frac{1}{2N_{\mathrm{Ft}}} \sum_{j=1}^{N_{\mathrm{Ft}}} y_{t+j} A_j - \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}} \mathbf{W}_t^{\top} \mathbf{\Theta} \mathbf{\Lambda} \mathbf{\Theta}^{\top} \mathbf{W}_t \right\|_F^2 \leq \frac{1}{\log d}$$

We note that for $d \ge \Omega(1)$ the first item holds with probability $1 - d^{-3}$ by Proposition 31, where we choose C = 5 and $u = \log d$, and the second item holds follows with probability $1 - d^{-3}$ by Proposition 32 where we choose C = 16. Given the events, we have

$$\begin{split} R(\boldsymbol{W}_{t}\hat{\boldsymbol{\Omega}}) &= \frac{1}{r_{s}} \left\| \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}) - \frac{\sqrt{r_{s}}}{\|\boldsymbol{\Lambda}\|_{\text{F}}} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \right\|_{F}^{2} + \left(1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_{t} \boldsymbol{\Lambda}^{\frac{1}{2}} \|_{F}^{2}}{\|\boldsymbol{\Lambda}\|_{\text{F}}^{2}}\right) \\ &\stackrel{(a)}{\leq} \frac{1}{r_{s}} \left\| \mathsf{L}(\boldsymbol{S}_{\text{glob}}) - \frac{\sqrt{r_{s}}}{\|\boldsymbol{\Lambda}\|_{\text{F}}} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \right\|_{F}^{2} + \left(1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_{t} \boldsymbol{\Lambda}^{\frac{1}{2}} \|_{F}^{2}}{\|\boldsymbol{\Lambda}\|_{\text{F}}^{2}}\right) \\ &\stackrel{(b)}{\leq} \frac{1}{\log d} + R(\boldsymbol{W}_{t} \boldsymbol{\Omega}_{*}). \end{split}$$

where we use the convexity of the cone of symmetric positive semi-definite matrices in (a) and the second event above in (b). By using (E.5), we have

$$\left\|\boldsymbol{S}_* - \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\mathrm{glob}})\right\|_F \leq \|\mathsf{L} - \mathsf{Id}\|_2 \|\boldsymbol{S}_* + \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\mathrm{glob}})\|_F.$$

Therefore,

$$\|S_*\|_F \leq rac{1 + \|\mathsf{L} - \mathsf{Id}\|_2}{1 - \|\mathsf{L} - \mathsf{Id}\|_2} \|\mathsf{\Pi} \circ \mathsf{L}(S_{\mathrm{glob}})\|_F,$$

and thus,

$$\left\|\boldsymbol{S}_* - \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\mathrm{glob}})\right\|_F \leq \frac{2\|\mathsf{L} - \mathsf{Id}\|_2 \|\boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\mathrm{glob}})\|_F}{1 - \|\mathsf{L} - \mathsf{Id}\|_2} \overset{(c)}{\leq} \frac{15r_s}{\sqrt{\log d}}$$

where we followed the reasoning in (a)-(b) to bound $\|\Pi \circ \mathsf{L}(S_{\mathrm{glob}})\|_F$ in (c) . Therefore,

$$R(\boldsymbol{W}_{t}\boldsymbol{\Omega}_{*}) = \frac{1}{r_{s}} \left\| \boldsymbol{S}_{*} - \frac{\sqrt{r_{s}}}{\|\boldsymbol{\Lambda}\|_{F}} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \right\|_{F}^{2} + \left(1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_{t} \boldsymbol{\Lambda}^{\frac{1}{2}} \|_{F}^{2}}{\|\boldsymbol{\Lambda}\|_{F}^{2}}\right)$$

$$\leq \frac{2}{r_s} \| \boldsymbol{S}_* - \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}) \|_F^2 + \frac{2}{r_s} \| \boldsymbol{\Pi} \circ \mathsf{L}(\boldsymbol{S}_{\text{glob}}) - \frac{\sqrt{r_s}}{\|\boldsymbol{\Lambda}\|_F} \boldsymbol{W}_t^\top \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^\top \boldsymbol{W}_t \|_F^2 \\
+ \left(1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_t \boldsymbol{\Lambda}^{\frac{1}{2}} \|_F^2}{\|\boldsymbol{\Lambda}\|_F^2} \right) \\
\leq \frac{O(1)}{\log d} + \left(1 - \frac{\|\boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{G}_t \boldsymbol{\Lambda}^{\frac{1}{2}} \|_F^2}{\|\boldsymbol{\Lambda}\|_F^2} \right).$$

F Deferred Proofs for Online SGD

F.1 Preliminaries

We consider

$$y_{t+1} = \frac{1}{\|\boldsymbol{\Lambda}\|_F} \sum_{j=1}^r \lambda_j \left(\left\langle \boldsymbol{\theta}_j, \boldsymbol{x}_{t+1} \right\rangle^2 - 1 \right) + \epsilon_{t+1} \text{ and } \hat{y}(\boldsymbol{W}_t; \boldsymbol{x}_{t+1}) = \frac{1}{\sqrt{r_s}} \sum_{j=1}^{r_s} \left\langle \boldsymbol{w}_{t,j}, \boldsymbol{x}_{t+1} \right\rangle^2 - 1,$$

where $\|\epsilon_{t+1}\|_{\psi_2} \leq \sigma$. We use $\hat{y}_{t+1} := \hat{y}(\boldsymbol{W}_t; \boldsymbol{x}_{t+1})$ and consider

- The loss function is $\mathcal{L}(\pmb{W}_t;(\pmb{x}_{t+1},y_{t+1})) = \frac{1}{16} ig(y_{t+1} \hat{y}_{t+1}ig)^2$
- The Euclidean gradient is $\nabla \mathcal{L}(\boldsymbol{W}_t) = \frac{-1}{4\sqrt{r_s}} (y_{t+1} \hat{y}_{t+1}) \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \boldsymbol{W}_t$. Therefore, we have

$$\nabla_{\operatorname{St}} \mathcal{L}(\boldsymbol{W}_t) = \frac{-1/4}{\sqrt{r_s}} \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top} \right) \left(y_{t+1} - \hat{y}_{t+1} \right) \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \boldsymbol{W}_t.$$

• We recall that $G_t = \Theta^\top W_t W_t^\top \Theta$.

Then, (SGD) reads

$$\widetilde{\boldsymbol{W}}_{t+1} = \boldsymbol{W}_{t} + \frac{\eta/4}{\sqrt{r_{s}}} \underbrace{\left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}\right) \left(y_{t+1} - \hat{y}_{t+1}\right) \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \boldsymbol{W}_{t}}_{:=\nabla_{St} \boldsymbol{L}_{t+1}}$$

$$\boldsymbol{W}_{t+1} = \widetilde{\boldsymbol{W}}_{t+1} \left(\boldsymbol{I}_{r_{s}} + \frac{\eta^{2}/16}{r_{s}} \underbrace{\nabla_{St} \boldsymbol{L}_{t+1}^{\top} \nabla_{St} \boldsymbol{L}_{t+1}}_{:=\boldsymbol{\mathcal{P}}_{t+1}}\right)^{-1/2}. \tag{SGD}$$

We observe that

$$\frac{\eta^{2}/16}{r_{s}} \mathcal{P}_{t+1} = \frac{\eta^{2}/16}{r_{s}} (y_{t+1} - \hat{y}_{t+1})^{2} \mathbf{W}_{t}^{\top} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^{\top} (\mathbf{I}_{d} - \mathbf{W}_{t} \mathbf{W}_{t}^{\top}) \mathbf{x}_{t+1} \mathbf{x}_{t+1}^{\top} \mathbf{W}_{t}$$

$$= \frac{\eta^{2}/16}{r_{s}} (y_{t+1} - \hat{y}_{t+1})^{2} \| (\mathbf{I}_{d} - \mathbf{W}_{t} \mathbf{W}_{t}^{\top}) \mathbf{x}_{t+1} \|_{2}^{2} \mathbf{W}_{t}^{\top} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^{\top} \mathbf{W}_{t}.$$

Let

$$c_{t+1}^2 \coloneqq \frac{\eta^2/16}{r_s} \| \boldsymbol{\mathcal{P}}_{t+1} \|_2 = \frac{\eta^2/16}{r_s} \big(y_{t+1} - \hat{y}_{t+1} \big)^2 \| \big(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top \big) \, \boldsymbol{x}_{t+1} \|_2^2 \| \boldsymbol{W}_t^\top \boldsymbol{x}_{t+1} \|_2^2.$$

We define $P_{t+1}\coloneqq \left(I_{r_s}+rac{\eta^2/16}{r_s}\mathcal{P}_{t+1}
ight)^{-1/2}$ and since \mathcal{P}_{t+1} is 1-rank, we have

$$P_{t+1}^2 = I_{r_s} - \frac{\eta^2/16}{r_s} \frac{\mathcal{P}_{t+1}}{1 + c_{t+1}^2}.$$

We let

$$oldsymbol{M}_t \coloneqq oldsymbol{\Theta}^ op oldsymbol{W}_t \ ext{ and } \ \hat{oldsymbol{M}}_{t+1} \coloneqq oldsymbol{\Theta}^ op \hat{oldsymbol{W}}_{t+1}.$$

We have

$$\hat{m{M}}_{t+1} = m{M}_t + rac{\eta/4}{\sqrt{r_s}}m{\Theta}^ op
abla_{ ext{St}}m{L}_{t+1}.$$

By recalling that $G_t = M_t M_t^{\top}$, we have

$$egin{aligned} oldsymbol{G}_{t+1} &= \hat{oldsymbol{M}}_{t+1}^{ op} \hat{oldsymbol{M}}_{t+1}^{ op} + \hat{oldsymbol{M}}_{t+1}^{ op} - oldsymbol{I}_{r_s}) \hat{oldsymbol{M}}_{t+1}^{ op} \ &= oldsymbol{G}_t + rac{\eta/4}{\sqrt{r_s}} oldsymbol{M}_t
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op} oldsymbol{\Theta}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1} oldsymbol{M}_t^{ op} + rac{\eta^2}{16r_s} oldsymbol{\Theta}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op} oldsymbol{M}_{t+1}^{ op} oldsymbol{M}_{t+1}^{ op} oldsymbol{M}_{t+1}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op} oldsymbol{M}_{t+1}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op} oldsymbol{L}_{t+1}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1}^{ op}
abla_{ ext{St}$$

We have

$$\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} = \frac{2}{\|\boldsymbol{\Lambda}\|_{\mathrm{F}}} \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} + \left(\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} - \mathbb{E}_{t} \left[\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \right] \right),$$

Therefore,

$$egin{aligned} oldsymbol{\Theta}^{ op}
abla_{ ext{St}} oldsymbol{L}_{t+1} oldsymbol{M}_t^{ op} &= rac{2}{\|oldsymbol{\Lambda}\|_{ ext{F}}} oldsymbol{\Theta}^{ op} \left(oldsymbol{I}_d - oldsymbol{W}_t oldsymbol{W}_t^{ op}
ight) oldsymbol{\Theta} oldsymbol{\Phi} oldsymbol{W}_t oldsymbol{W}_t^{ op} oldsymbol{\Theta} \ &+ oldsymbol{\Theta}^{ op} \left(
abla_{ ext{St}} oldsymbol{L}_{t+1} - \mathbb{E}_t \left[
abla_{ ext{St}} oldsymbol{L}_{t+1}
ight]
ight) oldsymbol{M}_t^{ op} \ &= rac{2}{\|oldsymbol{\Lambda}\|_{ ext{F}}} \left(oldsymbol{I}_r - oldsymbol{G}_t
ight) oldsymbol{\Lambda} oldsymbol{G}_t + oldsymbol{\Theta}^{ op} \left(
abla_{ ext{St}} oldsymbol{L}_{t+1} - \mathbb{E}_t \left[
abla_{ ext{St}} oldsymbol{L}_{t+1}
ight]
ight) oldsymbol{M}_t^{ op}. \end{aligned}$$

Hence, we have

$$G_{t+1} = G_t + \frac{\eta/2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}} (\mathbf{\Lambda}G_t + G_t\mathbf{\Lambda} - 2G_t\mathbf{\Lambda}G_t)$$

$$+ \frac{\eta/2}{\sqrt{r_s}} \mathrm{Sym} \left(\mathbf{\Theta}^\top \left(\nabla_{\mathrm{S}t} \mathbf{L}_{t+1} - \mathbb{E}_t \left[\nabla_{\mathrm{S}t} \mathbf{L}_{t+1}\right]\right) \mathbf{M}_t^\top\right)$$

$$+ \frac{\eta^2}{16r_s} \mathbf{\Theta}^\top \nabla_{\mathrm{S}t} \mathbf{L}_{t+1} \nabla_{\mathrm{S}t} \mathbf{L}_{t+1}^\top \mathbf{\Theta} - \frac{\eta^2}{16r_s} \frac{\hat{\mathbf{M}}_{t+1} \mathbf{\mathcal{P}}_{t+1} \hat{\mathbf{M}}_{t+1}^\top}{1 + c_{t+1}^2}$$

On the other hand,

$$\begin{split} \hat{M}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \hat{M}_{t+1}^{\top} &= \left(\boldsymbol{M}_{t} + \frac{\eta/4}{\sqrt{r_{s}}} \boldsymbol{\Theta}^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \right) \boldsymbol{\mathcal{P}}_{t+1} \left(\boldsymbol{M}_{t} + \frac{\eta/4}{\sqrt{r_{s}}} \boldsymbol{\Theta}^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \right)^{\top} \\ &= \boldsymbol{M}_{t} \boldsymbol{\mathcal{P}}_{t+1} \boldsymbol{M}_{t}^{\top} + \frac{\eta/2}{\sqrt{r_{s}}} \mathrm{Sym} \left(\boldsymbol{\Theta}^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \boldsymbol{M}_{t}^{\top} \right) + \frac{\eta^{2}}{16r_{s}} \boldsymbol{\Theta}^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \boldsymbol{\Theta}. \end{split}$$

We collect the higher order terms in a single term defined as follows:

$$\begin{split} R_{\text{so}}[\boldsymbol{G}_t] &\coloneqq \frac{\eta^2}{16r_s} \boldsymbol{\Theta}^\top \mathbb{E}_t \left[\nabla_{\text{St}} \boldsymbol{L}_{t+1} \nabla_{\text{St}} \boldsymbol{L}_{t+1}^\top \right] \boldsymbol{\Theta} - \frac{\eta^2}{16r_s} \boldsymbol{M}_t \mathbb{E}_t \left[\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \boldsymbol{M}_t^\top \\ &- \frac{\eta^3}{32r_s^{3/2}} \mathrm{Sym} \left(\boldsymbol{\Theta}^\top \mathbb{E}_t \left[\frac{\nabla_{\text{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \boldsymbol{M}_t^\top \right) \\ &- \frac{\eta^4}{256r_s^2} \boldsymbol{\Theta}^\top \mathbb{E}_t \left[\frac{\nabla_{\text{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\text{St}} \boldsymbol{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] \boldsymbol{\Theta}. \end{split}$$

We collect the noise terms in a single term defined as follows:

$$\frac{\eta/2}{\sqrt{r_s}} \boldsymbol{\nu}_{t+1} \coloneqq \frac{\eta/2}{\sqrt{r_s}} \operatorname{Sym} \left(\boldsymbol{\Theta}^\top \left(\nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} - \mathbb{E}_t \left[\nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} \right] \right) \boldsymbol{M}_t^\top \right) \\ - \frac{\eta^2}{16r_s} \boldsymbol{M}_t \left(\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \right) \boldsymbol{M}_t^\top$$

$$\begin{split} &+\frac{\eta^{2}}{16r_{s}}\boldsymbol{\Theta}^{\top}\left(\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}^{\top}-\mathbb{E}_{t}\left[\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}^{\top}\right]\right)\boldsymbol{\Theta}\\ &-\frac{\eta^{3}}{32r_{s}^{3/2}}\mathrm{Sym}\left(\boldsymbol{\Theta}^{\top}\left(\frac{\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}}{1+c_{t+1}^{2}}-\mathbb{E}_{t}\left[\frac{\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}}{1+c_{t+1}^{2}}\right]\right)\boldsymbol{M}_{t}^{\top}\right)\\ &-\frac{\eta^{4}}{256r_{s}^{2}}\boldsymbol{\Theta}^{\top}\left(\frac{\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}^{\top}}{1+c_{t+1}^{2}}-\mathbb{E}_{t}\left[\frac{\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}\nabla_{\mathsf{St}}\boldsymbol{L}_{t+1}^{\top}}{1+c_{t+1}^{2}}\right]\right)\boldsymbol{\Theta}. \end{split}$$

With these definitions in hand, we have

$$G_{t+1} = G_t + \frac{\eta/2}{\|\mathbf{\Lambda}\|_{\mathrm{F}}\sqrt{r_s}} \left(\mathbf{\Lambda}G_t + G_t\mathbf{\Lambda} - 2G_t\mathbf{\Lambda}G_t\right) + R_{\mathrm{so}}[G_t] + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1}.$$

F.2 Including second-order terms and monotone bounds

For C > 1, we define

$$\Lambda_{\ell_1} := \Lambda - C \|\Lambda\|_{\mathcal{F}} \frac{\eta d}{\sqrt{r_s}} I_r \text{ and } \Lambda_{u_1} := \Lambda + C \|\Lambda\|_{\mathcal{F}} \frac{\eta d}{\sqrt{r_s}} I_r.$$
(F.2)

We recall the definition of effective learning rate $\eta=\frac{\eta/2}{\|\Lambda\|_F\sqrt{r_s}}$. By Proposition 17, we have

$$G_{t+1} \succeq G_t + \eta \left(\mathbf{\Lambda}_{\ell_1} G_t + G_t \mathbf{\Lambda}_{\ell_1} - 2G_t \mathbf{\Lambda} G_t \right) - \frac{C}{2} \eta^2 \|\mathbf{\Lambda}\|_F^2 r_s I_r + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1}, \quad (F.3)$$

$$G_{t+1} \leq G_t + \eta \left(\Lambda_{u_1} G_t + G_t \Lambda_{u_1} - 2G_t \Lambda G_t \right) + \frac{C}{2} \eta^2 \| \Lambda \|_F^2 r_s I_r + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1}.$$
 (F.4)

F.2.1 Heavy tailed case - $\alpha \in [0, 0.5)$

Proposition 12. We consider $\alpha \in [0,0.5)$, $\frac{r_s}{r} \to (0,\infty]$ and $\eta \ll \frac{1}{d \log^4 d} \sqrt{\frac{r_s}{r}}$. We define

$$V_t^- := 2\Lambda^{\frac{1}{2}} G_t \Lambda^{\frac{1}{2}} - \Lambda_{\ell_t}$$
 and $V_t^+ := 2\Lambda^{\frac{1}{2}} G_t \Lambda^{\frac{1}{2}} - \Lambda_{\eta_t}$.

For $d \geq \Omega(1)$, we have

$$\Lambda + \frac{0.1r^{-\alpha}}{\log^4 d} I_r \succ \Lambda_{u_1} \succ \Lambda \succ \Lambda_{\ell_1} \succ \Lambda - \frac{0.1r^{-\alpha}}{\log^4 d} I_r$$
 (F.5)

and

$$V_{t+1}^{-} \succeq V_{t}^{-} \left(I_{r} + \frac{\eta}{1 - 1.1\eta} V_{t}^{-} \right)^{-1} + \eta \Lambda_{\ell_{1}}^{2} - C \eta^{2} \| \mathbf{\Lambda} \|_{F}^{2} r_{s} \mathbf{\Lambda} + \frac{\eta}{\sqrt{r_{s}}} \Lambda^{\frac{1}{2}} \nu_{t+1} \Lambda^{\frac{1}{2}}$$

$$V_{t+1}^{+} \preceq V_{t}^{+} \left(I_{r} + \frac{\eta}{1 + 1.1\eta} V_{t}^{+} \right)^{-1} + \eta \Lambda_{u_{1}}^{2} + C \eta^{2} \| \mathbf{\Lambda} \|_{F}^{2} r_{s} \mathbf{\Lambda} + \frac{\eta}{\sqrt{r_{s}}} \Lambda^{\frac{1}{2}} \nu_{t+1} \Lambda^{\frac{1}{2}}$$
(F.6)

where the bounding iterations are monotone in the sense defined in Proposition 4.

Proof. We first note that since $\|\mathbf{\Lambda}\|_{\mathrm{F}} \approx r^{\frac{1}{2}-\alpha}$ for $\alpha \in [0, 0.5)$, we have

$$\|\mathbf{\Lambda}\|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_s}} \ll \frac{r^{-\alpha}}{\log^4 d}$$

Therefore, (F.5) holds for $d \ge \Omega(1)$, which implies

$$\|\mathbf{V}_{t}^{-}\|_{2} \vee \|\mathbf{V}_{t}^{+}\|_{2} \le 1 + \frac{0.1r^{-\alpha}}{\log^{4} d}, \text{ for all } t \in \mathbb{N}.$$
 (F.7)

Therefore, the monotonicity follows from Proposition 4.

For the remaining part, we introduce the following notation, $K_t := \Lambda^{\frac{1}{2}} G_t \Lambda^{\frac{1}{2}}$. For the lower bound, by (F.3), we have

$$\boldsymbol{K}_{t+1} \succeq \boldsymbol{K}_t + \frac{\eta}{2} \left(\boldsymbol{\Lambda}_{\ell_1}^2 - (2\boldsymbol{K}_t - \boldsymbol{\Lambda}_{\ell_1})^2 \right) - \frac{C}{2} \eta^2 \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2 r_s \boldsymbol{\Lambda} + \frac{\eta/2}{\sqrt{r_s}} \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}^{\frac{1}{2}}.$$

By multiplying both sides with 2 and subtracting Λ_{ℓ_1} from both sides, we have

$$\begin{split} & \boldsymbol{V}_{t+1}^{-} \succeq \boldsymbol{V}_{t}^{-} - \eta(\boldsymbol{V}_{t}^{-})^{2} + \eta\boldsymbol{\Lambda}_{\ell_{1}}^{2} - C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}^{\frac{1}{2}} \\ & \stackrel{(a)}{\succeq} \boldsymbol{V}_{t}^{-} - \frac{\eta}{1 - 1.1\eta}(\boldsymbol{V}_{t}^{-})^{2}\left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta}\boldsymbol{V}_{t}^{-}\right)^{-1} + \eta\boldsymbol{\Lambda}_{\ell_{1}}^{2} \\ & - C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}^{\frac{1}{2}} \\ & = \boldsymbol{V}_{t}^{-}\left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta}\boldsymbol{V}_{t}^{-}\right)^{-1} + \eta\boldsymbol{\Lambda}_{\ell_{1}}^{2} - C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}^{\frac{1}{2}}, \end{split}$$

where we used (F.7) for (a).

For the upper bound, by (F.4), we have

$$\boldsymbol{K}_{t+1} \preceq \boldsymbol{K}_{t} + \frac{\eta}{2} \left(\boldsymbol{\Lambda}_{u_{1}}^{2} - (2\boldsymbol{K}_{t} - \boldsymbol{\Lambda}_{u_{1}})^{2} \right) + \frac{C}{2} \eta^{2} \|\boldsymbol{\Lambda}\|_{F}^{2} r_{s} \boldsymbol{\Lambda} + \frac{\eta/2}{\sqrt{r_{s}}} \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}^{\frac{1}{2}}.$$

By multiplying both sides with 2 and subtracting Λ_{u_1} from both sides, we get

$$\begin{split} \boldsymbol{V}_{t+1}^{+} & \preceq \boldsymbol{V}_{t}^{+} - \eta(\boldsymbol{V}_{t}^{+})^{2} + \eta\boldsymbol{\Lambda}_{u_{1}}^{2} + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}^{\frac{1}{2}} \\ & \stackrel{(b)}{\preceq} \boldsymbol{V}_{t}^{+} - \frac{\eta}{1+1.1\eta}(\boldsymbol{V}_{t}^{+})^{2}\left(\boldsymbol{I}_{r} + \frac{\eta}{1+1.1\eta}\boldsymbol{V}_{t}^{+}\right)^{-1} + \eta\boldsymbol{\Lambda}_{u_{1}}^{2} \\ & + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}^{\frac{1}{2}} \\ & = \boldsymbol{V}_{t}^{+}\left(\boldsymbol{I}_{r} + \frac{\eta}{1+1.1\eta}\boldsymbol{V}_{t}^{+}\right)^{-1} + \eta\boldsymbol{\Lambda}_{u_{1}}^{2} + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1}\boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}}, \end{split}$$

where we used (F.7) for (b).

F.2.2 Light tailed case - $\alpha > 0.5$

We introduce the submatrix notation

$$\boldsymbol{G}_t = : \begin{bmatrix} \boldsymbol{G}_{t,11} & \boldsymbol{G}_{t,12} \\ \boldsymbol{G}_{t,12}^\top & \boldsymbol{G}_{t,22} \end{bmatrix} \quad \boldsymbol{\nu}_t = : \begin{bmatrix} \boldsymbol{\nu}_{t,11} & \boldsymbol{\nu}_{t,12} \\ \boldsymbol{\nu}_{t,12}^\top & \boldsymbol{\nu}_{t,22} \end{bmatrix} \quad \boldsymbol{\Lambda} = : \begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Lambda}_{22} \end{bmatrix} \quad \boldsymbol{\Lambda}_{\ell_1} = : \begin{bmatrix} \boldsymbol{\Lambda}_{\ell_1,11} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Lambda}_{\ell_1,22} \end{bmatrix},$$

where $G_{t,11}, \nu_{t,11}, \Lambda_{11}, \Lambda_{\ell_1,11} \in \mathbb{R}^{r_u \times r_u}$ for $r_u < r$. Similarly, we define the block matrices of Λ_{u_1} as $\Lambda_{u_1,11} \in \mathbb{R}^{r_u \times r_u}$ and $\Lambda_{u_1,22}$. We can write iterations (F.3) and (F.4) for the left-top submatrix as:

$$G_{t+1,11} \succeq G_{t,11} + \eta \left(\mathbf{\Lambda}_{\ell_1,11} G_{t,11} + G_{t,11} \mathbf{\Lambda}_{\ell_1,11} - 2G_{t,11} \mathbf{\Lambda}_{11} G_{t,11} - 2G_{t,12} \mathbf{\Lambda}_{22} G_{t,12}^{\top} \right)$$

$$- \frac{C}{2} \eta^2 \|\mathbf{\Lambda}\|_{F}^2 r_s \mathbf{I}_{r_u} + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1,11}$$
(F.8)

$$G_{t+1,11} \leq G_{t,11} + \eta \left(\Lambda_{u_1,11} G_{t,11} + G_{t,11} \Lambda_{u_1,11} - 2G_{t,11} \Lambda_{11} G_{t,11} - 2G_{t,12} \Lambda_{22} G_{t,12}^{\top} \right)$$

$$+ \frac{C}{2} \eta^2 \| \Lambda \|_{F}^2 r_s I_{r_u} + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1,11}.$$
(F.9)

The following statement is analogous to Proposition 12 in the case $\alpha > 0.5$.

Proposition 13. We consider $\alpha > 0.5$, $r_s \approx 1$, and

$$\eta \ll \frac{1}{d \log^3 d} \frac{1}{r_u^{2+\alpha}}$$
 and $r_u = \lceil \log^{2.5} d \rceil$.

We define $V_t^+\coloneqq 2oldsymbol{\Lambda}_{11}^{rac{1}{2}}oldsymbol{G}_{t,11}oldsymbol{\Lambda}_{11}^{rac{1}{2}}-oldsymbol{\Lambda}_{u_1,11}$ and

$$\boldsymbol{V}_{t}^{-} \coloneqq 2 \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} \boldsymbol{G}_{t,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} - \left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right).$$

For $d \geq \Omega(1)$, we have

$$\Lambda_{11} + \frac{0.1}{r_u^{2+\alpha} \log^3 d} I_{r_u} > \Lambda_{u_1,11} > \Lambda_{11} > \Lambda_{\ell_1,11} > \Lambda_{11} - \frac{0.1}{r_u^{2+\alpha} \log^3 d} I_{r_u}$$
 (F.10)

and

$$\begin{split} \boldsymbol{V}_{t+1}^{-} &\succeq \boldsymbol{V}_{t}^{-} \Big(\boldsymbol{I}_{r_{u}} + \frac{\eta}{1-1.1\eta} \boldsymbol{V}_{t}^{-} \Big)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{2} - C \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \Big(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \Big)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} \\ &+ \frac{\eta}{\sqrt{r_{s}}} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} \\ &\boldsymbol{V}_{t+1}^{+} \preceq \boldsymbol{V}_{t}^{+} \left(\boldsymbol{I}_{r_{u}} + \frac{\eta}{1+1.1\eta} \boldsymbol{V}_{t}^{+} \right)^{-1} + \eta \boldsymbol{\Lambda}_{u_{1},11}^{2} + C \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \boldsymbol{\Lambda}_{11} + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}}. \end{split}$$

where the bounding iterations are monotone in the sense defined in Proposition 4.

Proof. We first note that since $r_s \approx 1$ and $\|\mathbf{\Lambda}\|_{\mathrm{F}} \approx 1$ for $\alpha > 0.5$, we have

$$\|\mathbf{\Lambda}\|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_s}} \ll \frac{1}{r_u^{2+\alpha} \log^3 d}.$$

Therefore, (F.10) holds for $d \ge \Omega(1)$, which implies

$$\|\mathbf{V}_{t}^{-}\|_{2} \vee \|\mathbf{V}_{t}^{+}\|_{2} \le 1 + \frac{0.1}{r_{u}^{2+\alpha} \log^{3} d} + \frac{1}{(r_{u}+1)^{\alpha}}, \text{ for all } t \in \mathbb{N}.$$
 (F.11)

For the remaining part, we introduce the following notation,

$$\boldsymbol{K}_t^- \coloneqq \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_u + 1)^\alpha} \boldsymbol{I}_{r_u}\right)^{\frac{1}{2}} \boldsymbol{G}_{t,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_u + 1)^\alpha} \boldsymbol{I}_{r_u}\right)^{\frac{1}{2}} \ \text{ and } \ \boldsymbol{K}_t^+ \coloneqq \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{G}_{t,11} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}}.$$

For the upper bound, since $\Lambda_{22} > 0$, by (F.9) we have

$$\boldsymbol{K}_{t+1}^{+} \preceq \boldsymbol{K}_{t}^{+} + \frac{\eta}{2} \left(\boldsymbol{\Lambda}_{u_{1},11}^{2} - (2\boldsymbol{K}_{t}^{+} - \boldsymbol{\Lambda}_{u_{1},11})^{2} \right) + \frac{C}{2} \eta^{2} \|\boldsymbol{\Lambda}\|_{F}^{2} r_{s} \boldsymbol{\Lambda}_{11} + \frac{\eta/2}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}}.$$

By multiplying both sides with 2 and subtracting $\Lambda_{u_1,11}$ from both sides, we get

$$\begin{split} \boldsymbol{V}_{t+1}^{+} & \leq \boldsymbol{V}_{t}^{+} - \eta(\boldsymbol{V}_{t}^{+})^{2} + \eta\boldsymbol{\Lambda}_{u_{1},11}^{2} + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda}_{11} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1,11}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \\ & \stackrel{(a)}{\leq} \boldsymbol{V}_{t}^{+} - \frac{\eta}{1+1.1\eta}(\boldsymbol{V}_{t}^{+})^{2}\left(\boldsymbol{I}_{r_{u}} + \frac{\eta}{1+1.1\eta}\boldsymbol{V}_{t}^{+}\right)^{-1} + \eta\boldsymbol{\Lambda}_{u_{1},11}^{2} \\ & + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda}_{11} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1,11}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \\ & = \boldsymbol{V}_{t}^{+}\left(\boldsymbol{I}_{r_{u}} + \frac{\eta}{1+1.1\eta}\boldsymbol{V}_{t}^{+}\right)^{-1} + \eta\boldsymbol{\Lambda}_{u_{1},11}^{2} + C\eta^{2}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2}r_{s}\boldsymbol{\Lambda}_{11} + \frac{\eta}{\sqrt{r_{s}}}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\boldsymbol{\nu}_{t+1,11}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}, \end{split}$$

where we used (F.11) for (a).

For the lower bound, we first observe that $G_{t,11}(I_{r_u}-G_{t,11})-G_{t,12}G_{t,12}^{\top}\succeq 0$ since it corresponds to the left-top submatrix of $G_t(I_r-G_t)$. Therefore, by (F.8)

$$\begin{aligned} G_{t+1,11} &\succeq G_{t,11} + \eta \Big(\mathbf{\Lambda}_{\ell_1,11} G_{t,11} + G_{t,11} \mathbf{\Lambda}_{\ell_1,11} - 2G_{t,11} \mathbf{\Lambda}_{11} G_{t,11} - 2G_{t,12} \mathbf{\Lambda}_{22} G_{t,12}^{\top} \Big) \\ &- \frac{C}{2} \eta^2 \| \mathbf{\Lambda} \|_{\mathrm{F}}^2 r_s \mathbf{I}_{r_u} + \frac{\eta/2}{\sqrt{r_s}} \boldsymbol{\nu}_{t+1,11} - \frac{2\eta}{(r_u + 1)^{\alpha}} \left(G_{t,11} (\mathbf{I}_{r_u} - G_{t,11}) - G_{t,12} G_{t,12}^{\top} \right) \\ &\succeq G_{t,11} + \eta \Big(\left(\mathbf{\Lambda}_{\ell_1,11} - \frac{1}{(r_u + 1)^{\alpha}} \mathbf{I}_{r_u} \right) G_{t,11} + G_{t,11} \left(\mathbf{\Lambda}_{\ell_1,11} - \frac{1}{(r_u + 1)^{\alpha}} \mathbf{I}_{r_u} \right) \Big) \\ &- 2\eta G_{t,11} \Big(\mathbf{\Lambda}_{11} - \frac{1}{(r_u + 1)^{\alpha}} \mathbf{I}_{r_u} \Big) G_{t,11} - \frac{C}{2} \eta^2 \| \mathbf{\Lambda} \|_{\mathrm{F}}^2 r_s \mathbf{I}_{r_u} + \frac{\eta/2}{\sqrt{r_s}} \boldsymbol{\nu}_{t+1,11}, \end{aligned}$$

where (b) follows by $\Lambda_{22} \preceq \frac{1}{(r_u+1)^{\alpha}} I_{r-r_u}$. Therefore, we have

$$\begin{split} \boldsymbol{K}_{t+1}^{-} \succeq \boldsymbol{K}_{t}^{-} + \frac{\eta}{2} \left(\left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{2} - \left(2\boldsymbol{K}_{t}^{-} - \left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right) \right)^{2} \right) \\ - \frac{C}{2} \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right) \\ + \frac{\eta/2}{\sqrt{r_{s}}} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}} \right)^{\frac{1}{2}}. \end{split}$$

By multiplying both sides with 2 and subtracting $\Lambda_{\ell_1,11}$ from both sides, we get

$$\begin{split} & \boldsymbol{V}_{t+1}^{-} \succeq \boldsymbol{V}_{t}^{-} - \eta(\boldsymbol{V}_{t}^{-})^{2} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{2} - C \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \\ & \quad + \frac{\eta}{\sqrt{r_{s}}} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \\ & \quad + C \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \\ & \quad + \frac{\eta}{\sqrt{r_{s}}} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \\ & \quad = \boldsymbol{V}_{t}^{-} \left(\boldsymbol{I}_{r_{u}} + \frac{\eta}{1-1.1\eta} \boldsymbol{V}_{t}^{-}\right)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{2} \\ & \quad - C \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right) \\ & \quad + \frac{\eta}{\sqrt{r_{s}}} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \left(\boldsymbol{\Lambda}_{11} - \frac{1}{(r_{u}+1)^{\alpha}} \boldsymbol{I}_{r_{u}}\right)^{\frac{1}{2}}, \end{split}$$

where we used (F.11) for (c). The monotonicity of the update follows the same argument in the heavy-tailed case.

F.3 Definitions and bounding systems

To avoid repetition in the derivations, we introduce the following unified notation:

$$\mathsf{rk} \in \{r, r_u\}, \quad \mathsf{G}_t \in \{G_t, G_{t+1}\}, \quad \mathsf{v}_t \in \{\nu_t, \nu_{t+1}\}.$$

where each variable will take its first value in the heavy-tailed case and its second value in the light-tailed case. To avoid repetition in the following sections, we make the following simplifications by slight abuse of notation:

$$\boldsymbol{\Lambda}_{\ell_1} \leftarrow \begin{cases} \boldsymbol{\Lambda}_{\ell_1}, & \alpha \in [0, 0.5) \\ \boldsymbol{\Lambda}_{\ell_1, 11} - \frac{1}{(r_u + 1)^{\alpha}} \boldsymbol{I}_{r_u}, & \alpha > 0.5 \end{cases} \text{ and } \boldsymbol{\Lambda}_{\ell_2} \leftarrow \begin{cases} \boldsymbol{\Lambda}, & \alpha \in [0, 0.5) \\ \boldsymbol{\Lambda}_{11} - \frac{1}{(r_u + 1)^{\alpha}} \boldsymbol{I}_{r_u}, & \alpha > 0.5 \end{cases}$$

and

$$\mathbf{\Lambda}_{u_1} \leftarrow \begin{cases} \mathbf{\Lambda}_{u_1}, & \alpha \in [0, 0.5) \\ \mathbf{\Lambda}_{u_1, 11}, & \alpha > 0.5 \end{cases} \text{ and } \mathbf{\Lambda}_{u_2} \leftarrow \begin{cases} \mathbf{\Lambda}, & \alpha \in [0, 0.5) \\ \mathbf{\Lambda}_{11}, & \alpha > 0.5. \end{cases}$$

The dimension of each block is $r_u < r$ for $\alpha > 0.5$ and r for $\alpha \in [0,0.5)$, from which readers can distinguish the light tailed case from the heavy tailed case. Throughout the proof, we will also use constants $\kappa_d \in o_d(1)$ and $\tilde{C} \in O(1)$ that will be specified later. Moreover, we make the following definitions:

- Noise sequence. For $\underline{\nu}_0 = 0$, we define the noise sequence $\underline{\nu}_{t+1} \coloneqq \underline{\nu}_t + \frac{\eta/2}{\sqrt{r_s}} \nu_{t+1}$.
- Reference sequence. For $T_0=rac{\kappa_d r_s}{d} I_{
 m rk},$ we define the reference sequence

$$T_{t+1} = T_t + 2(1 - 2\kappa_d)\eta \left(\Lambda_{\ell_1} T_t - \frac{3\kappa_d + 1}{\kappa_d (1 - 2\kappa_d)} \Lambda_{u_2} T_t^2 \right).$$

• **Bounding systems.** We define the lower and upper bounding recursions as

$$\underline{\boldsymbol{V}}_{t+1} = \underline{\boldsymbol{V}}_{t} \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta(1+2\kappa_{d})}{1-1.2\eta} \underline{\boldsymbol{V}}_{t} \right)^{-1} + \frac{\eta(1+2\kappa_{d})}{1-1.2\eta} \left(\frac{\boldsymbol{\Lambda}_{\ell_{1}}^{2}}{(1+2\kappa_{d})^{2}} - \tilde{C}\eta \|\boldsymbol{\Lambda}\|_{F}^{2} r_{s} \boldsymbol{\Lambda}_{\ell_{1}} \right)$$
(F.12)

$$\bar{\mathbf{V}}_{t+1} = \bar{\mathbf{V}}_t \left(\mathbf{I}_{\mathsf{rk}} + \frac{\eta(1 - 2\kappa_d)}{1 + 1.2\eta} \bar{\mathbf{V}}_t \right)^{-1} + \frac{\eta(1 - 2\kappa_d)}{1 + 1.2\eta} \left(\frac{\mathbf{\Lambda}_{u_1}^2}{(1 - 2\kappa_d)^2} + \tilde{C}\eta \|\mathbf{\Lambda}\|_{\mathrm{F}}^2 r_s \mathbf{\Lambda}_{u_1} \right).$$
(F.13)

where the iterates $\{\underline{V}_t\}_{t\in\mathbb{N}}$ and $\{\bar{V}_t\}_{t\in\mathbb{N}}$ are functions of the bounding sequences $\{\underline{G}_t\}_{t\in\mathbb{N}}$ and $\{\bar{G}_t\}_{t\in\mathbb{N}}$ as following:

$$\underline{\underline{K}}_t := \Lambda_{\ell_2}^{\frac{1}{2}} \underline{\underline{G}}_t \Lambda_{\ell_2}^{\frac{1}{2}} \text{ and } \underline{\underline{V}}_t = 2\underline{\underline{K}}_t - \frac{\Lambda_{\ell_1}}{1 + 2\kappa_d} \text{ and } \underline{\underline{G}}_0 \preceq \mathsf{G}_0 - \underline{T}_0,$$

$$ar{m{K}}_t \coloneqq m{\Lambda}_{u_2}^{rac{1}{2}} ar{m{G}}_t m{\Lambda}_{u_2}^{rac{1}{2}} \;\; ext{and} \;\; ar{m{V}}_t = 2ar{m{K}}_t - rac{m{\Lambda}_{u_1}}{1-2m{\kappa}_d} \;\;\; ext{and} \;\; ar{m{G}}_0 \succeq m{\mathsf{G}}_0 + m{T}_0.$$

• Stopping times. We define a sequence of events $\{\mathcal{E}_t\}_{t\geq 0}$

$$\mathcal{E}_{t} \coloneqq \begin{cases} \left\{ -\kappa_{d}r^{\frac{-\alpha}{2}}\boldsymbol{T}_{t} \leq \underline{\boldsymbol{\nu}}_{t} \leq \kappa_{d}r^{\frac{-\alpha}{2}}\boldsymbol{T}_{t} \right\} \cap \left\{ -\kappa_{d}^{2}r^{-\alpha}\boldsymbol{T}_{t} \leq \boldsymbol{\Lambda}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_{t}\boldsymbol{\Lambda}^{\frac{1}{2}} \leq \kappa_{d}^{2}r^{-\alpha}\boldsymbol{T}_{t} \right\}, \quad \alpha \in [0, 0.5) \\ \left\{ -\kappa_{d}r_{u}^{\frac{-\alpha}{2}}\boldsymbol{T}_{t} \leq \underline{\boldsymbol{\nu}}_{t} \leq \kappa_{d}r_{u}^{\frac{-\alpha}{2}}\boldsymbol{T}_{t} \right\} \cap \left\{ \frac{-\kappa_{d}^{2}}{4}r_{u}^{-\alpha}\boldsymbol{T}_{t} \leq \boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_{t}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \leq \frac{\kappa_{d}^{2}}{4}r_{u}^{-\alpha}\boldsymbol{T}_{t} \right\}, \quad \alpha > 0.5. \end{cases}$$

We define the stopping times

 $\mathcal{T}_{\text{noise}}(\omega) \coloneqq \inf \{ t \ge 0 \mid \omega \not\in \mathcal{E}_t \} \wedge d^3 \text{ and } \mathcal{T}_{\text{bounded}} \coloneqq \inf \{ t \ge 0 \mid \|\underline{\boldsymbol{G}}_t\|_2 \vee \|\bar{\boldsymbol{G}}_t\|_2 \ge 1.5 \},$ and

$$\mathcal{T}_{\text{bad}} := \mathcal{T}_{\text{noise}} \wedge \mathcal{T}_{\text{bounded}} \wedge \{t \ge 0 : ||T_t||_2 > 1.2\kappa_d\}. \tag{F.14}$$

The main result of this section is the following:

Proposition 14. Let κ_d satisfy (F.15) and consider d large enough so that $\kappa_d \leq \frac{1}{50}$. Under the learning rate conditions considered in Propositions 12 and 13, we have for $\tilde{C} > 1 \vee \Omega(1)$

$$\underline{G}_{t \wedge \mathcal{T}_{bad}} + T_{t \wedge \mathcal{T}_{bad}} + \underline{
u}_{t \wedge \mathcal{T}_{bad}} \preceq \mathsf{G}_{t \wedge \mathcal{T}_{bad}} \preceq \bar{G}_{t \wedge \mathcal{T}_{bad}} - T_{t \wedge \mathcal{T}_{bad}} + \underline{
u}_{t \wedge \mathcal{T}_{bad}}$$

Before starting the proof, we provide an auxiliary statement.

Lemma 4. We consider the learning rate conditions considered in Propositions 12 and 13 with

$$\kappa_d \ll \begin{cases} \frac{1}{\log d}, & \alpha \in [0, 0.5) \\ \frac{r_u^{-1}}{\log d}, & \alpha > 0.5. \end{cases}$$
(F.15)

The event \mathcal{E}_t implies for $d \geq \Omega(1)$ and $t \leq \mathcal{T}_{bounded} \wedge \{t : ||T_t||_2 > 1.2\kappa_d\}$ that

1.
$$-3\kappa_d \Lambda_{\ell_1}^{\frac{1}{2}} T_t \Lambda_{\ell_1}^{\frac{1}{2}} \preceq \Lambda_{\ell_1} \underline{\nu}_t + \underline{\nu}_t \Lambda_{\ell_1}$$

2.
$$\Lambda_{u_1} \nu_{\star} + \nu_{\star} \Lambda_{u_1} \prec 3 \kappa_d \Lambda_{u_1}^{\frac{1}{2}} T_t \Lambda_{u_1}^{\frac{1}{2}}$$

3.
$$\left(\Lambda_{\ell_2}^{\frac{1}{2}}\underline{\nu}_t\Lambda_{\ell_2}^{\frac{1}{2}}\right)^2 \leq \frac{\kappa_d^2}{4}\Lambda_{\ell_1}T_t\Lambda_{\ell_1}$$

4.
$$\left(\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \underline{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\right)^2 \preceq \frac{\kappa_d^2}{4} \boldsymbol{\Lambda}_{u_1} \boldsymbol{T}_t \boldsymbol{\Lambda}_{u_1}$$

Proof. For notational convenience, we define $\widetilde{\nu}_t := T_t^{\frac{-1}{2}} \underline{\nu}_t T_t^{\frac{-1}{2}}$. We initially observe that \mathcal{E}_t implies for $\alpha \in [0.0.5)$ that

$$\begin{split} & \boldsymbol{\Lambda} \widetilde{\boldsymbol{\nu}}_t + \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda} \preceq \kappa_d^2 \boldsymbol{\Lambda} + \frac{1}{\kappa_d^2} \boldsymbol{\Lambda}^{\frac{-1}{2}} \left(\boldsymbol{\Lambda}^{\frac{1}{2}} \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}^{\frac{1}{2}} \right)^2 \boldsymbol{\Lambda}^{\frac{-1}{2}} \preceq \kappa_d^2 \boldsymbol{\Lambda} + \kappa_d^2 r^{-\alpha} \boldsymbol{I}_r \\ & \boldsymbol{\Lambda} \widetilde{\boldsymbol{\nu}}_t + \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda} \succeq -\kappa_d^2 \boldsymbol{\Lambda} - \frac{1}{\kappa_d^2} \boldsymbol{\Lambda}^{\frac{-1}{2}} \left(\boldsymbol{\Lambda}^{\frac{1}{2}} \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}^{\frac{1}{2}} \right)^2 \boldsymbol{\Lambda}^{\frac{-1}{2}} \succeq -\kappa_d^2 \boldsymbol{\Lambda} - \kappa_d^2 r^{-\alpha} \boldsymbol{I}_r. \end{split}$$

For $\alpha > 0.5$, these bounds become

$$\begin{split} & \boldsymbol{\Lambda}_{11} \widetilde{\boldsymbol{\nu}}_t + \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{11} \preceq \frac{\kappa_d^2}{4} \boldsymbol{\Lambda}_{11} + \frac{4}{\kappa_d^2} \boldsymbol{\Lambda}_{11}^{\frac{-1}{2}} \left(\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \right)^2 \boldsymbol{\Lambda}_{11}^{\frac{-1}{2}} \preceq \frac{\kappa_d^2}{4} \boldsymbol{\Lambda}_{11} + \frac{\kappa_d^2}{4} r_u^{-\alpha} \boldsymbol{I}_{r_u} \\ & \boldsymbol{\Lambda}_{11} \widetilde{\boldsymbol{\nu}}_t + \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{11} \succeq -\frac{\kappa_d^2}{4} r_u^{-1} \boldsymbol{\Lambda}_{11} - \frac{4}{\kappa_d^2} \boldsymbol{\Lambda}_{11}^{\frac{-1}{2}} \left(\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \widetilde{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \right)^2 \boldsymbol{\Lambda}_{11}^{\frac{-1}{2}} \succeq -\frac{\kappa_d^2}{4} r_u^{-1} \boldsymbol{\Lambda}_{11} - \frac{\kappa_d^2}{4} r_u^{-\alpha} \boldsymbol{I}_{r_u}. \end{split}$$

In the following, we will use these bounds. For the first item, we have

$$\begin{split} \boldsymbol{\Lambda}_{\ell_{1}} \widetilde{\boldsymbol{\nu}}_{t} + \widetilde{\boldsymbol{\nu}}_{t} \boldsymbol{\Lambda}_{\ell_{1}} \succeq \begin{cases} -\kappa_{d} \left(\kappa_{d} \boldsymbol{\Lambda} + \kappa_{d} r^{-\alpha} \boldsymbol{I}_{r} + r^{\frac{-\alpha}{2}} \boldsymbol{C} \| \boldsymbol{\Lambda} \|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_{s}}} \boldsymbol{I}_{r} \right), & \alpha \in [0, 0.5) \\ -\kappa_{d} \left(\frac{\kappa_{d}}{4} \boldsymbol{\Lambda}_{11} + \frac{\kappa_{d}}{4} r_{u}^{-\alpha} \boldsymbol{I}_{r_{u}} + r_{u}^{\frac{-\alpha}{2}} \boldsymbol{C} \| \boldsymbol{\Lambda} \|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_{s}}} \boldsymbol{I}_{r_{u}} \right), & \alpha > 0.5 \end{cases} \\ \succeq -3\kappa_{d} \boldsymbol{\Lambda}_{\ell_{1}}. \end{split}$$

For the second item, we have

$$\Lambda_{u_1} \widetilde{\boldsymbol{\nu}}_t + \widetilde{\boldsymbol{\nu}}_t \Lambda_{u_1} \preceq \begin{cases}
\kappa_d \left(\kappa_d \boldsymbol{\Lambda} + \kappa_d r^{-\alpha} \boldsymbol{I}_r + r^{\frac{-\alpha}{2}} C \|\boldsymbol{\Lambda}\|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_s}} \boldsymbol{I}_r \right), & \alpha \in [0, 0.5) \\
\kappa_d \left(\frac{\kappa_d}{4} \Lambda_{11} + \frac{\kappa_d}{4} r_u^{-\alpha} \boldsymbol{I}_{r_u} + r_u^{\frac{-\alpha}{2}} C \|\boldsymbol{\Lambda}\|_{\mathrm{F}} \frac{\eta d}{\sqrt{r_s}} \boldsymbol{I}_{r_u} \right), & \alpha > 0.5
\end{cases}$$

$$\preceq 3\kappa_d \Lambda_{u_1}.$$

For the third item, we immediately observe that $(\Lambda_{\ell_2}^{\frac{1}{2}}\underline{\nu}_t\Lambda_{\ell_2}^{\frac{1}{2}})^2 \preceq \Lambda_{\ell_2}^{\frac{1}{2}}\underline{\nu}_t^2\Lambda_{\ell_2}^{\frac{1}{2}}$. Therefore,

$$\boldsymbol{T}_{t}^{-\frac{1}{2}}\boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}}\boldsymbol{\nu}_{t}^{2}\boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}}\boldsymbol{T}_{t}^{-\frac{1}{2}} \preceq 1.2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}}\tilde{\boldsymbol{\nu}}_{t}^{2}\boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}} \preceq 1.2\kappa_{d}^{3}\mathsf{rk}^{-\alpha}\boldsymbol{\Lambda}_{\ell_{2}} \preceq \frac{\kappa_{d}^{2}}{4}\boldsymbol{\Lambda}_{\ell_{1}}^{2}.$$

For the fourth item, we observe that

$$\left(\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\right)^2 = \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{\Lambda}_{u_2}\underline{\boldsymbol{\nu}}_t\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \preceq \left(1 + \frac{0.1\mathsf{rk}^{-\alpha}}{\log^3 d}\right)\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_t^2\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}.$$

Therefore

$$\left(1 + \frac{0.1 \mathrm{rk}^{-\alpha}}{\log^3 d}\right) \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \underline{\boldsymbol{\nu}}_t^2 \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \boldsymbol{T}_t^{\frac{-1}{2}} \preceq 1.25 \kappa_d \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \widetilde{\boldsymbol{\nu}}_t^2 \boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \preceq 1.25 \kappa_d^3 \mathrm{rk}^{-\alpha} \boldsymbol{\Lambda}_{u_2} \preceq \frac{\kappa_d^2}{4} \boldsymbol{\Lambda}_{u_1}^2.$$

F.3.1 Proof of Proposition 14

Proof. For the proof, we introduce the following notations

$$\underline{\boldsymbol{\zeta}}_t \coloneqq 2\boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}} \ \ \text{and} \ \ \bar{\boldsymbol{\zeta}}_t \coloneqq 2\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}} \ \ \text{and} \ \ \underline{\boldsymbol{B}}_t \coloneqq 2\boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}}\boldsymbol{T}_t\boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}} \ \ \text{and} \ \ \bar{\boldsymbol{B}}_t \coloneqq 2\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}\boldsymbol{T}_t\boldsymbol{\Lambda}_{u_2}^{\frac{1}{2}}.$$

Using this notation, we obtain:

$$\underline{\boldsymbol{B}}_{t+1} \preceq \underline{\boldsymbol{B}}_{t} + \frac{2(1 - 2\kappa_{d})\eta}{1 - 1.1\eta} \left(\boldsymbol{\Lambda}_{\ell_{1}} \underline{\boldsymbol{B}}_{t} - \frac{1.5\kappa_{d} + 0.5}{\kappa_{d}(1 - 2\kappa_{d})} \underline{\boldsymbol{B}}_{t}^{2} \right) + 10\eta^{2} \kappa_{d} \boldsymbol{\Lambda}_{\ell_{1}}$$
 (F.16)

$$\bar{B}_{t+1} \leq \bar{B}_t + \frac{2(1 - 2\kappa_d)\eta}{1 + 1.1\eta} \left(\Lambda_{u_1} \bar{B}_t - \frac{1.5\kappa_d + 0.5}{\kappa_d (1 - 2\kappa_d)} \bar{B}_t^2 \right) + 10\eta^2 \kappa_d \Lambda_{u_1}$$
 (F.17)

Before proceeding with the proof, we observe that the following inequalities hold:

$$\|\mathbf{\Lambda}_{\ell_2}^{-1}\mathbf{\Lambda}_{\ell_1}\|_2 \leq 1, \ \|\mathbf{\Lambda}_{\ell_1}^{-1}\mathbf{\Lambda}_{\ell_2}\|_2 \leq \frac{1}{1 - \frac{0.1}{\log^4 d}} \ \text{ and } \ \|\mathbf{\Lambda}_{u_2}^{-1}\mathbf{\Lambda}_{u_1}\|_2 \leq 1, \ \|\mathbf{\Lambda}_{u_1}^{-1}\mathbf{\Lambda}_{u_2}\|_2 \leq \frac{1}{1 - \frac{0.1}{\log^4 d}}.$$

These bounds will be used in the following whenever we apply Propositions 29 and 30, without explicitly restating them each time. We will establish the upper and lower bounds simultaneously for $rk \in \{r, r_u\}$.

Upper bound proof: We will use proof by induction. Specifically, we will show that for $t < T_{bad}$,

$$V_t^+ \leq \bar{V}_t + \bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1 - 2\kappa_d} \Rightarrow V_{t+1}^+ \leq \bar{V}_{t+1} + \bar{\zeta}_{t+1} - \bar{B}_{t+1} + \frac{2\kappa_d \Lambda_{u_1}}{1 - 2\kappa_d}.$$
 (F.18)

Since the base case holds at t=0 and $\mathcal{T}_{bad}>0$, it remains to prove (F.18). By (F.6), we have

$$\begin{split} & \boldsymbol{V}_{t+1}^{+} \preceq \left(\bar{\boldsymbol{V}}_{t} + \bar{\boldsymbol{\zeta}}_{t} - \bar{\boldsymbol{B}}_{t} + \frac{2\kappa_{d}\boldsymbol{\Lambda}_{u_{1}}}{1 - 2\kappa_{d}} \right) \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 + 1.1\eta} (\bar{\boldsymbol{V}}_{t} + \bar{\boldsymbol{\zeta}}_{t} - \bar{\boldsymbol{B}}_{t} + \frac{2\kappa_{d}\boldsymbol{\Lambda}_{u_{1}}}{1 - 2\kappa_{d}}) \right)^{-1} \\ & + \eta \boldsymbol{\Lambda}_{u_{1}}^{2} + C\eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \boldsymbol{\Lambda}_{u_{2}} + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}} \\ & = \bar{\boldsymbol{V}}_{t} \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 + 1.1\eta} \bar{\boldsymbol{V}}_{t} \right)^{-1} + \eta \boldsymbol{\Lambda}_{u_{1}}^{2} + C\eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \boldsymbol{\Lambda}_{u_{2}} + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}_{u_{2}}^{\frac{1}{2}} \\ & + \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 + 1.1\eta} \bar{\boldsymbol{V}}_{t} \right)^{-1} (\bar{\boldsymbol{\zeta}}_{t} - \bar{\boldsymbol{B}}_{t} + \frac{2\kappa_{d}\boldsymbol{\Lambda}_{u_{1}}}{1 - 2\kappa_{d}}) \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 + 1.1\eta} (\bar{\boldsymbol{V}}_{t} + \bar{\boldsymbol{\zeta}}_{t} - \bar{\boldsymbol{B}}_{t} + \frac{2\kappa_{d}\boldsymbol{\Lambda}_{u_{1}}}{1 - 2\kappa_{d}}) \right)^{-1} . \end{split}$$

$$(F.19)$$

By using Proposition 29, we have for $t < T_{\text{bad}}$

$$\begin{split} & \left(I_{\mathsf{rk}} + \frac{\eta}{1+1.1\eta} \bar{V}_t \right)^{-1} (\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d}) \left(I_{\mathsf{rk}} + \frac{\eta}{1+1.1\eta} (\bar{V}_t + \bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d}) \right)^{-1} \\ & \preceq \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) - \frac{\eta}{1+1.1\eta} \bar{V}_t \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) \\ & - \frac{\eta}{1+1.1\eta} \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) \bar{V}_t - \frac{\eta}{1+1.1\eta} \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right)^2 + \frac{\eta^2 \kappa_d^2}{(1+1.1\eta)^2} \bar{V}_t^4 \\ & + \frac{2\eta^2/\kappa_d^2}{(1+1.1\eta)^2} \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right)^2 + \frac{\eta^2 \kappa_d^2}{(1+1.1\eta)^2} \bar{V}_t \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right)^2 \bar{V}_t \\ & + \frac{\eta^2}{(1+1.1\eta)^2} \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) \bar{V}_t \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) + \eta^3 \tilde{C}_1 \Lambda_{u_1} \\ & + \frac{\eta^2}{(1+1.1\eta)^2} \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right)^3 + \frac{\eta^2}{(1+1.1\eta)^2} \bar{V}_t \left(\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \right) \bar{V}_t \end{aligned}$$

for some $\tilde{C}_1 = O(1)$. We have the following: First:

$$\begin{split} \kappa_d^2 \bar{\boldsymbol{V}}_t^4 + \kappa_d^2 \bar{\boldsymbol{V}}_t \left(\bar{\boldsymbol{\zeta}}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right)^2 \bar{\boldsymbol{V}}_t + \bar{\boldsymbol{V}}_t \left(\bar{\boldsymbol{\zeta}}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right) \bar{\boldsymbol{V}}_t \overset{(a)}{\leq} 4\kappa_d \bar{\boldsymbol{V}}_t^2 \\ \overset{(b)}{\leq} \frac{1}{15} \frac{1 + 1.1 \eta}{1 + 1.2 \eta} \bar{\boldsymbol{V}}_t^2, \end{split}$$

where (a) follows by $\|\bar{\pmb{V}}_t\|_2 \le 5$ and $\left\|\bar{\pmb{\zeta}}_t - \bar{\pmb{B}}_t + \frac{2\kappa_d \pmb{\Lambda}_{u_1}}{1-2\kappa_d}\right\|_2 \le 2.5\kappa_d$, and (b) follows by $\kappa_d \le \frac{1}{50}$. Second,

$$\begin{split} &\frac{2}{\kappa_d^2} \left(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right)^2 + \left(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right) \bar{\boldsymbol{V}}_t \left(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right) \\ &+ \left(\bar{\zeta}_t + \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2c_d} \right)^3 \overset{(c)}{\leq} \frac{3}{\kappa_d^2} \left(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \right)^2 \overset{g}{\leq} \frac{9}{\kappa_d^2} \left(\bar{\zeta}_t^2 + \bar{\boldsymbol{B}}_t^2 \right) + \frac{36}{(1 - 2\kappa_d)^2} \boldsymbol{\Lambda}_{u_1}^2, \end{split}$$

where (c) follows by $\|\bar{V}_t\|_2 \le 5$ and $\|\bar{\zeta}_t - \bar{B}_t + \frac{2\kappa_d\Lambda_{u_1}}{1-2\kappa_d}\|_2 \le 2.5\kappa_d$. Third:

$$\begin{split} &-\bar{\boldsymbol{V}}_t \bigg(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \bigg) - \bigg(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \bigg) \, \bar{\boldsymbol{V}}_t - \bigg(\bar{\zeta}_t - \bar{\boldsymbol{B}}_t + \frac{2\kappa_d \boldsymbol{\Lambda}_{u_1}}{1 - 2\kappa_d} \bigg)^2 \\ &+ \frac{9\eta/\kappa_d^2}{1 + 1.1\eta} \, \big(\bar{\zeta}_t^2 + \bar{\boldsymbol{B}}_t^2 \big) \end{split}$$

$$\leq -2(\bar{K}_{t}\bar{\zeta}_{t} + \bar{\zeta}_{t}\bar{K}_{t}) + 2(\bar{K}_{t}\bar{B}_{t} + \bar{B}_{t}\bar{K}_{t}) - \frac{4\kappa_{d}}{1 - 2\kappa_{d}}(\bar{K}_{t}\Lambda_{u_{1}} + \Lambda_{u_{1}}\bar{K}_{t}) + 3(\bar{\zeta}_{t}^{2} + \bar{B}_{t}^{2}) \\
+ (\Lambda_{u_{1}}\bar{\zeta}_{t} + \bar{\zeta}_{t}\Lambda_{u_{1}}) - (\Lambda_{u_{1}}\bar{B}_{t} + \bar{B}_{t}\Lambda_{u_{1}}) + \frac{4\kappa_{d}(1 - \kappa_{d})}{(1 - 2\kappa_{d})^{2}}\Lambda_{u_{1}}^{2} \\
\leq 8\kappa_{d}\bar{K}_{t}^{2} - \frac{4\kappa_{d}}{1 - 2\kappa_{d}}(\bar{K}_{t}\Lambda_{u_{1}} + \Lambda_{u_{1}}\bar{K}_{t}) + \frac{4\kappa_{d}(1 - \kappa_{d})}{(1 - 2\kappa_{d})^{2}}\Lambda_{u_{1}}^{2} \\
- (2 - 4\kappa_{d})\Lambda_{u_{1}}\bar{B}_{t} + (3 + \frac{1}{\kappa_{d}})\bar{B}_{t}^{2} \\
= 2\kappa_{d}\bar{V}_{t}^{2} + \frac{2\kappa_{d}}{1 - 2\kappa_{d}}\Lambda_{u_{1}}^{2} - (2 - 4\kappa_{d})\Lambda_{u_{1}}\bar{B}_{t} + (3 + \frac{1}{\kappa_{d}})\bar{B}_{t}^{2},$$

where we used Proposition 22, and the second and fourth items in Lemma 4 in (d). Therefore, we have

$$\begin{split} (\textbf{F.19}) & \preceq \bar{V}_t \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1+1.1\eta} \bar{V}_t \right)^{-1} + \frac{2\kappa_d \eta}{1+1.1\eta} \bar{V}_t^2 + \frac{1}{10} \frac{(1-2\kappa_d)\eta^2}{(1+1.1\eta)(1+1.2\eta)} \bar{V}_t^2 + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} \\ & + \frac{\eta}{1+1.1\eta} \left(\frac{\Lambda_{u_1}^2}{1-2\kappa_d} + \tilde{C}\eta \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2 r_s \boldsymbol{\Lambda}_{u_2} \right) + \bar{\zeta}_{t+1} - \bar{\boldsymbol{B}}_t \\ & - \frac{2(1-2\kappa_d)\eta}{1+1.1\eta} \left(\boldsymbol{\Lambda}_{u_1} \bar{\boldsymbol{B}}_t - \frac{1.5\kappa_d + 0.5}{\kappa_d (1-2\kappa_d)} \bar{\boldsymbol{B}}_t^2 \right) \\ & \overset{(e)}{\preceq} \bar{\boldsymbol{V}}_t \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta(1-2\kappa_d)}{1+1.2\eta} \bar{\boldsymbol{V}}_t \right)^{-1} + \frac{\eta}{1+1.2\eta} \left(\frac{\Lambda_{u_1}^2}{1-2\kappa_d} + \tilde{C}\eta \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2 r_s \boldsymbol{\Lambda}_{u_1} \right) \\ & + \bar{\zeta}_{t+1} + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} - \bar{\boldsymbol{B}}_{t+1} \\ & \preceq \bar{\boldsymbol{V}}_{t+1} + \bar{\zeta}_{t+1} + \frac{2\kappa_d \Lambda_{u_1}}{1-2\kappa_d} - \bar{\boldsymbol{B}}_{t+1}, \end{split}$$

where we used Proposition 30 and (F.17) in (e).

Lower bound proof: Similar to the upper bound proof, here we will show that for $t < T_{bad}$,

$$\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}} \leq \boldsymbol{V}_{t}^{-} \quad \Rightarrow \quad \underline{\boldsymbol{V}}_{t+1} + \underline{\boldsymbol{\zeta}}_{t+1} + \underline{\boldsymbol{B}}_{t+1} - \frac{2\kappa_{d}}{1 + 2\kappa_{d}}\boldsymbol{\Lambda}_{\ell_{1}} \leq \boldsymbol{V}_{t+1}^{-}. (F.20)$$

Since the base case holds at t=0 and $\mathcal{T}_{bad}>0$, it remains to prove (F.20). By (F.6), we have

$$\begin{split} & \boldsymbol{V}_{t+1}^{-} \succeq \left(\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right) \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}})\right)^{-1} + \eta\boldsymbol{\Lambda}_{\ell_{1}}^{2} \\ & - C\eta^{2} \|\boldsymbol{\Lambda}\|_{F}^{2} r_{s} \boldsymbol{\Lambda}_{\ell_{2}} + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}} \\ & = \underline{\boldsymbol{V}}_{t} \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t}\right)^{-1} + \eta\boldsymbol{\Lambda}_{\ell_{1}}^{2} - C\eta^{2} \|\boldsymbol{\Lambda}\|_{F}^{2} r_{s} \boldsymbol{\Lambda}_{\ell_{2}} + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1} \boldsymbol{\Lambda}_{\ell_{2}}^{\frac{1}{2}} \\ & + \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t}\right)^{-1} (\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}) \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}})\right)^{-1} \end{split}$$

By using Proposition 29, we have for $t < T_{\text{bad}}$

$$\begin{split} & \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t}\right)^{-1} \!\! \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right) \! \left(\boldsymbol{I}_{\mathsf{rk}} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}})\right)^{-1} \\ & \succeq \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right) - \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t} \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right) - \frac{\eta^{2}\kappa_{d}^{2}}{(1 - 1.1\eta)^{2}} \underline{\boldsymbol{V}}_{t}^{4} \\ & - \frac{\eta}{1 - 1.1\eta} \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right) \underline{\boldsymbol{V}}_{t} - \frac{\eta}{1 - 1.1\eta} \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\kappa_{d}}\right)^{2} \end{split}$$

$$\begin{split} &-\frac{2\eta^2/\kappa_d^2}{(1-1.1\eta)^2}\left(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\right)^2 - \frac{\eta^2\kappa_d^2}{(1-1.1\eta)^2}\underline{\boldsymbol{V}}_t\left(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\right)^2\underline{\boldsymbol{V}}_t \\ &+ \frac{\eta^2}{(1-1.1\eta)^2}\Big(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\Big)\underline{\boldsymbol{V}}_t\left(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\right) - \eta^3\tilde{C}_2\boldsymbol{\Lambda}_{\ell_1} \\ &+ \frac{\eta^2}{(1-1.1\eta)^2}\left(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\right)^3 + \frac{\eta^2}{(1-1.1\eta)^2}\underline{\boldsymbol{V}}_t\left(\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d\boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\right)\underline{\boldsymbol{V}}_t \end{split}$$

for some $\tilde{C}_2 = O(1)$. We have the following: First:

$$\begin{split} -\mathbf{K}_{d}^{2}\underline{\boldsymbol{V}}_{t}^{4} - \mathbf{K}_{d}^{2}\underline{\boldsymbol{V}}_{t} \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\mathbf{K}_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\mathbf{K}_{d}}\right)^{2}\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{V}}_{t} \left(\underline{\boldsymbol{\zeta}}_{t} + \underline{\boldsymbol{B}}_{t} - \frac{2\mathbf{K}_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1 + 2\mathbf{K}_{d}}\right)\underline{\boldsymbol{V}}_{t} \overset{(f)}{\succeq} - 3.2\mathbf{K}_{d}\underline{\boldsymbol{V}}_{t}^{2} \\ & \overset{(g)}{\succeq} - \frac{1}{15}\frac{1 - 1.1\mathbf{\eta}}{1 - 1.2\mathbf{\eta}}\underline{\boldsymbol{V}}_{t}^{2} \,. \end{split}$$

where (f) follows by $\|\underline{\boldsymbol{V}}_t\|_2 \leq 5$ and $\left\|\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d \Lambda_{\ell_1}}{1+2\kappa_d}\right\|_2 \leq 2.5\kappa_d$, and (g) follows by $\kappa_d \leq \frac{1}{50}$. Second:

$$\begin{split} &-\frac{2}{\kappa_{d}^{2}}\bigg(\underline{\zeta}_{t}+\underline{\boldsymbol{B}}_{t}-\frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1+2\kappa_{d}}\bigg)^{2}+\bigg(\underline{\zeta}_{t}+\underline{\boldsymbol{B}}_{t}-\frac{2c_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1+2\kappa_{d}}\bigg)\underline{\boldsymbol{V}}_{t}\left(\underline{\zeta}_{t}+\underline{\boldsymbol{B}}_{t}-\frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1+2\kappa_{d}}\right)\\ &+\bigg(\underline{\zeta}_{t}+\underline{\boldsymbol{B}}_{t}-\frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1+2\kappa_{d}}\bigg)^{3}\overset{(h)}{\succeq}\frac{-3}{\kappa_{d}^{2}}\left(\underline{\zeta}_{t}+\underline{\boldsymbol{B}}_{t}-\frac{2\kappa_{d}\boldsymbol{\Lambda}_{\ell_{1}}}{1+2\kappa_{d}}\right)^{2}\succeq\frac{-9}{\kappa_{d}^{2}}\left(\underline{\zeta}_{t}^{2}+\underline{\boldsymbol{B}}_{t}^{2}\right)-36\boldsymbol{\Lambda}_{\ell_{1}}^{2}, \end{split}$$

where (h) follows by $\|\underline{\boldsymbol{V}}_t\|_2 \leq 5$ and $\|\underline{\boldsymbol{\zeta}}_t + \underline{\boldsymbol{B}}_t - \frac{2\kappa_d \boldsymbol{\Lambda}_{\ell_1}}{1+2\kappa_d}\|_2 \leq 2.5\kappa_d$. Third:

$$\begin{split} &-\underline{V}_{t}\left(\underline{\zeta}_{t}+\underline{B}_{t}-\frac{2\kappa_{d}\Lambda_{\ell_{1}}}{1+2\kappa_{d}}\right)-\left(\underline{\zeta}_{t}+\underline{B}_{t}-\frac{2\kappa_{d}\Lambda_{\ell_{1}}}{1+2\kappa_{d}}\right)\underline{V}_{t}-\left(\underline{\zeta}_{t}+\underline{B}_{t}-\frac{2\kappa_{d}\Lambda_{\ell_{1}}}{1+2\kappa_{d}}\right)^{2} \\ &-\frac{9\eta/\kappa_{d}^{2}}{1-1.1\eta}\left(\underline{\zeta}_{t}^{2}+\underline{B}_{t}^{2}\right) \\ &\succeq -2\left(\underline{K}_{t}\underline{\zeta}_{t}+\underline{\zeta}_{t}\underline{K}_{t}\right)-2\left(\underline{K}_{t}\underline{B}_{t}+\underline{B}_{t}\underline{K}_{t}\right)+\frac{4\kappa_{d}}{1+2\kappa_{d}}\left(\underline{K}_{t}\Lambda_{\ell_{1}}+\Lambda_{\ell_{1}}\underline{K}_{t}\right)-3\left(\underline{\zeta}_{t}^{2}+\underline{B}_{t}^{2}\right) \\ &+\left(\Lambda_{\ell_{1}}\underline{\zeta}_{t}+\underline{\zeta}_{t}\Lambda_{\ell_{1}}\right)+\left(\Lambda_{\ell_{1}}\underline{B}_{t}+\underline{B}_{t}\Lambda_{\ell_{1}}\right)-\frac{4\kappa_{d}(1+\kappa_{d})\Lambda_{\ell_{1}}^{2}}{(1+2\kappa_{d})^{2}} \\ &\stackrel{(i)}{\succeq}-8c_{d}\underline{K}_{t}^{2}+\frac{4\kappa_{d}}{1+2\kappa_{d}}\left(\underline{K}_{t}\Lambda_{\ell_{1}}+\Lambda_{\ell_{1}}\underline{K}_{t}\right)-\frac{4\kappa_{d}(1+\kappa_{d})\Lambda_{\ell_{1}}^{2}}{(1+2\kappa_{d})^{2}} \\ &+(2-4\kappa_{d})\Lambda_{\ell_{1}}\underline{B}_{t}-\left(3+\frac{1}{\kappa_{d}}\right)\underline{B}_{t}^{2} \\ &=-2\kappa_{d}\underline{V}_{t}^{2}-\frac{2\kappa_{d}\Lambda_{\ell_{1}}^{2}}{(1+2\kappa_{d})}+(2-4\kappa_{d})\underline{B}_{t}\Lambda_{\ell_{1}}-\left(3+\frac{1}{\kappa_{d}}\right)\underline{B}_{t}^{2}, \end{split}$$

where we used Proposition 22, and the first and third items in Lemma 4 in (i). Therefore, we have

$$\begin{split} & (\mathbf{F}.21) \succeq \underline{V}_t \left(I_r + \frac{\eta}{1 - 1.1\eta} \underline{V}_t \right)^{-1} - \frac{2\kappa_d \eta}{1 - 1.1\eta} \underline{V}_t^2 - \frac{(1 + 2\kappa_d)\eta^2/15}{(1 - 1.1\eta)(1 - 1.2\eta)} \underline{V}_t^2 - \frac{2\kappa_d}{1 + 2\kappa_d} \mathbf{\Lambda}_{\ell_1} \\ & + \frac{\eta}{1 - 1.1\eta} \left(\frac{\mathbf{\Lambda}_{\ell_1}^2}{1 + 2\kappa_d} - \tilde{C}\eta \|\mathbf{\Lambda}\|_{\mathrm{F}}^2 r_s \mathbf{\Lambda}_{\ell_2} \right) + \underline{\zeta}_{t+1} + \underline{B}_t \\ & + \frac{2(1 - 2\kappa_d)\eta}{1 - 1.1\eta} \left(\mathbf{\Lambda}_{\ell_1} \underline{B}_t - \frac{1.5\kappa_d + 0.5}{\kappa_d (1 - 2\kappa_d)} \underline{B}_t^2 \right) \\ & \succeq \underline{V}_t \left(I_r + \frac{\eta(1 + 2\kappa_d)}{1 - 1.2\eta} \underline{V}_t \right)^{-1} + \frac{\eta}{1 - 1.1\eta} \left(\frac{\mathbf{\Lambda}_{\ell_1}^2}{1 + 2\kappa_d} - \tilde{C}\eta \|\mathbf{\Lambda}\|_{\mathrm{F}}^2 r_s \mathbf{\Lambda}_{\ell_1} \right) + \underline{\zeta}_{t+1} \\ & + \underline{B}_{t+1} - \frac{2\kappa_d}{1 + 2\kappa_d} \mathbf{\Lambda}_{\ell_1} \end{split}$$

$$= \underline{V}_{t+1} + \underline{\zeta}_{t+1} + \underline{B}_{t+1} - \frac{2\kappa_d}{1 + 2\kappa_d} \Lambda_{\ell_1},$$

where we used Proposition 30 and (F.16) in (j).

F.4 Analysis of the bounding systems

F.4.1 Lower bounding system

In this section, we consider (F.12). For notational convenience, we multiply both sides by the factor $(1 + 2\kappa_d)$ and use a generic learning rate η , i.e.,

$$\underline{\boldsymbol{V}}_{t+1} = \underline{\boldsymbol{V}}_t (\boldsymbol{I}_{\mathsf{rk}} + \eta \underline{\boldsymbol{V}}_t)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_1}^2 - \tilde{C} \eta \| \boldsymbol{\Lambda} \|_{\mathrm{F}}^2 r_s \boldsymbol{\Lambda}_{\ell_1} \right), \text{ where } \underline{\boldsymbol{V}}_t = 2 \boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}} \underline{\boldsymbol{G}}_t \boldsymbol{\Lambda}_{\ell_2}^{\frac{1}{2}} - \boldsymbol{\Lambda}_{\ell_1}.$$

The main result of this section is stated in Proposition 15. To establish it, we first prove an auxiliary result, Lemma 5. For the following, we define

$$\hat{\boldsymbol{\Lambda}} \coloneqq \sqrt{\boldsymbol{\Lambda}_{\ell_1}^2 - \tilde{C}\boldsymbol{\eta}\|\boldsymbol{\Lambda}\|_{\mathrm{F}}^2 r_s \boldsymbol{\Lambda}_{\ell_1}} = \mathrm{diag}(\{\hat{\lambda}_i\}_{i=1}^r), \quad \boldsymbol{D}_t \coloneqq \frac{\boldsymbol{\Lambda}_{\ell_2}^{-1} \hat{\boldsymbol{\Lambda}}\left(\frac{\boldsymbol{A}_{t,11}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}}\right)}{2} - \frac{1.1 \kappa_d r_s}{d} \boldsymbol{I}_{\mathsf{rk}}.$$

By Corollary 3, we have

$$\underline{\boldsymbol{G}}_{t} = \frac{1}{2} \left(\frac{\boldsymbol{\Lambda}_{\ell_{1}}}{\boldsymbol{\Lambda}_{\ell_{2}}} + \frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} \frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_{2}}} \right) - \frac{1}{4} \frac{\boldsymbol{A}_{t,12}^{-1} \hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_{2}}} \left(\frac{\hat{\boldsymbol{\Lambda}}_{\ell_{1}11}}{2} - \boldsymbol{I}_{rk} \right) + \frac{\left(\hat{\boldsymbol{\Lambda}} - \boldsymbol{\Lambda}_{\ell_{1}} \right)}{2 \boldsymbol{\Lambda}_{\ell_{2}}} + \underline{\boldsymbol{G}}_{0} \right)^{-1} \frac{\hat{\boldsymbol{\Lambda}} \boldsymbol{A}_{t,12}^{-1}}{\boldsymbol{\Lambda}_{\ell_{2}}}, \tag{F.22}$$

where $A_{t,11}$, $A_{t,12}$, and $A_{t,22}$ are defined as in (R.1) with $\hat{\Lambda}$. For $\alpha = 0$, we will consider $\{\underline{G}_t\}_{t \in \mathbb{N}}$ in the basis of \underline{G}_0 without writing explicitly, which will imply that $\{\underline{G}_t\}_{t \in \mathbb{N}}$ is diagonal due to the rotational symmetry for $\alpha = 0$.

We further decompose $\{\underline{G}_t\}_{t \in N}$ and related matrices to isolate their top-left submatrices of dimension $\mathsf{rk}_\star \in \{r_\star, r_{u_\star}\}$, where $r_\star < r$ and $r_{u_\star} < r_u$ which we will denote as $\mathsf{rk}_\star < \mathsf{rk}$. The decompositions are as follows:

$$\underline{\boldsymbol{G}}_t \coloneqq \begin{bmatrix} \underline{\boldsymbol{G}}_{t,11} & \underline{\boldsymbol{G}}_{t,12} \\ \underline{\boldsymbol{G}}_{t,12}^\top & \underline{\boldsymbol{G}}_{t,22} \end{bmatrix}, \quad \hat{\boldsymbol{\Lambda}} \coloneqq \begin{bmatrix} \hat{\boldsymbol{\Lambda}}_{11} & \boldsymbol{0} \\ \boldsymbol{0} & \hat{\boldsymbol{\Lambda}}_{22} \end{bmatrix}, \quad \boldsymbol{D}_t \coloneqq \begin{bmatrix} \boldsymbol{D}_{t,1} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{D}_{t,2} \end{bmatrix}, \quad \boldsymbol{Z}_{1:\mathsf{rk}} \coloneqq \begin{bmatrix} \boldsymbol{Z}_{1:\mathsf{rk}_\star} \\ \boldsymbol{Z}_2 \end{bmatrix},$$

where $\underline{G}_{t,11}, D_{t,1}, \hat{\Lambda}_{11} \in \mathbb{R}^{\mathsf{rk}_\star \times \mathsf{rk}_\star}.$ We define

$$\Gamma_t \coloneqq \boldsymbol{D}_t + \frac{1}{1.05} \boldsymbol{Z}_{1:\mathsf{rk}} \boldsymbol{Z}_{1:\mathsf{rk}}^\top = \begin{bmatrix} \frac{1}{1.05} \boldsymbol{Z}_{1:\mathsf{rk}_\star} \boldsymbol{Z}_{1:\mathsf{rk}_\star}^\top + \boldsymbol{D}_{t,1} & \frac{1}{1.05d} \boldsymbol{Z}_{1:\mathsf{rk}_\star} \boldsymbol{Z}_2^\top \\ \frac{1}{1.05} \boldsymbol{Z}_2 \boldsymbol{Z}_{1:\mathsf{rk}_\star}^\top & \frac{1}{1.05} \boldsymbol{Z}_2 \boldsymbol{Z}_2^\top + \boldsymbol{D}_{t,2} \end{bmatrix}$$

and

$$\boldsymbol{\Gamma}_t^{-1} \coloneqq \begin{bmatrix} (\boldsymbol{\Gamma}_t^{-1})_{11} & (\boldsymbol{\Gamma}_t^{-1})_{12} \\ (\boldsymbol{\Gamma}_t^{-\top})_{12} & (\boldsymbol{\Gamma}_t^{-1})_{22} \end{bmatrix}$$

whenever Γ_t is invertible. Lemma 5 is stated as follows:

Lemma 5. We consider the following setting:

$$\alpha \in [0, 0.5): \quad \frac{r_s}{r} \to \varphi \in (0, \infty), \quad \eta \ll \frac{1}{d r^{1-\alpha} \log^4 d}, \quad \kappa_d = \frac{1}{\log^{3.5} d},$$

$$\alpha > 0.5: \qquad r_s \approx 1, \qquad \qquad \eta \ll \frac{1}{d r_u^{2+\alpha} \log^3 d}, \quad \kappa_d = \frac{1}{r_u \log^{2.5} d}.$$

 G_{init} implies the following:

• For $\alpha \geq 0$ and $K \leq \operatorname{rk}_{\star} \leq \operatorname{rk}$, we have for $\eta t \leq \frac{1}{2}(K+1)^{\alpha}\log\Big(\frac{d\log^{1.5}d}{r_s}\Big)$,

$$\boldsymbol{D}_{t} \succeq \frac{\boldsymbol{\Lambda}_{\ell_{2}}^{-1} \hat{\boldsymbol{\Lambda}} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}} \right)}{2} - \frac{1.2 \kappa_{d} r_{s}}{d} \boldsymbol{I}_{\mathsf{rk}}, \quad \boldsymbol{D}_{t,2} \succeq \frac{\log^{3} d - 1}{\log^{3} d} \left(\frac{0.5 r_{s}}{d \log^{1.5} d} \right)^{\left(\frac{K+1}{\mathsf{rk}_{\star} + 1} \right)^{\alpha}} \boldsymbol{I}_{\mathsf{rk} - \mathsf{rk}_{\star}}.$$

• For $r_{\star} = \lfloor r_s \left(1 - \log^{\frac{-1}{2}} d\right) \wedge r \rfloor$ and $r_{u_{\star}} = r_s$, and $\eta t \leq \frac{1}{2} (\mathsf{rk}_{\star} + 1)^{\alpha} \log\left(\frac{d \log^{1.5} d}{r_s}\right)$, we have

$$\Gamma_{t} \succeq \frac{\boldsymbol{\Lambda}_{\ell_{2}}^{-1} \hat{\boldsymbol{\Lambda}} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}} \right)}{2} + \begin{bmatrix} \frac{C_{1} r_{s}}{d \log^{4.5} d} \boldsymbol{I}_{\mathsf{rk}_{\star}} & 0\\ 0 & -\frac{C_{2} r_{s}}{d \log^{2} d} \boldsymbol{I}_{\mathsf{rk} - \mathsf{rk}_{\star}} \end{bmatrix} \succeq 0, \tag{F.23}$$

where

$$C_1 = \begin{cases} \frac{1}{10}, & \alpha \in [0, 0.5) \\ \left(\frac{1}{1.1r_s^6} - \frac{1.3}{\sqrt{\log d}}\right), & \alpha > 0.5 \end{cases} \quad \text{and} \quad C_2 = \begin{cases} 2.1 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2, & \alpha \in [0, 0.5) \\ 2, & \alpha > 0.5. \end{cases}$$

Within the same time interval, we have

$$\boldsymbol{\Gamma}_{t}^{-1} \preceq \left(\frac{\boldsymbol{\Lambda}_{\ell_{2}}^{-1} \hat{\boldsymbol{\Lambda}} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}} \right)}{2} + \begin{bmatrix} \frac{C_{1} r_{s}}{d \log^{4.5} d} \boldsymbol{I}_{\mathsf{rk}_{\star}} & 0\\ 0 & \frac{-C_{2} r_{s}}{d \log^{2} d} \boldsymbol{I}_{\mathsf{rk} - \mathsf{rk}_{\star}} \end{bmatrix} \right)^{-1}. \tag{F.24}$$

• For $\alpha > 0$, we have

$$(\boldsymbol{\Gamma}_t^{-1})_{11} \preceq \left(\boldsymbol{D}_{t,1} + \frac{1}{2} \boldsymbol{Z}_{1:\mathsf{rk}_\star} \boldsymbol{Z}_{1:\mathsf{rk}_\star}^\top \right)^{-1},$$

for $0.001 > \delta \ge \log^{\frac{-1}{4}} d$

$$\begin{cases} r_{\star} = \lfloor r_{s} \left(1 - \delta \right) \wedge r \rfloor & \text{and } \eta t \leq \frac{1}{2} \left(r_{s} \left(1 - \sqrt{\delta} \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_{s}} \right), & \alpha \in (0, 0.5) \\ r_{u_{\star}} = r_{s} & \text{and } \eta t \leq \frac{1}{2} r_{s}^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_{s}} \right), & \alpha > 0.5 \end{cases}$$

provided that

$$\begin{cases} d \ge \Omega_{\beta}(1) \lor \exp\left(2.5\alpha^{-8}\right), & \alpha \in (0, 0.5) \\ d \ge \Omega_{r_*}(1), & \alpha > 0.5. \end{cases}$$

Proof. For the first part of the first item, by (R.2), we have

$$\boldsymbol{D}_{t} = \frac{\boldsymbol{\Lambda}_{\ell_{2}}^{-1}\hat{\boldsymbol{\Lambda}}\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}}\right)}{2} - \frac{\eta}{2}\boldsymbol{\Lambda}_{\ell_{2}}^{-1}\hat{\boldsymbol{\Lambda}}^{2} - \frac{1.1\mathsf{K}_{d}r_{s}}{d}\boldsymbol{I}_{\mathsf{rk}} \overset{(a)}{\succeq} \frac{\boldsymbol{\Lambda}_{\ell_{2}}^{-1}\hat{\boldsymbol{\Lambda}}\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}}\right)}{2} - \frac{1.2\mathsf{K}_{d}r_{s}}{d}\boldsymbol{I}_{\mathsf{rk}},$$

where (a) follows $\eta \ll \frac{\kappa_d r_s}{d}$. Moreover, since $\Lambda_{\ell_2}^{-1} \hat{\Lambda} \succeq (1 - \frac{1}{\log^3 d}) I_{rk}$, by (R.4), we have

$$D_{t,2} \succeq \left(1 - \frac{1}{\log^3 d}\right) \left[\frac{\left(\boldsymbol{I}_{\mathsf{rk-rk_{\star}}} - \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^t}{\left(\boldsymbol{I}_{\mathsf{rk-rk_{\star}}} + \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^t - \left(\boldsymbol{I}_{\mathsf{rk-rk_{\star}}} - \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^t} - \frac{1.3 \kappa_d r_s}{d} \boldsymbol{I}_{\mathsf{rk-rk_{\star}}} \right] (F.25)$$

We observe that

$$\begin{split} \frac{\left(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^{t}}{\left(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} + \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^{t} - \left(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - \eta \hat{\boldsymbol{\Lambda}}_{22}\right)^{t}} \succeq \frac{\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}}}{\exp\left(\frac{2t\eta \hat{\boldsymbol{\Lambda}}_{22}}{1 - \eta \hat{\boldsymbol{\Lambda}}_{22}}\right) - \boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}}} \\ & \stackrel{(b)}{\succeq} \left(\frac{r_{s}}{d \log^{1.5} d}\right)^{(1 + 2\eta \mathsf{rk}_{\star}^{-\alpha})\left(\frac{K + 1}{\mathsf{rk}_{\star} + 1}\right)^{\alpha}} \boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} \\ & \stackrel{(c)}{\succeq} \left(\frac{0.9r_{s}}{d \log^{1.5} d}\right)^{\left(\frac{K + 1}{\mathsf{rk}_{\star} + 1}\right)^{\alpha}} \boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}}, \end{split}$$

where we use $\eta t \leq \frac{1}{2}(K+1)^{\alpha}\log\left(\frac{d\log^{1.5}d}{r_s}\right)$ in (b) and $d \geq \Omega(1)$ in (c). By (F.25), the first item follows.

For the second item, by Proposition 24 with $\varepsilon = \frac{1}{\log^2 d}$ for $\alpha \in [0, 0.5)$ and $\varepsilon = \frac{1}{r_u \log^2 d}$ for $\alpha > 0.5$, we have

$$\boldsymbol{\Gamma}_t \succeq \frac{\boldsymbol{\Lambda}_{\ell_2}^{-1}\hat{\boldsymbol{\Lambda}}\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}}\right)}{2} - \frac{1.2\kappa_d r_s}{d}\boldsymbol{I}_{\mathsf{rk}} + \frac{1}{1.05}\begin{bmatrix}\varepsilon\boldsymbol{Z}_{1:\mathsf{rk}_\star}\boldsymbol{Z}_{1:\mathsf{rk}_\star}^\top & 0\\ 0 & \frac{-\varepsilon}{1-\varepsilon}\boldsymbol{Z}_2\boldsymbol{Z}_2^\top\end{bmatrix}.$$

For $\alpha \in [0, 0.5)$, since $\kappa_d = \frac{1}{\log^{3.5} d}$, by (H.1), we have

$$\frac{\varepsilon}{1.05} \boldsymbol{Z}_{1:r_{\star}} \boldsymbol{Z}_{1:r_{\star}}^{\top} - \frac{1.2\kappa_{d}r_{s}}{d} \boldsymbol{I}_{r_{\star}} \succeq \frac{r_{s}}{d\log^{3}d} \frac{1}{6.25} \boldsymbol{I}_{r_{\star}} - \frac{1.2\kappa_{d}r_{s}}{d} \boldsymbol{I}_{r_{\star}} \succeq \frac{1}{10} \frac{r_{s}}{d\log^{3}d} \boldsymbol{I}_{r_{\star}}.$$

Similarly by (H.3), we have

$$\frac{-\varepsilon}{1-\varepsilon} \frac{1}{1.05} \mathbf{Z}_2 \mathbf{Z}_2^{\top} - \frac{1.2 \kappa_d r_s}{d} \mathbf{I}_{r-r_{\star}} \succeq -\left(1 + \frac{1}{\sqrt{\varphi}}\right)^2 \frac{2.1 r_s}{d \log^2 d} \mathbf{I}_{r-r_{\star}}.$$

For $\alpha > 0.5$, since $\kappa_d = \frac{1}{r_u \log^{2.5} d}$ and $r_u = \lceil \log^{2.5} d \rceil$, by (L.1), we have

$$\frac{\varepsilon}{1.05} \boldsymbol{Z}_{1:r_{u_{\star}}} \boldsymbol{Z}_{1:r_{u_{\star}}}^{\intercal} - \frac{1.2 \kappa_{d} r_{s}}{d} \boldsymbol{I}_{r_{u_{\star}}} \succeq \frac{r_{s}}{d \log^{4.5} d} \left(\frac{1}{1.1 r_{s}^{6}} - \frac{1.3}{\sqrt{\log d}} \right) \boldsymbol{I}_{r_{u_{\star}}}.$$

Similarly by (L.2),

$$\frac{-\varepsilon}{1-\varepsilon}\frac{1}{1.05d}\boldsymbol{Z}_2\boldsymbol{Z}_2^\top - \frac{1.2\kappa_d r_s}{d}\boldsymbol{I}_{r_u-r_{u_\star}} \succeq \frac{-2r_s}{d\log^2 d}\boldsymbol{I}_{r_u-r_{u_\star}}.$$

Therefore, we have (F.23). By Proposition 25, we have (F.24).

For the last item, we have

$$(\mathbf{\Gamma}_t^{-1})_{11} = \left(\mathbf{D}_{t,1} + \frac{1}{1.05} \mathbf{Z}_{1:\mathsf{rk}_{\star}} \left(\mathbf{I}_{r_s} + \frac{1}{1.05} \mathbf{Z}_2^{\top} \mathbf{D}_{t,2}^{-1} \mathbf{Z}_2 \right)^{-1} \mathbf{Z}_{1:\mathsf{rk}_{\star}}^{\top} \right)^{-1}. \tag{F.26}$$

For $\alpha \in (0,0.5)$, if $r_{\star} = r$, the statement follows. If not by the first item, for $K = \lfloor (1 - \sqrt{\delta})r_s \rfloor$ and $r_{\star} = \lfloor r_s(1 - \delta) \rfloor$, we have

$$\frac{K+1}{r_{\star}+1} \leq \frac{1-\sqrt{\delta}}{1-\delta} + \frac{2}{r_s} \leq 1 - 0.9\sqrt{\delta} \ \Rightarrow \ \left(\frac{K+1}{r_{\star}+1}\right)^{\alpha} \leq 1 - \alpha 0.9\sqrt{\delta}.$$

Therefore,

$$\boldsymbol{D}_{t,2} \succeq \left(\frac{0.5r_s}{d\log^{1.5} d}\right)^{1-\alpha 0.9\sqrt{\delta}} \boldsymbol{I}_{r-r_\star} \stackrel{(d)}{\succeq} \frac{0.5r_s}{d\log^{1.5} d} \left(\frac{d}{r_s}\right)^{\log^{-1/4} d} \boldsymbol{I}_{r-r_\star} \stackrel{(e)}{\succeq} \frac{r_s \log d}{d} \boldsymbol{I}_{r-r_\star},$$

where we used $d \ge \Omega(1) \lor \exp(2.5\alpha^{-8})$ in (d) and $d \ge \Omega_{\beta}(1)$ in (e). By (F.26) and (H.3), we have the statement for $\alpha \in (0,0.5)$. For $\alpha > 0.5$, $K = r_{\star} = r_s$, we have

$$\left(\frac{K+1}{r_{\star}+1}\right)^{\alpha} \le \left(1+\frac{1}{r_s+1}\right)^{0.5} \le 1-\frac{1}{2(r_s+1)}.$$

Therefore,

$$\boldsymbol{D}_{t,2} \succeq \left(\frac{0.5r_s}{d\log^{1.5}d}\right)^{1-\frac{1}{2(r_s+1)}} \boldsymbol{I}_{r_u-r_{u_\star}} \succeq \frac{r_s\log^8d}{d} \boldsymbol{I}_{r-r_\star}$$

for $d \ge \Omega_{r_s}(1)$. By (F.26) and (L.2), we have the statement for $\alpha > 0.5$.

Proposition 15. *Let*

$$\underline{\boldsymbol{G}}_{0} = (1 + 2\kappa_{d}) \left(\mathsf{G}_{0} - \frac{\kappa_{d} r_{s}}{d} \boldsymbol{I}_{\mathsf{rk}} \right),$$

Under the parameter choice in Lemma 5, G_{init} *guarantees that:*

• We have $\Omega(-\log^{\frac{-1}{2}}d)I_{\mathsf{rk}} \preceq \underline{G}_t$ whenever

$$\eta t \leq \begin{cases} \frac{1}{2} \left(r_s \left(1 - \log^{\frac{-1}{2}} d \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha > 0.5 \end{cases}.$$

• Let Λ_{11} be the $\operatorname{rk}_{\star} \times \operatorname{rk}_{\star}$ dimensional top-left sub-matrix of Λ . Given $0.001 \geq \delta \geq \log^{\frac{-1}{4}} d$ and $\operatorname{rk}_{\star} = \left\{ r_{\star} = \lfloor r_{s} (1 - \delta) \wedge r \rfloor, r_{u_{\star}} = r_{s} \right\}$, we have

$$\underline{G}_{t,11} \succeq \frac{1 - \frac{10}{\log^3 d}}{\frac{1 \cdot 2}{C_{lb}} \frac{d}{r_s} \exp\left(-2\eta t \mathbf{\Lambda}_{11}\right) + 1}$$

and

$$\|\hat{\mathbf{\Lambda}}\|_F^2 - \|\hat{\mathbf{\Lambda}}_1^{\frac{1}{2}}\underline{G}_{t,11}\hat{\mathbf{\Lambda}}_1^{\frac{1}{2}}\|_F^2 \leq \sum_{i=(\mathsf{rk}, \land r)+1}^r \hat{\lambda}_i^2 + \sum_{i=1}^{\mathsf{rk}_\star} \hat{\lambda}_i^2 \left(1 - \frac{1 - \frac{10}{\log^3 d}}{\frac{1\cdot 2}{C_h} \frac{d}{r_\star} \exp\left(-2\eta t \lambda_i\right) + 1}\right)^2,$$

for

$$C_{lb} = \frac{1}{15} \begin{cases} \delta^2, & \alpha \in [0, 0.5) \\ \frac{1}{r_s^6}, & \alpha > 0.5 \end{cases} \quad and \quad \eta t \le \begin{cases} \frac{1}{2} \left(r_s \left(1 - \sqrt{\delta} \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha > 0.5. \end{cases}$$

• For $\delta = \log^{\frac{-1}{4}} d$, we define

$$\mathcal{T}_{lb} := \inf \left\{ n \ge 0 \ \middle| \ \|\hat{\mathbf{\Lambda}}\|_F^2 - \|\hat{\mathbf{\Lambda}}_1^{\frac{1}{2}} \underline{G}_{t,11} \hat{\mathbf{\Lambda}}_1^{\frac{1}{2}} \|_F^2 \le \sum_{j=(r_s \wedge r)+1}^r \lambda_j^2 + \frac{3\|\hat{\mathbf{\Lambda}}\|_F^2}{\log^{\frac{1}{8}} d} \right\}.$$

Then,

$$\mathcal{T}_{lb} \le \begin{cases} \frac{1}{2\eta} \left(r_s \left(1 - \log^{\frac{-1}{8}} d \right) \wedge r \right)^{\alpha} \log \left(\frac{20d \log^{\frac{3}{4}} (1 + d/r_s)}{r_s} \right), & \alpha \in [0, 0.5) \\ \frac{1}{2\eta} r_s^{\alpha} \log \left(\frac{20d \log^{\frac{3}{4}} d}{r_s} \right), & \alpha > 0.5. \end{cases}$$

Proof. For $\alpha > 0.5$, we assume that d is large enough to guarantee that $\left(\frac{1}{1.1r_s^6} - \frac{1.3}{\sqrt{\log d}}\right) > 0$. We observe that

$$rac{oldsymbol{\Lambda}_{\ell_2}^{-1}\hat{oldsymbol{\Lambda}}\left(rac{oldsymbol{A}_{t,11}}{oldsymbol{A}_{t,12}}-oldsymbol{I}_{\mathsf{rk}}
ight)}{2}+rac{oldsymbol{\Lambda}_{\ell_2}^{-1}\left(\hat{oldsymbol{\Lambda}}-oldsymbol{\Lambda}_{\ell_1}
ight)}{2}+oldsymbol{\underline{G}}_0\succeq oldsymbol{\Gamma}_t,$$

where we used (E.1), $\Lambda_{\ell_2} \succeq \Lambda_{\ell_1}$ and $\eta \|\Lambda\|_F^2 r_s I_{\mathsf{rk}} \ll \frac{\kappa_d r_s}{d} \Lambda_{\ell_1}$.

For the first item, by using $\mathsf{rk}_\star = \left\{ r_\star = \lfloor r_s \left(1 - \log^{\frac{-1}{2}} d \right) \wedge r \rfloor, r_{u_\star} = r_s \right\}$, we define

$$m{D}_{lb} \coloneqq egin{bmatrix} rac{C_1 r_s}{d \log^{4.5} d} m{I}_{\mathsf{rk}_{\star}} & 0 \ 0 & rac{-C_2 r_s}{d \log^2 d} m{I}_{\mathsf{rk} - \mathsf{rk}_{\star}} \end{bmatrix} \ \ ext{and} \ \ ilde{m{D}}_{lb} \coloneqq rac{m{\Lambda}_{\ell_2}}{\hat{m{\Lambda}}} m{D}_{lb}.$$

We introduce submatrix notation for block-diagonal matrices. Specifically, we write

$$\tilde{\boldsymbol{D}}_{lb} = \begin{bmatrix} \tilde{\boldsymbol{D}}_{lb,1} & \boldsymbol{0} \\ \boldsymbol{0} & \tilde{\boldsymbol{D}}_{lb,2} \end{bmatrix} \text{ and } \frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} \pm \boldsymbol{I}_{\mathsf{rk}} = \begin{bmatrix} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} \pm \boldsymbol{I}_{\mathsf{rk}}\right)_{11} & \boldsymbol{0} \\ \boldsymbol{0} & \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} \pm \boldsymbol{I}_{\mathsf{rk}}\right)_{22} \end{bmatrix},$$

where the block dimensions of each submatrix match those of D_{lb} . We start with proving the lower bound part. By the second item in Lemma 5, we have

$$\underline{\boldsymbol{G}}_{t} \succeq \frac{1}{2} \sqrt{\frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_{2}}}} \left(\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{\mathsf{rk}} \right) - \boldsymbol{A}_{t,12}^{-1} \left(\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\mathsf{rk}} \right) + 2\tilde{\boldsymbol{D}}_{lb} \right)^{-1} \boldsymbol{A}_{t,12}^{-1} \right) \sqrt{\frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_{2}}}}$$

$$\begin{split} &+\frac{\frac{\boldsymbol{\Lambda}_{\ell_1}}{\boldsymbol{\Lambda}_{\ell_2}}-\frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_2}}}{2} \\ &\geq \frac{1}{2}\sqrt{\frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_2}}}\left(\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}}+\boldsymbol{I}_{\mathsf{rk}}\right)-\boldsymbol{A}_{t,12}^{-1}\left(\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}}-\boldsymbol{I}_{\mathsf{rk}}\right)+2\tilde{\boldsymbol{D}}_{lb}\right)^{-1}\boldsymbol{A}_{t,12}^{-1}\right)\sqrt{\frac{\hat{\boldsymbol{\Lambda}}}{\boldsymbol{\Lambda}_{\ell_2}}}(F.27) \end{split}$$

where we used $\Lambda_{\ell_1} \succ \hat{\Lambda}$ in the second line. We have

$$\left(\frac{A_{t,22}}{A_{t,12}} + I_{\mathsf{rk}}\right) - A_{t,12}^{-1} \left(\left(\frac{A_{t,22}}{A_{t,12}} - I_{\mathsf{rk}}\right) + 2\tilde{D}_{lb}\right)^{-1} A_{t,12}^{-1} \\
= \frac{(A_{t,22} + A_{t,12})(A_{t,22} - A_{t,12} + 2\tilde{D}_{lb}A_{t,12}) - I_{\mathsf{rk}}}{A_{t,12}(A_{t,22} - A_{t,12} + 2\tilde{D}_{lb}A_{t,12})} \\
\stackrel{(a)}{\succeq} \frac{2\tilde{D}_{lb} \left(\frac{A_{t,22}}{A_{t,12}} + I_{\mathsf{rk}}\right)}{\frac{A_{t,22}}{A_{t,12}} - I_{\mathsf{rk}} + 2\tilde{D}_{lb}} \\
= \begin{bmatrix} \frac{2\tilde{D}_{lb,1} \left(\frac{A_{t,22}}{A_{t,12}} + I_{\mathsf{rk}}\right)_{11}}{\left(\frac{A_{t,22}}{A_{t,12}} + I_{\mathsf{rk}}\right)_{11}} & 0 \\ \frac{2\tilde{D}_{lb,2} \left(\frac{A_{t,22}}{A_{t,12}} + I_{\mathsf{rk}}\right)_{22}}{\left(\frac{A_{t,22}}{A_{t,12}} - I_{\mathsf{rk}}\right)_{22} + 2\tilde{D}_{lb,2}} \end{bmatrix}, \quad (F.28)$$

where we used $A_{t,22}^2 - A_{t,12}^2 \succ I_{\text{rk}}$ (by (R.2)) and $A_{t,22} - A_{t,12} + 2\tilde{D}_{lb}A_{t,12} \succ 0$ (by (F.23)) in (a). Since $\frac{A_{t,22}}{A_{t,12}} - I_{\text{rk}} \succ 0$ and $\tilde{D}_{lb,1} \succ 0$, it is enough to look at the bottom-right submatrix in (F.28) for the lower bound part. We have

$$\frac{2\tilde{\boldsymbol{D}}_{lb,2} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{rk}\right)_{22}}{\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{rk}\right)_{22} + 2\tilde{\boldsymbol{D}}_{lb,2}} = \frac{2\tilde{\boldsymbol{D}}_{lb,2} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{rk}\right)_{22}}{\left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{rk}\right)_{22} - 2\boldsymbol{I}_{rk-rk_{\star}} + 2\tilde{\boldsymbol{D}}_{lb,2}}.$$
(F.29)

Note that by (R.4).

$$\begin{split} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{\mathsf{rk}}\right)_2 &\succeq \frac{2(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} + \eta\hat{\boldsymbol{\Lambda}}_2)^t}{(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} + \eta\hat{\boldsymbol{\Lambda}}_2)^t - (\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - \eta\hat{\boldsymbol{\Lambda}}_2)^t} \\ &\succeq 2\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} + \frac{2(\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - \eta^2\hat{\boldsymbol{\Lambda}}_2^2)^t \exp\left(-2t\eta\hat{\boldsymbol{\Lambda}}_2\right)}{\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - (\boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}} - \eta^2\hat{\boldsymbol{\Lambda}}_2^2)^t \exp\left(-2t\eta\hat{\boldsymbol{\Lambda}}_2\right)} \\ &\stackrel{(b)}{\succeq} \left(2 + \frac{0.9r_s}{d\log^{1.5}d}\right) \boldsymbol{I}_{\mathsf{rk}-\mathsf{rk}_{\star}}, \end{split}$$

where we use $\eta t \leq \frac{1}{2} (\mathsf{rk}_\star + 1)^\alpha \log \left(\frac{d \log^{1.5} d}{r_s} \right)$ in (b). Hence, for $d \geq \Omega(1)$

$$(F.29) \succeq \frac{2\left(2 + \frac{0.9r_s}{d\log^{1.5}d}\right)\tilde{\boldsymbol{D}}_{lb,2}}{\frac{0.9r_s}{d\log^{1.5}d}\boldsymbol{I}_{rk-rk_{\star}} + \tilde{\boldsymbol{D}}_{lb,2}} \stackrel{(c)}{\succeq} \frac{12\tilde{\boldsymbol{D}}_{lb,2}}{\frac{r_s}{d\log^{1.5}d}\boldsymbol{I}_{rk-rk_{\star}}} \stackrel{(d)}{\succeq} \frac{-15C_2}{\log^{0.5}d}\boldsymbol{I}_{rk-rk_{\star}}, \tag{F.30}$$

where we used $\tilde{D}_{lb,2} \succeq \frac{-1.1C_2r_s}{d\log^2 d} I_{rk}$ in (c) and (d). The first item follows from (F.30).

For the second and third items, let $\left(\frac{A_{t,22}}{A_{t,12}}\pm I_{\rm rk}\right)_{11}$ denote the ${\rm rk}_\star\times{\rm rk}_\star$ dimensional top-left submatrices with ${\rm rk}_\star=\left\{r_\star=\left\lfloor r_s\left(1-\delta\right)\wedge r\right\rfloor, r_{u_\star}=r_s\right\}$. By using the third item in Lemma 5, we immediately observe that for $\alpha>0$, $\underline{G}_{t,11}\succeq 0$ and

$$\underline{\boldsymbol{G}}_{t,11} \stackrel{(e)}{\succeq} \left(1 - \frac{10}{\log^3 d}\right) \frac{1}{2} \frac{\frac{2C_{\text{b}}r_s}{d} \left(\frac{\boldsymbol{A}_{t,12}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{\text{rk}}\right)_{11}}{\left(\frac{\boldsymbol{A}_{t,11}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\text{rk}}\right)_{11} + \frac{2C_{\text{b}}r_s}{d} \boldsymbol{I}_{\text{rk}_{\star}}}, \tag{F.31}$$

for

$$C_{\mathrm{lb}} = \frac{1}{15} \begin{cases} \delta^2, & \alpha \in (0, 0.5) \\ \frac{1}{r_s^6}, & \alpha > 0.5 \end{cases} \quad \text{and} \quad \eta t \leq \begin{cases} \frac{1}{2} \left(r_s \left(1 - \sqrt{\delta} \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha \in (0, 0.5) \\ \frac{1}{2} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right), & \alpha > 0.5, \end{cases}$$

where we used $\Lambda_{\ell_2} \succeq \hat{\Lambda} \succeq \left(1 - \frac{0.5}{\log^4 d}\right) \Lambda_{\ell_2}$, and followed the steps in (F.27)- (F.28) with (H.1) and (L.1) to obtain (e). Then, by (R.4), we have

$$\frac{1}{2} \frac{\frac{2C_{\text{lb}}r_s}{d} \left(\frac{\boldsymbol{A}_{t,22}}{\boldsymbol{A}_{t,12}} + \boldsymbol{I}_{\text{rk}}\right)_{11}}{\left(\frac{\boldsymbol{A}_{t,11}}{\boldsymbol{A}_{t,12}} - \boldsymbol{I}_{\text{rk}}\right)_{11} + \frac{2C_{\text{lb}}r_s}{d}\boldsymbol{I}_{\text{rk}_{\star}}} \succeq \frac{\boldsymbol{I}_{\text{rk}_{\star}}}{\left(\frac{1}{C_{\text{lb}}}\frac{d}{r_s} - 1\right)\frac{(\boldsymbol{I}_{\text{rk}_{\star}} - \eta\hat{\boldsymbol{\Lambda}}_1)^t}{(\boldsymbol{I}_{\text{rk}_{\star}} + \eta\hat{\boldsymbol{\Lambda}}_1)^t} + \boldsymbol{I}_{\text{rk}_{\star}}}}{\frac{1}{C_{\text{lb}}}\frac{d}{r_s}\exp\left(-2\eta t\hat{\boldsymbol{\Lambda}}_1\right) + \boldsymbol{I}_{\text{rk}_{\star}}}}.$$

Consequently, by observing $\hat{\Lambda} \succeq \Lambda_{\ell_1} - \tilde{C} \eta I_{rk}$ and using the lower bounds for Λ_{ℓ_1} in Propositions 12 and 13, we have

$$\underline{G}_{t,11} \succeq \frac{1 - \frac{10}{\log^3 d}}{\frac{1.2}{C_{\text{lb}}} \frac{d}{r_s} \exp\left(-2\eta t \mathbf{\Lambda}_{11}\right) + 1},$$

where Λ_{11} denotes the $\mathsf{rk}_{\star} \times \mathsf{rk}_{\star}$ dimensional top-left sub-matrix of Λ . Therefore,

$$\|\hat{\mathbf{\Lambda}}\|_F^2 - \|\hat{\mathbf{\Lambda}}_1^{\frac{1}{2}}\underline{\mathbf{G}}_{t,11}\hat{\mathbf{\Lambda}}_1^{\frac{1}{2}}\|_F^2 \leq \sum_{i=(\mathsf{rk}_* \wedge r)+1}^r \hat{\lambda}_i^2 + \sum_{i=1}^{\mathsf{rk}_\star} \hat{\lambda}_i^2 \left(1 - \frac{1 - \frac{10}{\log^3 d}}{\frac{1\cdot 2}{C_{\mathsf{lb}}}\frac{d}{C_s}\exp\left(-2\eta t\lambda_i\right) + 1}\right)^2, (F.32)$$

which proves the second item for $\alpha>0$. Moreover, since (F.22) is in the eigenbasis of \underline{G}_0 , the arguments in (F.27)-(F.28) and the condition in (H.2) extend (F.31) to $\alpha=0$ in the eigenbasis of \underline{G}_0 for $d\geq\Omega(1)$. Given (F.31), we can extend (F.32) to $\alpha=0$ as the Frobenious norm is basis independent.

For the third item, for $\alpha>0.5$ and $t\geq \frac{1}{2\eta}r_s^{\alpha}\log\left(\frac{20d\log^{\frac{3}{4}}d}{r_s}\right)$, we have

$$(\textbf{F.32}) \leq \sum_{i=(r_2 \wedge r)+1}^r \hat{\lambda}_i^2 + \frac{\|\hat{\mathbf{\Lambda}}\|_F^2}{\log^{\frac{1}{2}} d} \leq \sum_{i=(r_2 \wedge r)+1}^r \lambda_i^2 + \frac{\|\hat{\mathbf{\Lambda}}\|_F^2}{\log^{\frac{1}{2}} d},$$

which gives us the corresponding bound for \mathcal{T}_{lb} .

For $\alpha \in [0, 0.5)$ and $t \geq \frac{1}{2\eta} \left(r_s \left(1 - \log^{\frac{-1}{8}} d \right) \wedge r \right)^{\alpha} \log \left(\frac{20d \log^{\frac{3}{4}} (1 + d/r_s)}{r_s} \right)$, we have

$$\begin{aligned} & (\text{F.32}) \leq \sum_{i=(r_s \wedge r)+1}^{r} \hat{\lambda}_i^2 + \sum_{i=\lfloor r_s \left(1 - \log \frac{-1}{8}d\right) \wedge r \rfloor + 1}^{r_s \wedge r} \hat{\lambda}_i^2 \\ & + \sum_{i=1}^{\lfloor r_s \left(1 - \log \frac{-1}{8}d\right) \wedge r \rfloor} \hat{\lambda}_i^2 \left(1 - \frac{1 - \frac{10}{\log^3 d}}{\frac{1.2}{C_{\text{lb}}} \frac{d}{r_s} \exp\left(-2\eta t \lambda_i\right) + 1}\right)^2 \\ & \leq \sum_{i=(r_s \wedge r)+1}^{r} \lambda_i^2 + \frac{3\|\hat{\mathbf{\Lambda}}\|_F^2}{\log^{\frac{1}{8}}d}, \end{aligned}$$

which gives us its bound for \mathcal{T}_{lb} .

F.4.2 Upper bounding system

In this section, we consider (F.13). For notational convenience, we multiply both sides by the factor $(1 - 2\kappa_d)$ and use a generic learning rate η , i.e.,

$$\bar{\pmb{V}}_{t+1} = \bar{\pmb{V}}_t(\pmb{I}_{\mathsf{rk}} + \eta \bar{\pmb{V}}_t)^{-1} + \eta \left(\pmb{\Lambda}_{u_1}^2 + \tilde{C}\eta \|\pmb{\Lambda}\|_{\mathrm{F}}^2 r_s \pmb{\Lambda}_{u_1}\right), \text{ where } \bar{\pmb{V}}_t = 2\pmb{\Lambda}_{u_2}^{\frac{1}{2}} \bar{\pmb{G}}_t \pmb{\Lambda}_{u_2}^{\frac{1}{2}} - \pmb{\Lambda}_{u_1}$$

The main result of this section is stated in Proposition 16. To establish it, we first prove an auxiliary result:

Lemma 6. The following statement holds:

- The reference sequence satisfies $T_t \succeq \frac{\kappa_d r_s}{d} I_{\mathsf{rk}}$ and $\{t \geq 0 : \|T_t\|_2 > 1.2\kappa_d\} = \infty$.
- For $r_{u_{\star}} = 2r_s$, we have

$$\begin{cases}
\bar{\boldsymbol{G}}_{0} = \frac{2.2\left(1 + \frac{1}{\sqrt{\varphi}}\right)^{2} r_{s}}{d} \boldsymbol{I}_{r} \succeq \left(1 - 2\kappa_{d}\right) \left(\boldsymbol{G}_{0} + \frac{\kappa_{d} r_{s}}{d} \boldsymbol{I}_{r}\right), & \alpha \in [0, 0.5) \\
\bar{\boldsymbol{G}}_{0} = \frac{5.5}{d} \begin{bmatrix} 2r_{s} \boldsymbol{I}_{r_{u_{\star}}} & 0 \\ 0 & r_{u} \boldsymbol{I}_{r-r_{u_{\star}}} \end{bmatrix} \succeq \left(1 - 2\kappa_{d}\right) \left(\boldsymbol{G}_{0,11} + \frac{\kappa_{d} r_{s}}{d} \boldsymbol{I}_{r_{u}}\right), & \alpha > 0.5
\end{cases}$$

provided that G_{init} holds.

For the following, we introduce $\hat{T}_t := \frac{\Lambda_{u_2}}{\Lambda_{\ell_1}} \frac{(3\kappa_d + 1)}{\kappa_d (1 - \kappa_d)} T_t$. Note that for $d \ge \Omega(1)$, we have

$$\hat{\boldsymbol{T}}_{t+1} = \hat{\boldsymbol{T}}_t + 2(1 - 2\kappa_d)\eta \boldsymbol{\Lambda}_{\ell_1} \hat{\boldsymbol{T}}_t \left(\boldsymbol{I}_{\mathsf{rk}} - \hat{\boldsymbol{T}}_t \right) \text{ and } \frac{\kappa_d}{1.1} \frac{\boldsymbol{\Lambda}_{\ell_1}}{\boldsymbol{\Lambda}_{u_2}} \leq \frac{\boldsymbol{T}_t}{\hat{\boldsymbol{T}}_t} \leq \kappa_d \boldsymbol{I}_{\mathsf{rk}}. \tag{F.34}$$

By Proposition 34, we have

 $1.1 \wedge \hat{T}_{0,ii} \exp(2\eta t \lambda_i) \geq \hat{T}_{t,ii}$

$$\geq \frac{1}{2} \begin{cases} 1 \wedge \hat{\boldsymbol{T}}_{0,ii} \exp\left(\frac{(1-2\kappa_{d})2\eta t \left(\lambda_{i} - \frac{0.1r^{-\alpha}}{\log^{4}d}\right)}{1+2(1-2\kappa_{d})\eta\lambda_{i}}\right), & \alpha \in [0, 0.5) \\ 1 \wedge \hat{\boldsymbol{T}}_{0,ii} \exp\left(\frac{(1-2\kappa_{d})2\eta t \left(\lambda_{i} - \frac{1}{(ru+1)^{\alpha}} - \frac{0.1}{r_{u}^{2+\alpha}\log^{4}d}\right)}{1+2(1-2\kappa_{d})\eta\lambda_{i}}\right), & \alpha > 0.5. \end{cases}$$
(F.35)

Proof of Lemma 6. For the first item, by Proposition 34, we have

$$\hat{T}_t \succeq \hat{T}_0 \overset{(a)}{\Rightarrow} T_t \succeq T_0 = \frac{\mathsf{\kappa}_d r_s}{d} I_{\mathsf{rk}},$$

where we multiplied each side with $\frac{\kappa_d(1-\kappa_d)}{3\kappa_d+1} \frac{\Lambda_{\ell_1}}{\Lambda_{\nu_2}}$ for (a). Moreover, by (F.34)-(F.35), we have

$$T_t \leq \kappa_d \hat{T}_t \leq 1.1 I_{\mathsf{rk}} \Rightarrow \{t \geq 0 : ||T_t||_2 > 1.2 \kappa_d\} = \infty.$$

The second item follows (E.2) and (H.3) (for $\alpha \in [0, 0.5)$) and (L.2) (for $\alpha > 0.5$).

Proposition 16. We consider $rk \in \{r, r_u\}$, where $r_u = \lceil \log^{2.5} d \rceil$, and

$$\alpha \in [0, 0.5): \quad \frac{r_s}{r} \to \varphi \in (0, \infty), \quad \eta \ll \frac{1}{d r^{1-\alpha} \log^4 d}, \quad \kappa_d = \frac{1}{\log^{3.5} d},$$

$$\alpha > 0.5: \qquad r_s \approx 1, \qquad \qquad \eta \ll \frac{1}{d r_u^{2+\alpha} \log^3 d}, \quad \kappa_d = \frac{1}{r_u \log^{2.5} d}$$

If \bar{G}_0 are taken as in (F.33), we have the following:

• $\{ar{G}_t\}_{n\in\mathbb{N}}$ is diagonal and satisfies

$$\frac{r_s}{d}\boldsymbol{I}_{\mathsf{rk}} \preceq \bar{\boldsymbol{G}}_{t+1} \preceq \bar{\boldsymbol{G}}_t + \eta \big((1+\kappa_d)\boldsymbol{\Lambda}_{u_1}\bar{\boldsymbol{G}}_t + (1+\kappa_d)\bar{\boldsymbol{G}}_t\boldsymbol{\Lambda}_{u_1} - 2\bar{\boldsymbol{G}}_t\boldsymbol{\Lambda}_{u_2}\bar{\boldsymbol{G}}_t \big) \preceq 1.1\boldsymbol{I}_{\mathsf{rk}}.$$

• For $\alpha \in [0, 0.5)$ and $d \ge \Omega(1)$, we have for $t \le \frac{1}{2\eta} r^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$:

$$- T_t^{-\frac{1}{2}} \bar{G}_j T_t^{-\frac{1}{2}} \leq \frac{11 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2}{\kappa_d} I_r \text{ for } 0 \leq j \leq t.$$

$$- T_t^{-\frac{1}{2}} \left(\eta \sum_{j=1}^t \bar{G}_{j-1} \right) T_t^{-\frac{1}{2}} \preceq \frac{5.5 \left(1 + \frac{1}{\sqrt{\varphi}} \right)^2}{\kappa_d} (2 \eta t \vee r^{\alpha}) I_r.$$

-
$$\bar{\boldsymbol{G}}_t \leq \left(1.1\bar{\boldsymbol{G}}_0 \exp\left(2\eta t \boldsymbol{\Lambda}\right) \wedge \boldsymbol{I}_r\right) + o_d(1)$$

$$- \| \mathbf{\Lambda} \|_{\mathrm{F}}^2 - Tr(\mathbf{\Lambda} \bar{\mathbf{G}}_t \mathbf{\Lambda}) \ge \sum_{i=1}^r \lambda_i^2 \left(1 - \frac{2.5 \left(1 + \frac{1}{\sqrt{\varphi}} \right)^2 r_s}{d} \exp(2\eta t \lambda_i) \right) - o_d(1).$$

• For
$$\alpha>0.5$$
 and $d\geq\Omega_{r_s}(1)$, we have for $t\leq \frac{1}{2\eta}r_s^{\alpha}\log\left(\frac{d\log^{1.5}d}{r_s}\right)$:

$$- T_t^{-\frac{1}{2}} \bar{G}_j T_t^{-\frac{1}{2}} \preceq \frac{26.4r_u}{\kappa_d} I_{r_u} \text{ for } 0 \leq j \leq t.$$

$$- T_t^{-\frac{1}{2}} \left(\eta \sum_{j=1}^t \bar{\boldsymbol{G}}_{j-1} \right) T_n^{-\frac{1}{2}} \preceq \frac{15r_u}{\kappa_d} (2r_s)^{\alpha} \log d\boldsymbol{I}_{r_u}.$$

$$- \bar{\boldsymbol{G}}_t \preceq (1.1\bar{\boldsymbol{G}}_0 \exp(2\eta t \boldsymbol{\Lambda}_{11}) \wedge \boldsymbol{I}_{r_u}) + o_d(1).$$

$$- \|\mathbf{\Lambda}_{11}\|_F^2 - Tr(\mathbf{\Lambda}_{11}\bar{\mathbf{G}}_{t,11}\mathbf{\Lambda}_{11}) \ge \sum_{i=1}^{r_s} \lambda_i^2 \left(1 - \frac{12.1r_s}{d} \exp\left(2\eta t \lambda_i\right)\right)_+ + \sum_{i=r_s+1}^{r_u} \lambda_i^2 - o_d(1).$$

Proof. Given that $\frac{r_s}{d} I_{\mathsf{rk}} \leq \bar{G}_t \leq 1.1 I_{\mathsf{rk}}$, we have

$$\bar{\boldsymbol{G}}_{t+1} \stackrel{(a)}{\leq} \bar{\boldsymbol{G}}_{t} + \eta \left(\boldsymbol{\Lambda}_{u_{1}} \bar{\boldsymbol{G}}_{t} + \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{1}} - 2 \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{2}} \bar{\boldsymbol{G}}_{t} \right) + 1.1 \tilde{C} \eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \boldsymbol{I}_{\mathsf{rk}}$$

$$\stackrel{(b)}{\leq} \bar{\boldsymbol{G}}_{t} + \eta \left((1 + \kappa_{d}) \boldsymbol{\Lambda}_{u_{1}} \bar{\boldsymbol{G}}_{t} + (1 + \kappa_{d}) \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{1}} - 2 \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{2}} \bar{\boldsymbol{G}}_{t} \right),$$

$$\bar{\boldsymbol{G}}_{t+1} \stackrel{(c)}{\succeq} \bar{\boldsymbol{G}}_{t} + \eta \left(\boldsymbol{\Lambda}_{u_{1}} \bar{\boldsymbol{G}}_{t} + \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{1}} - 2 \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{2}} \bar{\boldsymbol{G}}_{t} \right) - 1.1 \tilde{\boldsymbol{C}} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \eta^{2} \boldsymbol{I}_{\mathsf{rk}}$$

$$\stackrel{(d)}{\succeq} \bar{\boldsymbol{G}}_{t} + \eta \left((1 - \kappa_{d}) \boldsymbol{\Lambda}_{u_{1}} \bar{\boldsymbol{G}}_{t} + (1 - \kappa_{d}) \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{2}} - 2 \bar{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{u_{2}} \bar{\boldsymbol{G}}_{t} \right)$$

where we use $-2\Lambda_{u_2} \preceq \bar{V}_t^3 \left(I_r + \eta \bar{V}_t \right)^{-1} \preceq 2\Lambda_{u_2}$ in (a) and (c), and we use $\eta \|\Lambda\|_{\mathrm{F}}^2 r_s \ll \kappa_d \mathsf{rk}^{-\alpha} \frac{r_s}{d}$ in (b) and (d). By Proposition 34, we have $\frac{r_s}{d} I_{\mathsf{rk}} \preceq \bar{G}_{t+1} \preceq 1.1 I_{\mathsf{rk}}$, hence, the induction hypothesis holds. Therefore, we have the first item.

By using the first item and Proposition 34, we can write for the given time horizons in second and third items that

$$\bar{\boldsymbol{G}}_{t} \leq \left(\bar{\boldsymbol{G}}_{0} \exp\left(2(1+\kappa_{d})\eta t \boldsymbol{\Lambda}_{u_{1}}\right) \wedge \left(1+(1+\kappa_{d})^{2} \eta^{2} \boldsymbol{\Lambda}_{u_{1}}^{2}\right) \boldsymbol{I}_{\mathsf{rk}}\right) \\
\leq \begin{cases} \left(1.1 \bar{\boldsymbol{G}}_{0} \exp\left(2\eta t \boldsymbol{\Lambda}\right) \wedge \boldsymbol{I}_{\mathsf{rk}}\right) + 2\eta^{2} \boldsymbol{I}_{\mathsf{rk}} \\ 1.2 \left(\bar{\boldsymbol{G}}_{0} \exp\left(2\eta t \boldsymbol{\Lambda}\right) \wedge \boldsymbol{I}_{\mathsf{rk}}\right), \end{cases} (F.36)$$

where both upper bounds in (F.36) are valid and will be used in different parts of the proof. The third sub-items immediately follow from the first bound.

For $\alpha \in [0, 0.5)$, we have

$$\hat{\boldsymbol{T}}_{t,ii} \geq \frac{1}{3} \left(1 \wedge \hat{\boldsymbol{T}}_{0,ii} \exp\left(2\eta t \lambda_{i}\right) \right) \Rightarrow \boldsymbol{T}_{t,ii} \geq \frac{1}{4} \left(\kappa_{d} \wedge \boldsymbol{T}_{0,ii} \exp\left(2\eta t \lambda_{i}\right) \right) \\
\stackrel{(e)}{\Rightarrow} \boldsymbol{T}_{t,ii} \geq \frac{\kappa_{d}}{4} \left(1 \wedge \frac{r_{s}}{d} \exp\left(2\eta t \lambda_{i}\right) \right),$$

where we used $T_0 = \frac{\kappa_d r_s}{d} I_{rk}$ in (e). Therefore by (F.33) and the second bound in (F.36), we have for $j \leq t$

$$\frac{\bar{G}_{j,ii}}{T_{t,ii}} \leq \frac{1.2 \left(1 \wedge \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2 \frac{2.2r_s}{d} \exp\left(2\eta t \lambda_i\right)\right)}{0.25 \kappa_d \left(1 \wedge \frac{r_s}{d} \exp\left(2\eta t \lambda_i\right)\right)} \leq \frac{11}{\kappa_d} \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2.$$

On the other hand, by using the second bound in (F.36).

$$\frac{\eta \sum_{j=0}^{t-1} \bar{G}_{j,ii}}{T_{t,ii}} \leq \frac{11 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2}{\kappa_d} \frac{\eta \left(t \wedge \frac{r_s}{d} \sum_{j=0}^{t-1} \exp(2\eta j \lambda_i)\right)}{\left(1 \wedge \frac{r_s}{d} \exp(2\eta t \lambda_i)\right)}$$
$$\leq \frac{5.5 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2}{\kappa_d} \begin{cases} \frac{1}{\lambda_i}, & 2\eta t \leq \frac{\log \frac{d}{r_s}}{\lambda_i} \\ 2\eta t, & 2\eta t > \frac{\log \frac{d}{r_s}}{\lambda_i} \end{cases}$$

$$\leq \frac{5.5\left(1+\frac{1}{\sqrt{\varphi}}\right)^2}{\mathsf{K}_d}(2\eta t \vee r^{\alpha}).$$

Lastly, by using the first bound in (F.36), we get

$$\|\mathbf{\Lambda}\|_{\mathrm{F}}^2 - \operatorname{Tr}(\mathbf{\Lambda}\bar{\mathbf{G}}_t\mathbf{\Lambda}) \ge \sum_{i=1}^r \lambda_i^2 \left(1 - \frac{2.5\left(1 + \frac{1}{\sqrt{\varphi}}\right)^2 r_s}{d} \exp\left(2\eta t \lambda_i\right)\right)_+ - 2\eta^2 \|\mathbf{\Lambda}\|_{\mathrm{F}}^2.$$

For $\alpha > 0.5$, we have for $i \leq 2r_s \log^{\frac{1}{\alpha}} d$ and $d \geq \Omega_{r_s}(1)$,

$$\hat{\boldsymbol{T}}_{t,ii} \geq \frac{1}{3} \big(1 \wedge \hat{\boldsymbol{T}}_{0,ii} \exp\left(2\eta t \lambda_i\right) \big) \Rightarrow \boldsymbol{T}_{t,ii} \geq \frac{1}{4} \big(\kappa_d \wedge \boldsymbol{T}_{0,ii} \exp\left(2\eta t \lambda_i\right) \big).$$

Therefore, we have

$$T_{t,ii} \ge \kappa_d \begin{cases} 0.25 \left(1 \wedge \frac{r_s}{d} \exp\left(2\eta t \lambda_i\right)\right), & i \le 2r_s \log^{\frac{1}{\alpha}} d \\ \frac{r_s}{d}, & i > 2r_s \log^{\frac{1}{\alpha}} d. \end{cases}$$

On the other hand, for $\eta t \leq \frac{1}{2} r_s^{\alpha} \log(\frac{d \log^{1.5} d}{r_s})$ and $i > 2r_s \log^{\frac{1}{\alpha}} d$, we have for $d \geq \Omega(1)$.

$$\bar{\boldsymbol{G}}_{t,ii} \leq 1.2 \left(\bar{\boldsymbol{G}}_{0,ii} \exp\left(2\eta t \lambda_i\right) \wedge 1 \right) \leq 1.2 \left(\bar{\boldsymbol{G}}_{0,ii} \exp\left(\frac{r_s^{\alpha} \log\left(\frac{d \log d}{r_s}\right)}{2^{\alpha} r_s^{\alpha} \log d} \right) \wedge 1 \right) \leq 1.5 \bar{\boldsymbol{G}}_{0,ii}.$$

Therefore, for $\eta t \leq \frac{1}{2} r_s^{\alpha} \log(\frac{d \log^{1.5} d}{r_s})$,

$$\begin{split} \frac{\bar{\boldsymbol{G}}_{j,ii}}{\boldsymbol{T}_{t,ii}} &\leq \begin{cases} \frac{1.2 \left(1 \wedge \frac{5.5 r_u r_s}{d} \exp\left(2 \eta t \lambda_i\right)\right)}{0.25 \kappa_d \left(1 \wedge \frac{r_s}{d} \exp\left(2 \eta t \lambda_i\right)\right)}, & i \leq 2 r_s \log^{\frac{1}{\alpha}} d \\ \frac{d}{\kappa_d r_s} \frac{8.25 r_u r_s}{d}, & i > 2 r_s \log^{\frac{1}{\alpha}} d \end{cases} \\ &\leq \frac{26.4 r_u}{\kappa_d}. \end{split}$$

Moreover, for $\eta t \leq \frac{1}{2} r_s^{\alpha} \log(\frac{d \log^{1.5} d}{r_s})$,

$$\begin{split} \frac{\eta \sum_{j=0}^{t-1} \bar{G}_{j,ii}}{T_{t,ii}} &\leq \begin{cases} \frac{26.4r_u}{\kappa_d} \frac{\eta \left(t \wedge \frac{r_s}{d} \sum_{j=0}^{t-1} \exp(2\eta j \lambda_i) \right)}{\left(1 \wedge \frac{r_s}{d} \exp\left(2\eta t \lambda_i \right) \right)}, & i \leq 2r_s \log^{\frac{1}{\alpha}} d \\ \frac{d}{\kappa_d r_s} \frac{8.25r_u r_s}{d} \eta t, & i > 2r_s \log^{\frac{1}{\alpha}} d \end{cases} \\ &\leq \frac{13.2r_u}{\kappa_d} \begin{cases} \frac{1}{\lambda_i}, & 2\eta t \leq \frac{\log \frac{d}{r_s}}{\lambda_i} \text{ and } i \leq 2r_s \log^{\frac{1}{\alpha}} d \\ 2\eta t, & \text{otherwise} \end{cases} \\ &\leq \frac{13.2r_u}{\kappa_d} (2\eta t \vee (2r_s)^{\alpha} \log d) \leq \frac{15r_u}{\kappa_d} (2r_s)^{\alpha} \log d. \end{split}$$

Finally for $\eta t \leq \frac{1}{2} r_s^{\alpha} \log(\frac{d \log^{1.5} d}{r_s})$ and $d \geq \Omega_{r_s}(1)$, by using the first bound in (F.36),

$$\begin{split} &\|\boldsymbol{\Lambda}_{11}\|_{F}^{2} - \text{Tr}(\boldsymbol{\Lambda}_{11}\bar{\boldsymbol{G}}_{t,11}\boldsymbol{\Lambda}_{11}) \geq \sum_{i=1}^{r_{s}} \lambda_{i}^{2} \left(1 - \frac{12.1r_{s}}{d} \exp\left(2\eta t \lambda_{i}\right)\right)_{+} \\ &+ \sum_{i=r_{s}+1}^{2r_{s}} \lambda_{i}^{2} \left(1 - \frac{12.1r_{s}}{d} \exp\left(2\eta t \lambda_{i}\right)\right)_{+} + \sum_{i=2r_{s}+1}^{r_{u}} \lambda_{i}^{2} \left(1 - \frac{6.05r_{u}r_{s}}{d} \exp\left(2\eta t \lambda_{i}\right)\right)_{+} - 2\eta^{2} \|\boldsymbol{\Lambda}\|_{F}^{2} \\ &\geq \sum_{i=1}^{r_{s}} \lambda_{i}^{2} \left(1 - \frac{12.1r_{s}}{d} \exp\left(2\eta t \lambda_{i}\right)\right)_{+} + \left(1 - \frac{1}{\log d}\right) \sum_{i=r_{s}+1}^{2r_{s}} \lambda_{i}^{2} \end{split}$$

$$\begin{split} & + \left(1 - 6.05 r_u \log^{\frac{1.5}{\sqrt{2}}} d\left(\frac{r_s}{d}\right)^{1 - \frac{1}{\sqrt{2}}}\right) \sum_{i = 2r_s + 1}^{r_u} \lambda_i^2 - 2\eta^2 \|\mathbf{\Lambda}\|_{\mathrm{F}}^2 \\ & \geq \sum_{i = 1}^{r_s} \lambda_i^2 \left(1 - \frac{12.1 r_s}{d} \exp\left(2\eta t \lambda_i\right)\right)_+ + \left(1 - 6.05 r_u \log^{\frac{1.5}{\sqrt{2}}} d\left(\frac{r_s}{d}\right)^{1 - \frac{1}{\sqrt{2}}}\right) \sum_{i = r_s + 1}^{r_u} \lambda_i^2 \\ & - \frac{(r_s + 1)^{1 + 2\alpha}}{\log d} - 2\eta^2 \|\mathbf{\Lambda}\|_{\mathrm{F}}^2, \end{split}$$

where we used the bounds for t, d in (f).

F.5 Bounds for the second-order terms

We recall

$$R_{so}[G_{t}] = \frac{\eta^{2}}{16r_{s}} \mathbf{\Theta}^{\top} \mathbb{E}_{t} \left[\nabla_{St} \mathbf{L}_{t+1} \nabla_{St} \mathbf{L}_{t+1}^{\top} \right] \mathbf{\Theta}$$

$$- \frac{\eta^{2}}{16r_{s}} \mathbf{M}_{t} \mathbb{E}_{t} \left[\frac{\mathbf{\mathcal{P}}_{t+1}}{1 + c_{t+1}^{2}} \right] \mathbf{M}_{t}^{\top} - \frac{\eta^{3}}{32r_{s}^{3/2}} \operatorname{Sym} \left(\mathbf{\Theta}^{\top} \mathbb{E}_{t} \left[\frac{\nabla_{St} \mathbf{L}_{t+1} \mathbf{\mathcal{P}}_{t+1}}{1 + c_{t+1}^{2}} \right] \mathbf{M}_{t}^{\top} \right)$$

$$- \frac{\eta^{4}}{256r_{s}^{2}} \mathbf{\Theta}^{\top} \mathbb{E}_{t} \left[\frac{\nabla_{St} \mathbf{L}_{t+1} \mathbf{\mathcal{P}}_{t+1} \nabla_{St} \mathbf{L}_{t+1}^{\top}}{1 + c_{t+1}^{2}} \right] \mathbf{\Theta}.$$
(F.37)

Proposition 17. For $\eta \ll d^{-1/2}$, there exists a universal constant C > 0 such that

$$-C\left(\frac{\eta^2 d}{r_s} G_t + \eta^2 I_r\right) \leq R_{so}[G_t] \leq C\left(\frac{\eta^2 d}{r_s} G_t + \eta^2 I_r\right).$$

Proof. We bound each term in (F.37). In the following, v denotes a generic unit norm vector with proper dimensionality. For the first term,

$$\begin{aligned} \mathbf{\Theta}^{\top} \mathbb{E}_t \Big[\nabla_{\mathsf{S}t} \boldsymbol{L}_{t+1} \nabla_{\mathsf{S}t} \boldsymbol{L}_{t+1}^{\top} \Big] \mathbf{\Theta} \\ &= \mathbf{\Theta}^{\top} \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top} \right) \mathbb{E}_t \left[(y_{t+1} - \hat{y}_{t+1})^2 \| \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top} \right) \mathbf{\Theta}. \end{aligned}$$

We have

$$\mathbb{E}_t \left[(y_{t+1} - \hat{y}_{t+1})^2 \| \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \|_2^2 \langle \boldsymbol{v}, \boldsymbol{x}_{t+1} \rangle^2 \right] \leq Cr_s.$$

Therefore,

$$0 \leq \frac{\eta^2}{16r_s} \mathbf{\Theta}^{\top} \mathbb{E}_t \left[\nabla_{\mathsf{St}} \mathbf{L}_{t+1} \nabla_{\mathsf{St}} \mathbf{L}_{t+1}^{\top} \right] \mathbf{\Theta} \leq C \eta^2 (\mathbf{I}_r - \mathbf{G}_t).$$

For the second term,

$$egin{aligned} oldsymbol{M}_t \mathbb{E}_t \left[rac{oldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2}
ight] oldsymbol{M}_t^ op \ &= oldsymbol{M}_t \mathbb{E}_t \left[rac{\left(y_{t+1} - \hat{y}_{t+1}
ight)^2 \lVert \left(oldsymbol{I}_d - oldsymbol{W}_t oldsymbol{W}_t^ op
ight) oldsymbol{x}_{t+1} \lVert_2^2 oldsymbol{W}_t^ op oldsymbol{x}_{t+1} oldsymbol{x}_{t+1}^ op oldsymbol{W}_t^ op oldsymbol{X}_{t+1} oldsymbol{W}_t^ op oldsym$$

We have

$$\mathbb{E}_t \left[\frac{\left(y_{t+1} - \hat{y}_{t+1} \right)^2 \| \left(\mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^{\top} \right) \mathbf{x}_{t+1} \|_2^2 \left\langle \mathbf{v}, \mathbf{W}_t^{\top} \mathbf{x}_{t+1} \right\rangle^2}{1 + c_{t+1}^2} \right] \le Cd.$$

Therefore,

$$0 \leq \frac{\eta^2}{16r_s} \boldsymbol{M}_t \mathbb{E}_t \left[\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \boldsymbol{M}_t^\top \leq C \frac{\eta^2 d}{r_s} \boldsymbol{G}_t.$$

For the third term by using Proposition 22,

$$\frac{\eta^{3}}{32r_{s}^{3/2}}\operatorname{Sym}\left(\boldsymbol{\Theta}^{\top}\mathbb{E}_{t}\left[\frac{\nabla_{\operatorname{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}}{1+c_{t+1}^{2}}\right]\boldsymbol{M}_{t}^{\top}\right) \\
\leq C\left(\frac{\eta^{4}}{r_{s}^{2}d}\boldsymbol{\Theta}^{\top}\mathbb{E}_{t}\left[\frac{\nabla_{\operatorname{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}}{1+c_{t+1}^{2}}\right]\mathbb{E}_{t}\left[\frac{\boldsymbol{\mathcal{P}}_{t+1}\nabla_{\operatorname{St}}\boldsymbol{L}_{t+1}^{\top}}{1+c_{t+1}^{2}}\right]\boldsymbol{\Theta} + \frac{\eta^{2}d}{r_{s}}\boldsymbol{G}_{t}\right)$$

We have

$$\begin{split} \boldsymbol{\Theta}^{\top} & \mathbb{E}_t \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \\ &= \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \\ &\times \mathbb{E}_t \left[\frac{\left(y_{t+1} - \hat{y}_{t+1} \right)^3}{1 + c_{t+1}^2} \| \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top} \right) \boldsymbol{x}_{t+1} \|_2^2 \| \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \boldsymbol{W}_t \right]. \end{split}$$

Then, by using Cauchy-Schwartz inequality, we can show that

$$\left\| \mathbb{E}_t \left[\frac{\left(y_{t+1} - \hat{y}_{t+1} \right)^3}{1 + c_{t+1}^2} \| \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top \right) \boldsymbol{x}_{t+1} \|_2^2 \| \boldsymbol{W}_t^\top \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^\top \boldsymbol{W}_t \right] \right\|_2 \le C dr_s.$$

Therefore,

$$\frac{\eta^4}{r_s^2 d} \boldsymbol{\Theta}^\top \mathbb{E}_t \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \mathbb{E}_t \left[\frac{\boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] \boldsymbol{\Theta} \leq C \eta^4 d(\boldsymbol{I}_r - \boldsymbol{G}_t).$$

We get

$$\frac{\eta^3}{32r_s^{3/2}} \operatorname{Sym}\left(\boldsymbol{\Theta}^{\top} \mathbb{E}_t \left[\frac{\nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \boldsymbol{M}_t^{\top} \right) \preceq C \left(\eta^4 d(\boldsymbol{I}_r - \boldsymbol{G}_t) + \frac{\eta^2 d}{r_s} \boldsymbol{G}_t \right).$$

By repeating the argument with the lower bound in Proposition 22, we can also show

$$\frac{\eta^3}{32r_s^{3/2}}\mathrm{Sym}\left(\boldsymbol{\Theta}^{\top}\mathbb{E}_t\left[\frac{\nabla_{\mathrm{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}}{1+c_{t+1}^2}\right]\boldsymbol{M}_t^{\top}\right)\succeq -C\left(\eta^4d(\boldsymbol{I}_r-\boldsymbol{G}_t)+\frac{\eta^2d}{r_s}\boldsymbol{G}_t\right).$$

For the last term, we write

$$\begin{split} & \boldsymbol{\Theta}^{\top} \mathbb{E}_{t} \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top}}{1 + c_{t+1}^{2}} \right] \boldsymbol{\Theta} \\ & = \boldsymbol{\Theta}^{\top} \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \\ & \times \mathbb{E}_{t} \left[\frac{\left(y_{t+1} - \hat{y}_{t+1} \right)^{4}}{1 + c_{t+1}^{2}} \| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{4} \| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{\Theta}. \end{split}$$

We have

$$\mathbb{E}_{t} \left[\frac{\left(y_{t+1} - \hat{y}_{t+1} \right)^{4}}{1 + c_{t+1}^{2}} \| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{4} \| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \|_{2}^{2} \left\langle \boldsymbol{v}, \boldsymbol{x}_{t+1} \right\rangle^{2} \right] \leq C d r_{s}^{2}.$$

Therefore,

$$0 \leq \frac{\eta^4}{r_s^2} \boldsymbol{\Theta}^\top \mathbb{E}_t \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] \boldsymbol{\Theta} \leq C \eta^4 d(\boldsymbol{I}_r - \boldsymbol{G}_t).$$

By using $G_t \succeq 0$ and $\eta \ll d^{-1/2}$, the result follows.

F.6 Noise characterization

To prove the noise concentration bound for both the heavy-tailed and light-tailed cases simultaneously, we introduce some new notation. Specifically, we define the submatrix notation:

$$\boldsymbol{\Theta} \eqqcolon \begin{bmatrix} \boldsymbol{\Theta}_1 & \boldsymbol{\Theta}_2 \end{bmatrix} \ \ \text{and} \ \ \boldsymbol{M}_t \eqqcolon \begin{bmatrix} \boldsymbol{M}_{t,1} \\ \boldsymbol{M}_{t,2} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Theta}_1^\top \boldsymbol{W}_t \\ \boldsymbol{\Theta}_2^\top \boldsymbol{W}_t \end{bmatrix},$$

where $\Theta_1 \in \mathbb{R}^{d \times r_u}$ and $M_{t,1} \in \mathbb{R}^{r_u \times r_s}$. We note that $G_{t,11} = M_{t,1} M_{t,1}^{\top}$. To unify the treatment of the heavy-tailed and light-tailed cases, we use the following notation to represent both cases:

$$\Theta \coloneqq \{\boldsymbol{\Theta}, \boldsymbol{\Theta}_1\} \quad \mathsf{M}_t \coloneqq \{\boldsymbol{M}_t, \boldsymbol{M}_{t,1}\}.$$

With the new notation, we have

$$\begin{split} \frac{\eta/2}{\sqrt{r_s}} \mathbf{v}_{t+1} &= \frac{\eta/2}{\sqrt{r_s}} \mathrm{Sym} \left(\Theta^\top \left(\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} - \mathbb{E}_t \left[\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \right] \right) \mathsf{M}_t^\top \right) \\ &- \frac{\eta^2}{16r_s} \mathsf{M}_t \left(\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[\frac{\boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \right) \mathsf{M}_t^\top \\ &+ \frac{\eta^2}{16r_s} \Theta^\top \left(\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top - \mathbb{E}_t \left[\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top \right] \right) \Theta \\ &- \frac{\eta^3}{32r_s^{3/2}} \mathrm{Sym} \left(\Theta^\top \left(\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1}}{1 + c_{t+1}^2} \right] \right) \mathsf{M}_t^\top \right) \\ &- \frac{\eta^4}{256r_s^2} \Theta^\top \left(\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top}{1 + c_{t+1}^2} - \mathbb{E}_t \left[\frac{\nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] \right) \Theta. \end{split}$$

For $\mathsf{rk} \in \{r, r_u\}$ and $T_1, T_2 \in \mathbb{R}^{\mathsf{rk} \times \mathsf{rk}}$ be a deterministic symmetric positive definite matrices, we define

$$\begin{split} &\mathcal{A}_{t+1}\left(\boldsymbol{T}_{1},\boldsymbol{T}_{2}\right) \equiv \left\{ \left\|\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}\mathsf{M}_{t}^{\top}\boldsymbol{T}_{2}\right\|_{2} \leq \frac{L^{2}}{2}\sqrt{\mathsf{Tr}\big(\boldsymbol{T}_{1}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)\mathsf{Tr}(\boldsymbol{T}_{2}^{2}\mathsf{G}_{t})}\right\} \\ &\mathcal{B}_{t+1}\left(\boldsymbol{T}_{1},\boldsymbol{T}_{2}\right) \equiv \left\{ \left\|\boldsymbol{T}_{1}\mathsf{M}_{t}\boldsymbol{\mathcal{P}}_{t+1}\mathsf{M}_{t}^{\top}\boldsymbol{T}_{2}\right\|_{2} \leq \frac{L^{4}d}{2}\sqrt{\mathsf{Tr}(\boldsymbol{T}_{1}^{2}\mathsf{G}_{t})\mathsf{Tr}(\boldsymbol{T}_{2}^{2}\mathsf{G}_{t})}\right\} \\ &\mathcal{C}_{t+1}\left(\boldsymbol{T}_{1},\boldsymbol{T}_{2}\right) \equiv \left\{ \left\|\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}^{\top}\boldsymbol{\Theta}\boldsymbol{T}_{2}\right\|_{2} \leq \frac{L^{4}r_{s}}{2}\sqrt{\mathsf{Tr}\big(\boldsymbol{T}_{1}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)\mathsf{Tr}\big(\boldsymbol{T}_{2}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)}\right\} \\ &\mathcal{D}_{t+1}\left(\boldsymbol{T}_{1},\boldsymbol{T}_{2}\right) \equiv \left\{ \left\|\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}\mathsf{M}_{t}^{\top}\boldsymbol{T}_{2}\right\|_{2} \leq \frac{L^{6}dr_{s}}{2}\sqrt{\mathsf{Tr}\big(\boldsymbol{T}_{1}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)\mathsf{Tr}\big(\boldsymbol{T}_{2}^{2}\mathsf{G}_{t}\big)}\right\} \\ &\mathcal{F}_{t+1}\left(\boldsymbol{T}_{1},\boldsymbol{T}_{2}\right) \equiv \left\{ \left\|\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}\nabla_{\mathsf{S}t}\boldsymbol{L}_{t+1}^{\top}\boldsymbol{\Theta}\boldsymbol{T}_{2}\right\|_{2} \leq \frac{L^{8}dr_{s}^{2}}{2}\sqrt{\mathsf{Tr}\big(\boldsymbol{T}_{1}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)\mathsf{Tr}\big(\boldsymbol{T}_{2}^{2}(\boldsymbol{I}_{\mathsf{rk}}-\mathsf{G}_{t})\big)}\right\}. \end{split}$$

We start with the following statement:

Proposition 18. Let $e_{t+1} := (y_{t+1} - \hat{y}_{t+1})$. There exists a universal constant C > 0 such that for $L \ge 2e(\sqrt{8} + \sigma)$, the following statements hold:

1. We have
$$\mathbb{P}_t \left[\mathcal{A}_{t+1} \left(T_1, T_2 \right) \cap \left\{ |e_{t+1}| \leq L \right\} \right] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8}+\sigma)}}$$
. Moreover,
$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(T_1 \Theta^\top \nabla_{St} \boldsymbol{L}_{t+1} \mathsf{M}_t^\top \boldsymbol{T}_2 \right) \right)^2 \right] \\ \leq C \left(Tr(T_2^2 \mathsf{G}_t) T_1 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) T_1 + Tr \left(T_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \right) T_2 \mathsf{G}_t T_2 \right).$$

2. We have
$$\mathbb{P}_t \left[\mathcal{B}_{t+1} \left(T_1, T_2 \right) \cap \left\{ |e_{t+1}| \leq L \right\} \right] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8}+\sigma)}}$$
. Moreover,
$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(T_1 \mathsf{M}_t \mathcal{P}_{t+1} \mathsf{M}_t^\top T_2 \right) \right)^2 \right] \\ \prec Cd^2 \left(T_t (T_2^2 \mathsf{G}_t) T_1 \mathsf{G}_t T_1 + T_t (T_1^2 \mathsf{G}_t) T_2 \mathsf{G}_t T_2 \right).$$

3. We have
$$\mathbb{P}_t \left[\mathcal{C}_{t+1} \left(\mathbf{T}_1, \mathbf{T}_2 \right) \cap \left\{ |e_{t+1}| \leq L \right\} \right] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8}+\sigma)}}$$
. Moreover,
$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(\mathbf{T}_1 \Theta^\top \nabla_{St} \mathbf{L}_{t+1} \nabla_{St} \mathbf{L}_{t+1}^\top \mathbf{T}_2 \right) \right)^2 \right]$$

$$\leq Cr_s^2 \left(Tr \left(\mathbf{T}_2^2 (\mathbf{I}_{\mathsf{rk}} - \mathsf{G}_t) \right) \mathbf{T}_1 (\mathbf{I}_{\mathsf{rk}} - \mathsf{G}_t) \mathbf{T}_1 + Tr \left(\mathbf{T}_1^2 (\mathbf{I}_{\mathsf{rk}} - \mathsf{G}_t) \right) \mathbf{T}_2 (\mathbf{I}_{\mathsf{rk}} - \mathsf{G}_t) \mathbf{T}_2 \right).$$

4. We have $\mathbb{P}_t \left[\mathcal{D}_{t+1} \left(T_1, T_2 \right) \cap \{ |e_{t+1}| \leq L \} \right] \geq 1 - 2e^{\frac{-L/e}{(7/2+\sigma)}}$. Moreover,

$$\mathbb{E}_{t} \left[\left(\operatorname{Sym} \left(\boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \nabla_{St} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \boldsymbol{\mathsf{M}}_{t}^{\top} \boldsymbol{T}_{2} \right) \right)^{2} \right]$$

$$\leq C d^{2} r_{s}^{2} \left(Tr(\boldsymbol{T}_{2}^{2} \boldsymbol{\mathsf{G}}_{t}) \boldsymbol{T}_{1} (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_{t}) \boldsymbol{T}_{1} + Tr(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_{t})) \boldsymbol{T}_{2} \boldsymbol{\mathsf{G}}_{t} \boldsymbol{T}_{2} \right).$$

5. We have $\mathbb{P}_t \left[\mathcal{F}_{t+1} \left(T_1, T_2 \right) \cap \{ |e_{t+1}| \leq L \} \right] \geq 1 - 2e^{\frac{-L/e}{(4\sqrt{2}+\sigma)}}$. Moreover,

$$\mathbb{E}_{t} \left[\left(\operatorname{Sym} \left(\boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \nabla_{St} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{St} \boldsymbol{L}_{t+1}^{\top} \boldsymbol{\Theta} \boldsymbol{T}_{2} \right) \right)^{2} \right]$$

$$\leq C d^{2} r_{s}^{4} \left(Tr \left(\boldsymbol{T}_{2}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \boldsymbol{T}_{1} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \boldsymbol{T}_{1} + Tr \left(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \boldsymbol{T}_{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \boldsymbol{T}_{2} \right).$$

Proof. First, we derive a concentration bound for $|e_{t+1}|$. By Corollary 6 and Proposition 33 we have

$$\mathbb{E}_{t}[|e_{t+1}|^{p}] \leq (\mathbb{E}_{t}[|y_{t+1}|^{p}]^{\frac{1}{p}} + \mathbb{E}_{t}[|\hat{y}_{t+1}|^{p}]^{\frac{1}{p}})^{p} \leq (\sqrt{8} + \sigma)^{p}p^{p} \text{ for } p \geq 2,$$

which implies $\mathbb{P}_t[|e_{t+1}| \geq u] \leq e^{\frac{-u/e}{(\sqrt{8}+\sigma)}}$ for $u \geq 2e(\sqrt{8}+\sigma)$. In the following, we prove each item separately.

First item. We define

$$\boldsymbol{T}_1 \boldsymbol{\Theta}^\top \nabla_{\mathrm{St}} \boldsymbol{L}_{t+1} \mathsf{M}_t^\top \boldsymbol{T}_2 = \underbrace{\boldsymbol{T}_1 \boldsymbol{\Theta}^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) e_{t+1} \boldsymbol{x}_{t+1}}_{\coloneqq \boldsymbol{u}_{t+1}} \underbrace{\boldsymbol{x}_{t+1}^\top \boldsymbol{W}_t \boldsymbol{W}_t^\top \boldsymbol{\Theta} \boldsymbol{T}_2}_{\coloneqq \boldsymbol{v}_{t+1}^\top}.$$

For u, L > 0

$$\begin{split} & \mathbb{P}_t \Big[\big\| \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \geq uL \sqrt{\mathrm{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \mathrm{Tr} (\boldsymbol{T}_2^2 \mathsf{G}_t)} \ \, \text{or} \ \, |\boldsymbol{e}_{t+1}| \geq L \Big] \\ & \leq \mathbb{P}_t \Big[\big\| \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \geq uL \sqrt{\mathrm{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \mathrm{Tr} (\boldsymbol{T}_2^2 \mathsf{G}_t)} \ \, \text{and} \ \, |\boldsymbol{e}_{t+1}| \leq L \Big] + \mathbb{P}_t \Big[|\boldsymbol{e}_{t+1}| \geq L \Big] \\ & \leq \mathbb{P}_t \Big[\big\| \mathbb{1}_{|\boldsymbol{e}_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \geq uL \sqrt{\mathrm{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \mathrm{Tr} (\boldsymbol{T}_2^2 \mathsf{G}_t)} \Big] + \mathbb{P}_t \Big[|\boldsymbol{e}_{t+1}| \geq L \Big]. \end{split}$$

We have for $p \geq 2$

$$\mathbb{E}_{t} \left[\left\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right\|_{2}^{p} \right]$$

$$\leq L^{p} \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \right\|_{2}^{p} \right] \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{p} \right]$$

$$\stackrel{(a)}{\leq} L^{p} \left(\frac{p}{2} \right)^{p} \left(3 \text{Tr} \left(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \text{Tr} (\boldsymbol{T}_{2}^{2} \mathsf{G}_{t}) \right)^{\frac{p}{2}},$$

where we used Corollary 7 in (a). By Proposition 33, we have for $u \ge 2e$

$$\mathbb{P}_t \Big[\big\| \mathbb{1}_{|e_{t+1}| \le t} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \ge uL \sqrt{\text{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \text{Tr} (\boldsymbol{T}_2^2 \mathsf{G}_t)} \Big] \le e^{-\frac{u}{e}}.$$

By choosing $u = \frac{L}{2}$, we have the probability bound.

For the variance bound, we have

$$\mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{T}_1 \boldsymbol{\Theta}^\top \nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathsf{M}}_t^\top \boldsymbol{T}_2 \right)^2 \right] = \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \right)^2 \right].$$

By using Proposition 22, we have

$$\mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right)^2 \right]$$

$$\leq \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \mathbb{E}_{t} \left[e_{t+1}^{2} \| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{\Theta} \boldsymbol{T}_{1} \\ + \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \mathbb{E}_{t} \left[e_{t+1}^{2} \| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{T}_{2} \\ \leq C \Big(\operatorname{Tr}(\boldsymbol{T}_{2}^{2} \mathsf{G}_{t}) \boldsymbol{T}_{1} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \boldsymbol{T}_{1} + \operatorname{Tr} \Big(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \Big) \boldsymbol{T}_{2} \mathsf{G}_{t} \boldsymbol{T}_{2} \Big),$$

where we used the Cauchy-Schwartz inequality in (b).

Second item. We define

$$\boldsymbol{T}_{1}\mathsf{M}_{t}\boldsymbol{\mathcal{P}}_{t+1}\mathsf{M}_{t}^{\top}\boldsymbol{T}_{2} = \underbrace{e_{t+1}^{2}\|\left(\boldsymbol{I}_{d}-\boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top}\right)\boldsymbol{x}_{t+1}\|_{2}^{2}\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top}\boldsymbol{x}_{t+1}}_{:=\boldsymbol{u}_{t+1}}\underbrace{\boldsymbol{x}_{t+1}^{\top}\boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top}\boldsymbol{\Theta}\boldsymbol{T}_{2}}_{:=\boldsymbol{v}_{t+1}^{\top}}.$$

We have for $p \ge 2$

$$\begin{split} & \mathbb{E}_{t} \left[\left\| \mathbb{1}_{\left| e_{t+1} \right| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right\|_{2}^{p} \right] \\ & \leq L^{2p} \mathbb{E}_{t} \left[\left\| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \right] \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \right]^{\frac{1}{2}} \\ & \times \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \right]^{\frac{1}{2}} \\ & \stackrel{(c)}{\leq} L^{2p} p^{2p} \left(3d \sqrt{\text{Tr}(\boldsymbol{T}_{1}^{2} \boldsymbol{G}_{t}) \text{Tr}(\boldsymbol{T}_{2}^{2} \boldsymbol{G}_{t})} \right)^{p}, \end{split}$$

where we use Corollary 7 in (c). By Proposition 33, we have for $u \ge (2e)^2$

$$\mathbb{P}_t \Big[\big\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \geq u L^2 3d \sqrt{\text{Tr}(\boldsymbol{T}_1^2 \mathsf{G}_t) \text{Tr}(\boldsymbol{T}_2^2 \mathsf{G}_t)} \Big] \leq e^{-\frac{u^{1/2}}{e}}.$$

By choosing $u=\frac{L^2}{6}$, we have the probability bound. For the variance bound, we have

$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(\boldsymbol{T}_1 \mathsf{M}_t \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_t^\top \boldsymbol{T}_2 \right) \right)^2 \right] = \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \right)^2 \right].$$

By using Proposition 22, we have

$$\mathbb{E}_{t} \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right)^{2} \right]$$

$$\preceq \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}$$

$$\times \mathbb{E}_{t} \left[e_{t+1}^{4} \| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \|_{2}^{4} \| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{T}_{1}$$

$$+ \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}$$

$$\times \mathbb{E}_{t} \left[e_{t+1}^{4} \| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \|_{2}^{4} \| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{\Theta} \boldsymbol{T}_{2}$$

$$\stackrel{(d)}{\leq} C d^{2} \left(\operatorname{Tr} (\boldsymbol{T}_{2}^{2} \boldsymbol{G}_{t}) \boldsymbol{T}_{1} \boldsymbol{G}_{t} \boldsymbol{T}_{1} + \operatorname{Tr} (\boldsymbol{T}_{1}^{2} \boldsymbol{G}_{t}) \boldsymbol{T}_{2} \boldsymbol{G}_{t} \boldsymbol{T}_{2} \right),$$

where we use the Cauchy-Schwartz inequality in (d).

Third item. We define

$$T_1 \Theta^ op
abla_{ ext{St}} oldsymbol{L}_{t+1}
abla_{ ext{St}} oldsymbol{L}_{t+1}^ op oldsymbol{T}_2 = \underbrace{e_{t+1}^2 \| oldsymbol{W}_t^ op oldsymbol{x}_{t+1} \|_2^2 oldsymbol{T}_1 oldsymbol{\Theta}^ op (oldsymbol{I}_d - oldsymbol{W}_t oldsymbol{W}_t^ op) oldsymbol{x}_{t+1}}_{:=oldsymbol{u}_{t+1}} \underbrace{oldsymbol{x}_{t+1}^ op (oldsymbol{I}_d - oldsymbol{W}_t oldsymbol{W}_t^ op) oldsymbol{x}_{t+1}}_{:=oldsymbol{u}_{t+1}}.$$

We have for $p \geq 2$

$$\begin{split} & \mathbb{E}_t \left[\left\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right\|_2^p \right] \\ & \leq L^{2p} \mathbb{E}_t \left[\left\| \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \right\|_2^{2p} \right] \mathbb{E}_t \left[\left\| \boldsymbol{T}_1 \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \boldsymbol{x}_{t+1} \right\|_2^{2p} \right]^{\frac{1}{2}} \\ & \times \mathbb{E}_t \left[\left\| \boldsymbol{T}_2 \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \boldsymbol{x}_{t+1} \right\|_2^{2p} \right]^{\frac{1}{2}} \end{split}$$

$$\overset{(e)}{\leq} L^{2p} p^{2p} \left(3r_s \sqrt{\text{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \text{Tr} \big(\boldsymbol{T}_2^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big)} \right)^p,$$

where we use Corollary 7 in (e). By Proposition 33, we have for $u \ge (2e)^2$

$$\mathbb{P}_t \Big[\big\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \big\|_2 \geq u L^2 3 r_s \sqrt{\mathrm{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \mathrm{Tr} \big(\boldsymbol{T}_2^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big)} \Big] \leq e^{-\frac{u^{1/2}}{e}}.$$

By choosing $u = \frac{L^2}{6}$, we have the probability bound. For the variance bound, we have

$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(\boldsymbol{T}_1 \boldsymbol{\Theta}^\top \nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} \nabla_{\operatorname{St}} \boldsymbol{L}_{t+1}^\top \boldsymbol{\Theta} \boldsymbol{T}_2 \right) \right)^2 \right] = \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \right)^2 \right].$$

By using Proposition 22, we have

$$\begin{split} & \mathbb{E}_{t} \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right)^{2} \right] \\ & \preceq \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \\ & \times \mathbb{E}_{t} \big[\boldsymbol{e}_{t+1}^{4} \| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{4} \| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{\Theta} \boldsymbol{T}_{1} \\ & + \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \\ & \times \mathbb{E}_{t} \big[\boldsymbol{e}_{t+1}^{4} \| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{4} \| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \|_{2}^{2} \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \big] \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{\Theta} \boldsymbol{T}_{2} \\ & \stackrel{(f)}{\preceq} C r_{s}^{2} \left(\operatorname{Tr} \left(\boldsymbol{T}_{2}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \boldsymbol{T}_{1} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \boldsymbol{T}_{1} + \operatorname{Tr} \left(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \boldsymbol{T}_{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \boldsymbol{T}_{2} \right), \end{split}$$

where we used the Cauchy-Schwartz inequality in (f).

Fourth item. We define

$$T_1 \Theta^\top \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_t^\top T_2 \\ = \underbrace{e_{t+1}^3 \| (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{x}_{t+1} \|_2^2 \| \boldsymbol{W}_t^\top \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{T}_1 \Theta^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{x}_{t+1}}_{:= \boldsymbol{u}_{t+1}} \underbrace{\boldsymbol{x}_{t+1}^\top \boldsymbol{W}_t \boldsymbol{W}_t^\top \Theta \boldsymbol{T}_2}_{:= \boldsymbol{u}_{t+1}}.$$

We have for $p \ge 2$

$$\begin{split} & \mathbb{E}_{t} \left[\left\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right\|_{2}^{p} \right] \\ & \leq L^{3p} \mathbb{E}_{t} \left[\left\| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \right\|_{2}^{p} \right] \\ & \times \mathbb{E}_{t} \left[\left\| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \left\| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{p} \right] \\ & \leq L^{3p} (2p)^{p} (\sqrt{3}d)^{p} p^{\frac{p}{2}} \left(\sqrt{3} \text{Tr} \left(\boldsymbol{T}_{1}^{2} \left(\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t} \right) \right) \right)^{\frac{p}{2}} \left(2p)^{p} (\sqrt{3}r_{s})^{p} p^{\frac{p}{2}} \left(\sqrt{3} \text{Tr} \left(\boldsymbol{T}_{2}^{2} \mathsf{G}_{t} \right) \right)^{\frac{p}{2}} \\ & = L^{3p} (12\sqrt{3})^{p} p^{3p} \left(dr_{s} \sqrt{\text{Tr} \left(\boldsymbol{T}_{1}^{2} \left(\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t} \right) \right) \text{Tr} \left(\boldsymbol{T}_{2}^{2} \mathsf{G}_{t} \right)} \right)^{p}. \end{split}$$

By Proposition 33, we have for $u \ge (2e)^3$

$$\mathbb{P}_t \Big[\big\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \big\|_2 \geq uL^3 12\sqrt{3} dr_s \sqrt{\mathrm{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \mathrm{Tr} \big(\boldsymbol{T}_2^2 \mathsf{G}_t \big)} \Big] \leq e^{-\frac{u^{1/3}}{e}}.$$

By choosing $u = \frac{L^3}{24\sqrt{3}}$, we have the probability bound. For the variance bound, we have

$$\mathbb{E}_t \left[\left(\operatorname{Sym} \left(\boldsymbol{T}_1 \boldsymbol{\Theta}^\top \nabla_{\operatorname{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \boldsymbol{\mathsf{M}}_t^\top \boldsymbol{T}_2 \right) \right)^2 \right] = \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \right)^2 \right].$$

By using Proposition 22, we have

$$\begin{split} & \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right)^2 \right] \\ & \leq \boldsymbol{T}_1 \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \\ & \times \mathbb{E}_t \left[e_{t+1}^6 \| (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \boldsymbol{x}_{t+1} \|_2^4 \| \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \|_2^4 \| \boldsymbol{T}_2 \boldsymbol{\Theta}^{\top} \boldsymbol{W}_t \boldsymbol{W}_t^{\top} \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^{\top} \right] \\ & \times (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^{\top}) \boldsymbol{\Theta} \boldsymbol{T}_1 \end{split}$$

$$\begin{split} &+ \boldsymbol{T}_{2}\boldsymbol{\Theta}^{\top}\boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top} \\ &\times \mathbb{E}_{t}\left[e_{t+1}^{6}\|(\boldsymbol{I}_{d} - \boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top})\boldsymbol{x}_{t+1}\|_{2}^{4}\|\boldsymbol{W}_{t}^{\top}\boldsymbol{x}_{t+1}\|_{2}^{4}\|\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}(\boldsymbol{I}_{d} - \boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top})\boldsymbol{x}_{t+1}\|_{2}^{2}\boldsymbol{x}_{t+1}\boldsymbol{x}_{t+1}^{\top}\right] \\ &\times \boldsymbol{W}_{t}\boldsymbol{W}_{t}^{\top}\boldsymbol{\Theta}\boldsymbol{T}_{2} \\ &\leq Cd^{2}r_{s}^{2}\Big(\mathrm{Tr}(\boldsymbol{T}_{2}^{2}\mathsf{G}_{t})\boldsymbol{T}_{1}(\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t})\boldsymbol{T}_{1} + \mathrm{Tr}\big(\boldsymbol{T}_{1}^{2}(\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t})\big)\boldsymbol{T}_{2}\mathsf{G}_{t}\boldsymbol{T}_{2}\Big). \end{split}$$

Fifth item. We define

$$\begin{split} \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \boldsymbol{\Theta} \boldsymbol{T}_{2} \\ &= \underbrace{e_{t+1}^{4} \| (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \|_{2}^{2} \| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \|_{2}^{4} \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1}}_{::= \boldsymbol{u}_{t+1}} \\ &\times \underbrace{\boldsymbol{x}_{t+1}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{\Theta} \boldsymbol{T}_{2}}_{::= \boldsymbol{v}_{t+1}^{\top}}. \end{split}$$

We have for $p \ge 2$

$$\begin{split} &\mathbb{E}_{t} \left[\left\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^{\top} \right\|_{2}^{p} \right] \\ &\leq L^{4p} \mathbb{E}_{t} \left[\left\| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \right\|_{2}^{2p} \left\| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{4p} \\ & \qquad \qquad \times \left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \right\|_{2}^{p} \left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \right\|_{2}^{p} \right] \\ &\leq L^{4p} \mathbb{E}_{t} \left[\left\| \left(\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top} \right) \boldsymbol{x}_{t+1} \right\|_{2}^{4p} \right]^{\frac{1}{2}} \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{1} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \right\|_{2}^{4p} \right]^{\frac{1}{4}} \\ & \qquad \times \mathbb{E}_{t} \left[\left\| \boldsymbol{T}_{2} \boldsymbol{\Theta}^{\top} (\boldsymbol{I}_{d} - \boldsymbol{W}_{t} \boldsymbol{W}_{t}^{\top}) \boldsymbol{x}_{t+1} \right\|_{2}^{4p} \right]^{\frac{1}{4}} \mathbb{E}_{t} \left[\left\| \boldsymbol{W}_{t}^{\top} \boldsymbol{x}_{t+1} \right\|_{2}^{4p} \right] \\ &\leq L^{4p} (2p)^{p} (\sqrt{3}d)^{p} (2p)^{p} \left(3 \text{Tr} \left(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \text{Tr} \left(\boldsymbol{T}_{2}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \right)^{\frac{p}{2}} (2p)^{2p} (\sqrt{3}r_{s})^{2p} \\ &= L^{4p} (2\sqrt{3})^{4p} p^{4p} \left(dr_{s}^{2} \sqrt{\text{Tr} \left(\boldsymbol{T}_{1}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right) \text{Tr} \left(\boldsymbol{T}_{2}^{2} (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_{t}) \right)} \right)^{p}. \end{split}$$

By Proposition 33, we have for $u \ge (2e)^4$

$$\mathbb{P}_t \Big[\Big\| \mathbb{1}_{|e_{t+1}| \leq L} \boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \Big\|_2 \geq u L^4 (2\sqrt{3})^4 dr_s^2 \sqrt{\text{Tr} \big(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big) \text{Tr} \big(\boldsymbol{T}_2^2 (\boldsymbol{I}_{\mathsf{rk}} - \mathsf{G}_t) \big)} \Big] \leq 2e^{-\frac{u^{1/4}}{e}}.$$

By choosing $u=\frac{L^4}{2(2\sqrt{3})^4}$, we have the probability bound. For the variance bound, we have

$$\mathbb{E}_{t}\left[\left(\operatorname{Sym}\left(\boldsymbol{T}_{1}\boldsymbol{\Theta}^{\top}\nabla_{\operatorname{St}}\boldsymbol{L}_{t+1}\boldsymbol{\mathcal{P}}_{t+1}\nabla_{\operatorname{St}}\boldsymbol{L}_{t+1}^{\top}\boldsymbol{\Theta}\boldsymbol{T}_{2}\right)\right)^{2}\right]=\mathbb{E}_{t}\left[\operatorname{Sym}\left(\boldsymbol{u}_{t+1}\boldsymbol{v}_{t+1}^{\top}\right)^{2}\right].$$

By using Proposition 22, we have

$$\begin{split} & \mathbb{E}_t \left[\operatorname{Sym} \left(\boldsymbol{u}_{t+1} \boldsymbol{v}_{t+1}^\top \right)^2 \right] \\ & \preceq \boldsymbol{T}_1 \boldsymbol{\Theta}^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \\ & \times \mathbb{E}_t \left[e_{t+1}^8 \| \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top \right) \boldsymbol{x}_{t+1} \|_2^4 \| \boldsymbol{W}_t^\top \boldsymbol{x}_{t+1} \|_2^8 \| \boldsymbol{T}_2 \boldsymbol{\Theta}^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^\top \right] \\ & \times (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{\Theta} \boldsymbol{T}_1 \\ & + \boldsymbol{T}_2 \boldsymbol{\Theta}^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \\ & \times \mathbb{E}_t \left[e_{t+1}^8 \| \left(\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top \right) \boldsymbol{x}_{t+1} \|_2^4 \| \boldsymbol{W}_t^\top \boldsymbol{x}_{t+1} \|_2^8 \| \boldsymbol{T}_1 \boldsymbol{\Theta}^\top (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{x}_{t+1} \|_2^2 \boldsymbol{x}_{t+1} \boldsymbol{x}_{t+1}^\top \right] \\ & \times (\boldsymbol{I}_d - \boldsymbol{W}_t \boldsymbol{W}_t^\top) \boldsymbol{\Theta} \boldsymbol{T}_2 \\ & \preceq C d^2 r_s^4 \Big(\operatorname{Tr} \left(\boldsymbol{T}_2^2 (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_t) \right) \boldsymbol{T}_1 (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_t) \boldsymbol{T}_1 + \operatorname{Tr} \left(\boldsymbol{T}_1^2 (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_t) \right) \boldsymbol{T}_2 (\boldsymbol{I}_{\mathsf{rk}} - \boldsymbol{\mathsf{G}}_t) \boldsymbol{T}_2 \Big). \end{split}$$

By recalling the definitions $\{T_t\}_{t\in\mathbb{N}}$, Λ , Λ_{11} in Sections F.2 and F.3, we define the event:

$$\mathcal{A}_{t+1} := \begin{cases} \mathcal{A}_{t+1} \left(\mathbf{T}_{t}^{\frac{-1}{2}}, \mathbf{T}_{t}^{\frac{-1}{2}} \right) \cap \mathcal{A}_{t+1} \left(\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{T}_{t}^{\frac{-1}{2}}, \mathbf{T}_{t}^{\frac{-1}{2}} \mathbf{\Lambda}^{\frac{1}{2}} \right) \cap \{e_{t+1} \leq L\}, & \alpha \in [0, 0.5) \\ \mathcal{A}_{t+1} \left(\mathbf{T}_{t}^{\frac{-1}{2}}, \mathbf{T}_{t}^{\frac{-1}{2}} \right) \cap \mathcal{A}_{t+1} \left(\mathbf{\Lambda}_{11}^{\frac{1}{2}} \mathbf{T}_{t}^{\frac{-1}{2}}, \mathbf{T}_{t}^{\frac{-1}{2}} \mathbf{\Lambda}_{11}^{\frac{1}{2}} \right) \cap \{e_{t+1} \leq L\}, & \alpha > 0.5. \end{cases}$$

We define the events \mathcal{B}_{t+1} , \mathcal{C}_{t+1} , \mathcal{D}_{t+1} , and \mathcal{F}_{t+1} in the same way. Based on these events, we define the clipped versions of the noise matrices:

$$\begin{split} \mathsf{A}_{t+1} &\coloneqq \operatorname{Sym}\left(\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{A}_{t+1}} - \mathbb{E}_{t} \left[\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{A}_{t+1}}\right]\right) \\ \mathsf{B}_{t+1} &\coloneqq \frac{\mathsf{M}_{t} \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{B}_{t+1}}}{1 + c_{t+1}^{2}} - \mathbb{E}_{t} \left[\frac{\mathsf{M}_{t} \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{B}_{t+1}}}{1 + c_{t+1}^{2}}\right] \\ \mathsf{C}_{t+1} &\coloneqq \Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \Theta \mathbb{1}_{\mathcal{C}_{t+1}} - \mathbb{E}_{t} \left[\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \Theta \mathbb{1}_{\mathcal{C}_{t+1}}\right] \\ \mathsf{D}_{t+1} &\coloneqq \operatorname{Sym}\left(\frac{\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{D}_{t+1}}}{1 + c_{t+1}^{2}} - \mathbb{E}_{t} \left[\frac{\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \mathsf{M}_{t}^{\top} \mathbb{1}_{\mathcal{D}_{t+1}}}{1 + c_{t+1}^{2}}\right]\right) \\ \mathsf{F}_{t+1} &\coloneqq \frac{\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \Theta \mathbb{1}_{\mathcal{F}_{t+1}}}{1 + c_{t+1}^{2}} - \mathbb{E}_{t} \left[\frac{\Theta^{\top} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1} \boldsymbol{\mathcal{P}}_{t+1} \nabla_{\mathsf{St}} \boldsymbol{L}_{t+1}^{\top} \Theta \mathbb{1}_{\mathcal{F}_{t+1}}}{1 + c_{t+1}^{2}}\right] (F.38) \\ \mathsf{Let} \, \mathsf{X} &\in \left\{\frac{\eta/2}{\sqrt{r_{s}}} \mathsf{A}, \frac{\eta^{2}/16}{r_{s}} \mathsf{B}, \frac{\eta^{2}/16}{r_{s}} \mathsf{C}, \frac{\eta^{3}/32}{r_{s}^{3/2}} \mathsf{D}, \frac{\eta^{4}/256}{r_{s}^{2}} \mathsf{F}\right\} \text{ and} \\ &\Gamma_{1} \coloneqq \left\{\boldsymbol{I}_{r}, \quad \alpha \in [0, 0.5) \\ \boldsymbol{I}_{r_{u}} \quad \alpha > 0.5 \right\} &\Gamma_{2} \coloneqq \left\{\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}, \quad \alpha \in [0, 0.5) \\ \boldsymbol{\Lambda}_{11}^{\frac{1}{2}}, \quad \alpha > 0.5. \right\} \end{split}$$

For $\ell \in \{1, 2\}$, we define:

$$\operatorname{Quad}_{k,t}^{(\ell)}(\mathsf{X}) \coloneqq \sum_{j=1}^k \mathbb{E}_{j-1} \left[\left(\mathsf{\Gamma}_\ell \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{X}_j \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{\Gamma}_\ell \right)^2 \right].$$

We have the following corollary.

Corollary 4. Let $\mathsf{rk} \in \{r, r_u\}$, $\mathsf{rk}_{\star} \in \{r, r_s\}$ and

$$\mathsf{S}_t \coloneqq \mathsf{\eta} \sum_{j=1}^t \mathsf{G}_{j-1} \ \text{ and } \mathsf{\eta} = \frac{\eta}{\sqrt{r_s} \|\mathbf{\Lambda}\|_{\mathrm{F}}} \ \text{ and } \ r_u = \lceil \log^{2.5} d \rceil.$$

Assume the following conditions hold:

•
$$m{T}_t \succeq rac{\kappa_d r_s}{d} m{I}_{\mathsf{rk}}$$
,

•
$$T_t^{\frac{-1}{2}} \mathsf{S}_t T_t^{\frac{-1}{2}} \preceq \frac{C_{ub}}{\kappa_d} \begin{cases} (2\eta t \vee r^{\alpha}) \mathbf{I}_r, & \alpha \in [0, 0.5) \\ r_s^{\alpha} \log d \mathbf{I}_{r_u}, & \alpha > 0.5, \end{cases}$$

•
$$T_t^{\frac{-1}{2}}\mathsf{G}_jT_t^{\frac{-1}{2}} \preceq \frac{C_{ub}}{\mathsf{K}_d}I_{\mathsf{rk}}$$
 for $j \leq t-1$.

Let

$$\begin{cases} \mathsf{p}_1 = 1 \ \ \textit{and} \ \ \mathsf{p}_2 = 1 - \alpha, & \alpha \in [0, 1) \\ \mathsf{p}_1 = 1 \ \ \textit{and} \ \ \mathsf{p}_2 = \frac{2 \log \log r_u}{\log r_u} & \alpha = 1 \\ \mathsf{p}_1 = 1 \ \ \textit{and} \ \ \mathsf{p}_2 = 0 & \alpha > 1. \end{cases}$$

For $\eta t \leq \frac{1}{2} r k_{\star}^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$, the following results hold:

(a) Quadratic variation bounds. We have:

$$\|\mathit{Quad}_{t,t}^{(\ell)}(\mathsf{X})\|_2 \leq \frac{C\log d}{\kappa_d^2} \frac{\|\mathbf{\Lambda}\|_{\mathrm{F}} \mathsf{rk}^{\mathsf{p}_\ell}_{\star} \mathsf{rk}^{\alpha}_{\star}}{r_s^{3/2}} \begin{cases} C_{ub} \eta d & \mathsf{X} = \frac{\eta/2}{\sqrt{r_s}} \mathsf{A} \\ C_{ub}^2 \eta^3 d^2, & \mathsf{X} = \frac{\eta^2/16}{r_s} \mathsf{B} \\ \eta^3 d^2, & \mathsf{X} = \frac{\eta^2/16}{r_s} \mathsf{C} \\ C_{ub} \eta^5 d^3, & \mathsf{X} = \frac{\eta^3/32}{r_s^{3/2}} \mathsf{D} \\ \eta^7 d^2, & \mathsf{X} = \frac{\eta^4/256}{r_s^2} \mathsf{F}. \end{cases}$$

(b) Operator norm bounds. For $L \ge 8\sqrt{2}e$, there exists C > 0 such that

$$\begin{split} \left\| \Gamma_{\ell} T_{t}^{\frac{-1}{2}} \mathsf{X}_{j} T_{t}^{\frac{-1}{2}} \Gamma_{\ell} \right\|_{2} &\leq r_{j,t}^{(\ell)}(\mathsf{X}) \coloneqq \begin{cases} \frac{\eta L^{2}}{2\sqrt{r_{s}}} \sqrt{Tr \left(\Gamma_{\ell}^{2} T_{t}^{-1}\right) Tr \left(\Gamma_{\ell}^{2} T_{t}^{\frac{-1}{2}} \mathsf{G}_{j-1} T_{t}^{\frac{-1}{2}}\right)}, & \mathsf{X} = \frac{\eta/2}{\sqrt{r_{s}}} \mathsf{A} \\ \frac{\eta^{2} L^{4}}{16 r_{s}} dTr \left(\Gamma_{\ell}^{2} T_{t}^{\frac{-1}{2}} \mathsf{G}_{j-1} T_{t}^{\frac{-1}{2}}\right), & \mathsf{X} = \frac{\eta^{2}/16}{r_{s}} \mathsf{B} \\ \frac{\eta^{2} L^{4}}{16} Tr \left(\Gamma_{\ell}^{2} T_{t}^{-1}\right), & \mathsf{X} = \frac{\eta^{2}/16}{r_{s}} \mathsf{C} \\ \frac{\eta^{3} L^{6}}{32\sqrt{r_{s}}} d\sqrt{Tr \left(\Gamma_{\ell}^{2} T_{t}^{-1}\right) Tr \left(\Gamma_{\ell}^{2} T_{t}^{\frac{-1}{2}} \mathsf{G}_{j-1} T_{t}^{\frac{-1}{2}}\right)}, & \mathsf{X} = \frac{\eta^{3}/32}{r_{s}^{3/2}} \mathsf{D} \\ \frac{\eta^{4} L^{8}}{256} dTr \left(\Gamma_{\ell}^{2} T_{t}^{-1}\right), & \mathsf{X} = \frac{\eta^{4}/256}{r_{s}^{2}} \mathsf{F} \end{cases} \\ \leq \frac{C}{\kappa_{d}} \frac{\mathsf{rk}^{\mathsf{P}\ell}}{r_{s}} \begin{cases} L^{2} \sqrt{C_{ub}} \eta \sqrt{d}, & \mathsf{X} = \frac{\eta^{2}/16}{r_{s}} \mathsf{B} \\ L^{4} \eta^{2} d, & \mathsf{X} = \frac{\eta^{2}/16}{r_{s}} \mathsf{C} \\ L^{6} \sqrt{C_{ub}} \eta^{3} d^{3/2}, & \mathsf{X} = \frac{\eta^{3}/32}{r_{s}^{3/2}} \mathsf{D} \\ L^{8} \eta^{4} d^{2}, & \mathsf{X} = \frac{\eta^{4}/256}{r_{s}^{2}} \mathsf{F}. \end{cases} \end{cases}$$

Proof. Quadratic variation bounds. We will use the variance bounds given in Proposition 18. For $X = \frac{\eta/2}{\sqrt{r_-}}A$, we have

$$\begin{aligned} \operatorname{Quad}_{t,t}^{(\ell)}(\frac{\eta/2}{\sqrt{r_s}}\mathsf{A}) & \preceq \frac{C\eta \|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} \left(\operatorname{Tr} \left(\Gamma_{\ell}^2 \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{S}_t \boldsymbol{T}_t^{\frac{-1}{2}} \right) \Gamma_{\ell} \boldsymbol{T}_t^{-1} \Gamma_{\ell} + \operatorname{Tr} \left(\Gamma_{\ell}^2 \boldsymbol{T}_t^{-1} \right) \Gamma_{\ell} \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{S}_t \boldsymbol{T}_t^{\frac{-1}{2}} \Gamma_{\ell} \right) \\ & \preceq \frac{C\eta \|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} \frac{C_{\mathsf{ub}} \mathsf{rk}^{\mathsf{p}\ell} \mathsf{rk}_{\star}^{\alpha} \log d}{\kappa_d^2} \frac{d}{r_s} \boldsymbol{I}_{\mathsf{rk}}. \end{aligned}$$

For $X = \frac{\eta^2/16}{r_s}B$, we have

$$\begin{split} \operatorname{Quad}_{t,t}^{(\ell)}(\tfrac{\eta^2/16}{r_s}\mathsf{B}) & \preceq \frac{CC_{\mathsf{ub}}\eta^3d^2\|\mathbf{\Lambda}\|_{\mathsf{F}}}{r_s^{3/2}} \sup_{j \leq t} \left(\operatorname{Tr} \big(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{G}_{j-1} \boldsymbol{T}_t^{\frac{-1}{2}} \big) \right) \mathsf{\Gamma}_\ell \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{S}_t \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{\Gamma}_\ell \\ & \preceq \frac{C\eta^3d^2\|\mathbf{\Lambda}\|_{\mathsf{F}}}{r_s^{3/2}} \frac{C_{\mathsf{ub}}^2 \mathsf{rk}^{\mathsf{p}\ell} \mathsf{rk}_\star^\alpha \log d}{\mathsf{\kappa}_d^2} \boldsymbol{I}_{\mathsf{rk}}. \end{split}$$

For $X = \frac{\eta^2/16}{r}$ C, we have

$$\operatorname{Quad}_{t,t}^{(\ell)}(\tfrac{\eta^2/16}{r_s}\mathsf{C}) \preceq C\eta^4 t \operatorname{Tr}(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{-1}) \mathsf{\Gamma}_\ell \boldsymbol{T}_t^{-1} \mathsf{\Gamma}_\ell \preceq \frac{C\eta^3 d^2 \|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{r_s^{3/2}} \frac{\mathsf{rk}^{\mathsf{p}_\ell} \mathsf{rk}_\star^\alpha \log d}{\mathsf{k}_d^2} \boldsymbol{I}_{\mathsf{rk}}.$$

For $X = \frac{\eta^3/32}{r_s^{3/2}}D$, we have

$$\begin{split} \operatorname{Quad}_{t,t}^{(\ell)}(\tfrac{\eta^3/32}{r_s^{3/2}}\mathsf{D}) \! & \preceq \! \frac{C\eta^5d^2\|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} \Big(\mathrm{Tr}(\mathsf{\Gamma}_{\ell}^2 \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{S}_t \boldsymbol{T}_t^{\frac{-1}{2}}) \mathsf{\Gamma}_{\ell} \boldsymbol{T}_t^{-1} \mathsf{\Gamma}_{\ell} + \mathrm{Tr}\big(\mathsf{\Gamma}_{\ell}^2 \boldsymbol{T}_t^{-1}\big) \mathsf{\Gamma}_{\ell} \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{S}_t \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{S}_t \boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{\Gamma}_{\ell} \Big) \\ & \preceq \frac{C\eta^5d^2\|\mathbf{\Lambda}\|_{\mathrm{F}}}{\sqrt{r_s}} \frac{C_{\mathsf{ub}}\mathsf{rk}^{\mathsf{p\ell}}\mathsf{rk}_{\star}^{\mathsf{q}} \log d}{\mathsf{k}_d^2} \frac{d}{r_s} \boldsymbol{I}_{\mathsf{rk}}. \end{split}$$

For $X = \frac{\eta^4/256}{r_s^2} F$, we have

$$\mathrm{Quad}_{t,t}^{(\ell)}(\tfrac{\eta^4/256}{r_s^2}\mathsf{F}) \preceq C\eta^8 d^2t \mathrm{Tr}(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{-1}) \mathsf{\Gamma}_\ell \boldsymbol{T}_t^{-1} \mathsf{\Gamma}_\ell \preceq \frac{C\eta^7 d^2 \|\boldsymbol{\Lambda}\|_{\mathrm{F}}}{r_s^{3/2}} \frac{\mathsf{rk}^{\mathsf{p}_\ell} \mathsf{rk}_\star^\alpha \log d}{\mathsf{\kappa}_d^2} \boldsymbol{I}_{\mathsf{rk}}.$$

Operator Norm Bounds. We will use the events defined in Proposition 18. For $X = \frac{\eta/2}{\sqrt{r_e}}A$, we have

$$r_{j,t}^{(\ell)}(\tfrac{\eta/2}{\sqrt{r_s}}\mathsf{A}) = \frac{\eta/2}{\sqrt{r_s}}L^2\sqrt{\mathrm{Tr}(\mathsf{\Gamma}_\ell^2\boldsymbol{T}_t^{-1})\mathrm{Tr}(\mathsf{\Gamma}_\ell^2\boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{G}_{j-1}\boldsymbol{T}_t^{\frac{-1}{2}})} \leq \frac{CL^2}{\mathsf{K}_d}\frac{\sqrt{C_{\mathrm{ub}}}\eta\sqrt{d}\mathsf{r}\mathsf{k}^{\mathsf{P}_\ell}}{r_s}.$$

For $X = \frac{\eta^2/16}{r_s}B$, we have

$$r_{j,t}^{(\ell)}(\tfrac{\eta^2/16}{r_s}\mathsf{B}) = \frac{\eta^2}{16r_s}L^4d\operatorname{Tr}(\mathsf{\Gamma}_\ell^2\boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{G}_{j-1}\boldsymbol{T}_t^{\frac{-1}{2}}) \leq \frac{CL^4}{\kappa_d}\frac{C_{\mathsf{ub}}\eta^2d\mathsf{rk}^{\mathsf{P}\ell}}{r_s}.$$

For $X = \frac{\eta^2/16}{r_s}C$, we have

$$r_{j,t}^{(\ell)}(\tfrac{\eta^2/16}{r_s}\mathsf{C}) = \frac{\eta^2}{16r_s}L^4\mathrm{Tr}(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{-1}) \leq \frac{CL^4}{\kappa_d}\frac{\eta^2 d\mathsf{rk}^{\mathsf{p}_\ell}}{r_s}.$$

For $X = \frac{\eta^3/32}{r_s^{3/2}}D$, we have

$$r_{j,t}^{(\ell)}(\tfrac{\eta^3/32}{r_s^{3/2}}\mathsf{D}) = \frac{\eta^3}{32\sqrt{r_s}}L^6d\sqrt{\mathrm{Tr}(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{G}_{j-1}\boldsymbol{T}_t^{\frac{-1}{2}})\mathrm{Tr}(\mathsf{\Gamma}_\ell^2 \boldsymbol{T}_t^{-1})} \leq \frac{CL^6}{\kappa_d}\frac{\sqrt{C_{\mathrm{ub}}}\eta^3d^{3/2}\mathsf{rk}^{\mathsf{p}_\ell}}{r_s}.$$

For $X = \frac{\eta^4/256}{r_s^2} F$, we have

$$r_{j,t}^{(\ell)}(\tfrac{\eta^4/256}{r_s^2}\mathsf{F}) = \frac{\eta^4}{256} L^8 d\mathrm{Tr}(\mathsf{\Gamma}_\ell^2 \pmb{T}_t^{-1}) \leq \frac{CL^8}{\mathsf{\kappa}_d} \frac{\eta^4 d^2 \mathsf{rk}^{\mathsf{p}_\ell}}{r_s}.$$

Proposition 19. Let $\{Y_t, t = 1, 2 \cdots \}$ be a symmetric-matrix martingale with difference sequence $\{X_t := Y_{t+1} - Y_t, t = 1, 2 \cdots \}$, whose values are symmetric matrices with dimension $r \leq d$. Let $\{T_t, t = 1, 2, \cdots \}$ be a deterministic sequence, whose values are positive semi-definite matrices with the same dimensionality. Assume that the difference sequence is uniformly bounded in the sense that for a predictable triangular sequence $\{r_{j,t}\}_{j \leq t}$, we have

$$\lambda_{max}(T_t^{-\frac{1}{2}} X_j T_t^{-\frac{1}{2}}) \le r_{j,t} \text{ for } j = 1, 2, \dots, t.$$

Define the predictable uniform bound and quadratic variation process of the martingale:

$$R_{k,t} := \max_{j \le k} r_{j,t} \text{ and } Quad_{k,t}(\mathsf{X}) := \sum_{j=1}^k \mathbb{E}_{j-1} \left[\left(T_t^{\frac{-1}{2}} \mathsf{X}_j T_t^{\frac{-1}{2}} \right)^2 \right] \text{ for } k \le t = 1, 2, \cdots.$$

Let $T \leq d^3$ be a bounded stopping time. Then, for any deterministic σ^2 , $\widetilde{L} > 0$

$$\mathbb{P}\left[\exists t \leq \mathcal{T}, \mathsf{Y}_t \not\preceq u \boldsymbol{T}_t \ and \ \max_{t \leq \mathcal{T}} \|Quad_{t,t}(\mathsf{X})\|_2 \leq \sigma^2 \ and \ \max_{t \leq \mathcal{T}} R_{t,t} \leq \widetilde{L}\right]$$

$$\leq d^4 \cdot \exp\left(\frac{-u^2/2}{\sigma^2 + \widetilde{L}u/3}\right).$$

Proof. We have

$$\begin{split} \mathcal{E}_{\text{target}} &\equiv \left\{\exists t \leq \mathcal{T}, \mathsf{Y}_t \not\preceq u \boldsymbol{T}_t \text{ and } \max_{t \leq \mathcal{T}} \|\mathsf{Quad}_{t,t}(\mathsf{X})\|_2 \leq \sigma^2 \text{ and } \max_{t \leq \mathcal{T}} R_{t,t} \leq \widetilde{L}\right\} \\ &\subseteq \bigcup_{t=0}^{\mathcal{T}} \left\{\exists k \leq t, \mathsf{Y}_k \not\preceq u \boldsymbol{T}_t \text{ and } \|\mathsf{Quad}_{t,t}(\mathsf{X})\|_2 \leq \sigma^2 \text{ and } R_{t,t} \leq \widetilde{L}\right\}. \end{split}$$

Therefore, we have

$$\mathbb{P}\big[\mathcal{E}_{\text{target}}\big] \leq \sum_{n=1}^{d^3} \mathbb{P}\left[\exists k \leq t, \ \mathsf{Y}_k \not\preceq u \boldsymbol{T}_t \ \text{ and } \ \|\mathsf{Quad}_{t,t}(\mathsf{X})\|_2 \leq \sigma^2 \text{ and } R_{t,t} \leq \widetilde{L}\right]. \tag{F.39}$$

In the following, we will bound the each term in the right hands-side of (F.39). By [Tro10, Lemma 6.7], we have for $k = 1, \dots, t$ and $\theta > 0$,

$$\mathbb{1}_{R_{k,t}\leq \widetilde{L}}\mathbb{E}_{k-1}\Big[e^{\frac{\theta}{L}\boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{X}_k\boldsymbol{T}_t^{\frac{-1}{2}}}\Big] \leq \mathbb{1}_{R_{k,t}\leq \widetilde{L}}\exp\left(\frac{e^{\theta}-\theta-1}{\widetilde{L}^2}\mathbb{E}_{k-1}\left[(\boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{X}_k\boldsymbol{T}_t^{\frac{-1}{2}})^2\right]\right). (F.40)$$

For notational convenience call $g(\theta) := e^{\theta} - \theta - 1$. We define a super martingale such that for 0 < k < t,

$$S_k \coloneqq \operatorname{Tr}\left(\exp\left(\frac{\theta}{L}\boldsymbol{T}_t^{\frac{-1}{2}}\mathsf{Y}_k\boldsymbol{T}_t^{\frac{-1}{2}} - \frac{g(\theta)}{\widetilde{L}^2}\mathrm{Quad}_{k,t}(\mathsf{X})\right)\right)\mathbbm{1}_{R_{k,t}\leq \widetilde{L}},$$

with initial values $R_{0,t}=0$, $Y_0=\operatorname{Quad}_{0,t}=0$, and thus, $S_0=r$. Note that by (F.40) and [Tro10, Corollary 3.3], we can show that $\mathbb{E}_{k-1}S_k \leq S_{k-1}$. We define a stopping time and an event

$$\mathcal{T}_{hit} := \{k \ge 0 \mid \lambda_{max}(\boldsymbol{T}_t^{\frac{-1}{2}} \mathsf{Y}_k \boldsymbol{T}_t^{\frac{-1}{2}}) \ge u\} \land t,$$

$$\mathcal{E}_{hit} := \{\mathcal{T}_{hit} \le t\} \cap \{\|\mathsf{Quad}_{t,t}(\mathsf{X})\|_2 \le \sigma^2 \text{ and } R_{t,t} \le \widetilde{L}\}.$$

We have

$$\begin{split} \mathbb{1}_{\mathcal{E}_{\text{hit}}} S_{\mathcal{T}_{\text{hit}}} & \geq \mathbb{1}_{\mathcal{E}_{\text{hit}}} \exp\left(\frac{\theta}{\widetilde{L}} u - \frac{g(\theta)}{\widetilde{L}^2} \sigma^2\right) \overset{(a)}{\Rightarrow} r \geq \mathbb{P}\left[\mathcal{E}_{\text{hit}}\right] \exp\left(\frac{\theta}{\widetilde{L}} u - \frac{g(\theta)}{\widetilde{L}^2} \sigma^2\right) \\ & \Rightarrow r \inf_{\theta > 0} \exp\left(-\theta \frac{u}{\widetilde{L}} + g(\theta) \frac{\sigma^2}{\widetilde{L}^2}\right) \geq \mathbb{P}\left[\mathcal{E}_{\text{hit}}\right], \end{split}$$

where we use Doob's optional sampling theorem in (a). Since the infimum is attained at $\theta > 0$ and the convex conjugate of $g(\theta)$ is $g^*(\theta) = (\theta + 1)\log(\theta + 1) - \theta$, we have

$$\mathbb{P}\left[\mathcal{E}_{\mathsf{hit}}\right] \leq r \cdot \exp\left(-\frac{\sigma^2}{\widetilde{L}^2} g^{\star}\left(\frac{u\widetilde{L}}{\sigma^2}\right)\right) \leq r \cdot \exp\left(\frac{-u^2/2}{\sigma^2 + \widetilde{L}u/3}\right),$$

where we used $g^*(\theta) \ge \frac{u^2/2}{1+u/3}$ in the last step. By $r \le d$ and (F.39), we have the statement. \square

Proposition 20. Let \mathbb{P}_0 denote the conditional probability conditioned on W_0 . We consider $r_u = \lceil \log^{2.5} d \rceil$, and

$$\alpha \in [0, 0.5): \frac{r_s}{r} \to \varphi, \quad \eta \asymp \frac{1}{dr^{\alpha} \log^{20}(1 + d/r_s)}, \quad \kappa_d = \frac{1}{\log^{3.5} d}, \qquad C_{ub} = 12 \left(1 + \frac{1}{\sqrt{\varphi}}\right)^2$$

$$\alpha > 0.5: \quad r_s \approx 1, \quad \eta \approx \frac{1}{d \, r_u^{4\alpha + 3} \log^{18} d}, \quad \kappa_d = \frac{1}{r_u \log^{2.5} d}, \quad C_{ub} = 2^{\alpha} 30 r_u.$$

For $\alpha \in [0, 0.5)$, we define $\mathcal{T} := \mathcal{T}_{bad} \wedge \frac{1}{2\eta} \left(r_s \left(1 - \log^{\frac{-1}{2}} d \right) \wedge r \right)^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$. We have for $d \geq \Omega_{\alpha, \varphi, \beta}(1)$

$$\begin{split} \mathbb{P}_{0} \left[\sup_{t \leq \mathcal{T}} & \| \boldsymbol{T}_{t}^{\frac{-1}{2}} \underline{\boldsymbol{\nu}}_{t} \boldsymbol{T}_{t}^{\frac{-1}{2}} \|_{2} \vee r^{\frac{\alpha}{2}} \| \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{T}_{t}^{\frac{-1}{2}} \underline{\boldsymbol{\nu}}_{t} \boldsymbol{T}_{t}^{\frac{-1}{2}} \boldsymbol{\Lambda}^{\frac{1}{2}} \|_{2} \geq \kappa_{d} r^{\frac{-\alpha}{2}} \quad and \quad \mathcal{G}_{init} \right] \\ & \leq 20 d^{4} \exp(-\log^{2} d). \end{split}$$

For
$$\alpha > 0.5$$
, we set $\mathcal{T} := \mathcal{T}_{bad} \wedge \frac{1}{2\eta} r_s^{\alpha} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$. We have for $d \geq \Omega_{\alpha, r_s}(1)$

$$\mathbb{P}_{0}\left[\sup_{t \leq \mathcal{T}} \|\boldsymbol{T}_{t}^{\frac{-1}{2}} \underline{\boldsymbol{\nu}}_{t} \boldsymbol{T}_{t}^{\frac{-1}{2}} \|_{2} \vee r_{u}^{\frac{\alpha}{2}} \|\boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{T}_{t}^{\frac{-1}{2}} \underline{\boldsymbol{\nu}}_{t} \boldsymbol{T}_{t}^{\frac{-1}{2}} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \|_{2} \geq \kappa_{d} r_{u}^{\frac{-\alpha}{2}} \text{ and } \mathcal{G}_{\textit{init}}\right] \\ \leq 20d^{4} \exp(-\log^{2} d).$$

 $\begin{aligned} \textit{Proof.} \ \ \text{For notational convenience, we introduce} \ \mathcal{X} &:= \big\{ \frac{\eta/2}{\sqrt{r_s}} \mathsf{A}, \, \frac{\eta^2/16}{r_s} \mathsf{B}, \, \frac{\eta^2/16}{r_s} \mathsf{C}, \, \frac{\eta^3/32}{r_s^{3/2}} \mathsf{D}, \, \frac{\eta^4/256}{r_s^2} \mathsf{F} \big\}. \end{aligned}$ For both cases, we will set the clip threshold to $L = \log^2 d$. We introduce the notation $R_t^{(\ell)} := \max_{\mathsf{X} \in \mathcal{X}} \max_{j \leq t} r_{j,t}^{(\ell)}(\mathsf{X})$ and $\|\mathsf{Quad}_t^{(\ell)}\|_2 := \max_{\mathsf{X} \in \mathcal{X}} \|\mathsf{Quad}_{t,t}^{(\ell)}(\mathsf{X})\|_2$ for $\ell = 1, 2$.

For $\alpha \in [0, 0.5)$, we can write for all $X \in \mathcal{X}$,

$$\begin{split} &\mathbb{P}_0\left[\sup_{t\leq\mathcal{T}} & \|\boldsymbol{T}_t^{\frac{-1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{T}_t^{\frac{-1}{2}}\|_2 \vee r^{\frac{\alpha}{2}}\|\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{T}_t^{\frac{-1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{T}_t^{\frac{-1}{2}}\boldsymbol{\Lambda}^{\frac{1}{2}}\|_2 \geq \kappa_d r^{\frac{-\alpha}{2}} \ \text{and} \ \mathcal{G}_{\text{init}}\right] \\ &\leq \sum_{\mathsf{X}\in\mathcal{X}} \mathbb{P}_0\bigg[\sup_{t\leq\mathcal{T}} \left\|\boldsymbol{T}_t^{\frac{-1}{2}}\big(\sum_{j\leq t}\mathsf{X}_j\big)\boldsymbol{T}_t^{\frac{-1}{2}}\right\|_2 \geq \frac{\kappa_d r^{\frac{-\alpha}{2}}}{10} \ \text{and} \ \mathcal{G}_{\text{init}}\bigg] \\ &+ \sum_{\mathsf{X}\in\mathcal{X}} \mathbb{P}_0\bigg[\sup_{t\leq\mathcal{T}} \left\|\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{T}_t^{\frac{-1}{2}}\big(\sum_{j\leq t}\mathsf{X}_j\big)\boldsymbol{T}_t^{\frac{-1}{2}}\boldsymbol{\Lambda}^{\frac{1}{2}}\right\|_2 \geq \frac{\kappa_d r^{-\alpha}}{10} \ \text{and} \ \mathcal{G}_{\text{init}}\bigg]. \end{split}$$

By Propositions 14 and 16 and Corollary 4, \mathcal{G}_{init} implies the events

$$\begin{split} \mathcal{E}_{\mathrm{ht},1} &\equiv \bigg\{ \max_{t \leq \mathcal{T}} \lVert \mathsf{Quad}_t^{(1)} \rVert_2 \leq \frac{O_{\alpha,\beta,\varphi}\left(r^{-\alpha}\right)}{\log^{12}d} \quad \text{and} \quad \max_{t \leq \mathcal{T}} R_t^{(1)} \leq \frac{O_{\alpha,\beta,\varphi}\left(r^{-\alpha}\right)}{\sqrt{d}\log^{12.5}d} \bigg\} \\ \mathcal{E}_{\mathrm{ht},2} &\equiv \bigg\{ \max_{t \leq \mathcal{T}} \lVert \mathsf{Quad}_t^{(2)} \rVert_2 \leq \frac{O_{\alpha,\beta,\varphi}\left(r^{-2\alpha}\right)}{\log^{12}d} \quad \text{and} \quad \max_{t \leq \mathcal{T}} R_t^{(2)} \leq \frac{O_{\alpha,\beta,\varphi}\left(r^{-2\alpha}\right)}{\sqrt{d}\log^{12.5}d} \bigg\}. \end{split}$$

Therefore, by using Proposition 19

$$\begin{split} \mathbb{P}_0 \Big[\sup_{t \leq \mathcal{T}} \Big\| \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \Big\|_2 &\geq \frac{\mathsf{\kappa}_d r^{\frac{-\alpha}{2}}}{10} \ \text{and} \ \mathcal{G}_{\mathsf{init}} \Big] \\ &\leq \mathbb{P}_0 \Big[\sup_{t \leq \mathcal{T}} \Big\| \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \Big\|_2 &\geq \frac{\mathsf{\kappa}_d r^{\frac{-\alpha}{2}}}{10} \ \text{and} \ \mathcal{E}_{\mathsf{ht},1} \Big] \\ &\leq 2d^4 \exp \left(-\log^2 d \right). \end{split}$$

Similarly,

$$\begin{split} & \mathbb{P}_0 \left[\sup_{t \leq \mathcal{T}} \left\| \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{T}_t^{\frac{-1}{2}} \left(\sum_{j \leq t} \mathsf{X}_j \right) \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{\Lambda}^{\frac{1}{2}} \right\|_2 \geq \frac{\mathsf{\kappa}_d r^{-\alpha}}{10} \quad \text{and} \quad \mathcal{G}_{\text{init}} \right] \\ & \leq \mathbb{P}_0 \left[\sup_{t \leq \mathcal{T}} \left\| \boldsymbol{\Lambda}^{\frac{1}{2}} \boldsymbol{T}_t^{\frac{-1}{2}} \left(\sum_{j \leq t} \mathsf{X}_j \right) \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{\Lambda}^{\frac{1}{2}} \right\|_2 \geq \frac{\mathsf{\kappa}_d r^{-\alpha}}{10} \quad \text{and} \quad \mathcal{E}_{\text{ht},2} \right] \\ & \leq 2d^4 \exp\left(-\log^2 d \right). \end{split}$$

For $\alpha > 0.5$, we can write for all $X \in \mathcal{X}$,

$$\begin{split} &\mathbb{P}_0\left[\sup_{t\leq\mathcal{T}} & \|\boldsymbol{T}_t^{\frac{-1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{T}_t^{\frac{-1}{2}}\|_2 \vee r_u^{\frac{\alpha}{2}} \|\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\boldsymbol{T}_t^{\frac{-1}{2}}\underline{\boldsymbol{\nu}}_t\boldsymbol{T}_t^{\frac{-1}{2}}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\|_2 \geq \kappa_d r_u^{\frac{-\alpha}{2}} \text{ and } \mathcal{G}_{\text{init}}\right] \\ &\leq \sum_{\mathsf{X}\in\mathcal{X}} \mathbb{P}_0\Big[\sup_{t\leq\mathcal{T}} \left\|\boldsymbol{T}_t^{\frac{-1}{2}}\big(\sum_{j\leq t}\mathsf{X}_j\big)\boldsymbol{T}_t^{\frac{-1}{2}}\right\|_2 \geq \frac{\kappa_d r_u^{\frac{-\alpha}{2}}}{10} \text{ and } \mathcal{G}_{\text{init}}\Big] \\ &+ \sum_{\mathsf{X}\in\mathcal{X}} \mathbb{P}_0\Big[\sup_{t\leq\mathcal{T}} \left\|\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\boldsymbol{T}_t^{\frac{-1}{2}}\big(\sum_{j\leq t}\mathsf{X}_j\big)\boldsymbol{T}_t^{\frac{-1}{2}}\boldsymbol{\Lambda}_{11}^{\frac{1}{2}}\right\|_2 \geq \frac{\kappa_d r_u^{-\alpha}}{10} \text{ and } \mathcal{G}_{\text{init}}\Big] \end{split}$$

By Propositions 14 and 16 and Corollary 4, \mathcal{G}_{init} implies the events

$$\begin{split} \mathcal{E}_{\mathrm{lt},1} &\equiv \bigg\{ \max_{t \leq \mathcal{T}} \lVert \mathrm{Quad}_{t}^{(1)} \rVert_{2} \leq \frac{O_{\alpha,r_{s}}(r_{u}^{-4\alpha})}{\log^{12}d} \quad \text{ and } \max_{t \leq \mathcal{T}} R_{t}^{(1)} \leq \frac{O_{r_{s}}(r_{u}^{-4\alpha-1})}{\sqrt{d}\log^{11.5}d} \bigg\} \\ \mathcal{E}_{\mathrm{lt},2} &\equiv \bigg\{ \max_{t \leq \mathcal{T}} \lVert \mathrm{Quad}_{t}^{(2)} \rVert_{2} \leq \frac{O_{\alpha,r_{s}}(r_{u}^{-4\alpha-(\alpha \wedge 1)})}{\log^{12}d\log^{-2}r_{u}} \quad \text{and } \max_{t \leq \mathcal{T}} R_{t}^{(2)} \leq \frac{O_{r_{s}}(r_{u}^{-4\alpha-1}r_{u}^{-(\alpha \wedge 1)})}{\sqrt{d}\log^{11.5}d\log^{-2}r_{u}} \bigg\}. \end{split}$$

Therefore, by using Proposition 19

$$\begin{split} \mathbb{P}_0 \Big[\sup_{t \leq \mathcal{T}} \Big\| \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \Big\|_2 &\geq \frac{\mathsf{\kappa}_d r_u^{\frac{-\alpha}{2}}}{10} \ \text{and} \ \mathcal{G}_{\mathsf{init}} \Big] \\ &\leq \mathbb{P}_0 \Big[\sup_{n \leq \mathcal{T}} \Big\| \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \Big\|_2 &\geq \frac{\mathsf{\kappa}_d r_u^{\frac{-\alpha}{2}}}{10} \ \text{and} \ \mathcal{E}_{\mathsf{lt},1} \Big] \\ &\leq 2d^4 \exp \big(-\log^2 d \big) \,. \end{split}$$

Similarly,

$$\begin{split} & \mathbb{P}_0 \bigg[\sup_{t \leq \mathcal{T}} \left\| \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \right\|_2 \geq \frac{\mathsf{\kappa}_d r_u^{-\alpha}}{10} \quad \text{and} \quad \mathcal{G}_{\text{init}} \bigg] \\ & \leq \mathbb{P}_0 \bigg[\sup_{t \leq \mathcal{T}} \left\| \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \boldsymbol{T}_t^{\frac{-1}{2}} \big(\sum_{j \leq t} \mathsf{X}_j \big) \boldsymbol{T}_t^{\frac{-1}{2}} \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \right\|_2 \geq \frac{\mathsf{\kappa}_d r_u^{-\alpha}}{10} \quad \text{and} \quad \mathcal{E}_{\text{lt},2} \bigg] \\ & \leq 2d^4 \exp\left(-\log^2 d\right). \end{split}$$

Corollary 5. Consider $\operatorname{rk}_{\star} = \{r_{\star} = \lfloor r_s (1 - \log^{-1/2} d) \wedge r \rfloor, r_{u_{\star}} = r_s \}$ and the parameters in Proposition 20. We have

$$\mathbb{P}_0\left[\mathcal{T}_{bad} \geq \frac{1}{2\eta} \mathsf{rk}_\star \log\left(\frac{d\log^{1.5} d}{r_s}\right) \text{ and } \mathcal{G}_{init}\right] \geq 1 - 20d^4 \exp(-\log^2 d).$$

Proof. By using the first items in Proposition 15 and 16 ,and Lemma 6, \mathcal{G}_{init} implies that

$$\mathcal{T}_{\mathsf{bad}} \geq \mathcal{T}_{\mathsf{noise}} \wedge \frac{1}{2\eta} \mathsf{rk}_{\star} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$$

On the other hand, within the (negation) of the events given in Proposition 20, we have

$$\mathcal{T}_{\text{noise}} > \mathcal{T}_{\text{bad}} \wedge \frac{1}{2\eta} \mathsf{rk}_{\star} \log \left(\frac{d \log^{1.5} d}{r_s} \right)$$

Therefore, the statement follows.

F.7 Stability near minima

In this section, we will establish that given (SGD) is near global minimum it will stay near global minimum. For the statement, we (re)introduce the block matrix notation: $\mathsf{rk}_\star = \{r_\star = \lfloor r_s (1 - \log^{-1/8} d) \land r \rfloor, r_{u_\star} = r_s\}$, we have

$$oldsymbol{G}_t = egin{bmatrix} oldsymbol{G}_{t,11} & oldsymbol{G}_{t,12} \ oldsymbol{G}_{t,12}^{ op} & oldsymbol{G}_{t,22} \end{bmatrix} \hspace{0.2cm} oldsymbol{
u}_t = egin{bmatrix} oldsymbol{
u}_{t,11} & oldsymbol{
u}_{t,12} \end{bmatrix} \hspace{0.2cm} oldsymbol{\Lambda} = egin{bmatrix} oldsymbol{\Lambda}_{11} & 0 \ 0 & oldsymbol{\Lambda}_{22} \end{bmatrix}, \hspace{0.2cm} oldsymbol{\Lambda}_{\ell_j} = egin{bmatrix} oldsymbol{\Lambda}_{\ell_j,11} & 0 \ 0 & oldsymbol{\Lambda}_{\ell_j,22} \end{bmatrix},$$

where $G_{t,11}, \nu_{t,11}, \Lambda_{11}, \Lambda_{\ell_j,11} \in \mathbb{R}^{\mathsf{rk}_{\star} \times \mathsf{rk}_{\star}}$ and Λ_{ℓ_j} is introduced (F.2). We define the following iterations:

• Given $\underline{\boldsymbol{G}}_0 = \boldsymbol{I}_{\mathsf{rk}_\star} - \frac{1}{\log d}$ diagonal and $\underline{\boldsymbol{V}}_t = 2\boldsymbol{\Lambda}_{\ell_2,11}^{\frac{1}{2}}\underline{\boldsymbol{G}}_t\boldsymbol{\Lambda}_{\ell_2,11}^{\frac{1}{2}} - \boldsymbol{\Lambda}_{\ell_1,11}$, we define

$$\begin{split} \underline{\boldsymbol{V}}_{t+1} &= \underline{\boldsymbol{V}}_t \left(\boldsymbol{I}_{\mathsf{rk}_\star} + \frac{\eta}{1-1.1\eta} \underline{\boldsymbol{V}}_t \right)^{-1} \\ &+ \frac{\eta}{1-1.1\eta} \left(\boldsymbol{\Lambda}_{\ell_1,11}^2 - \frac{8.1}{\mathsf{rk}_\star^2 \log d} \boldsymbol{\Lambda}_{\ell_2,11} - \frac{O(1)}{\log^2 d} \boldsymbol{\Lambda}_{\ell_2,11}^2 \right). \end{split}$$

• For $\underline{\boldsymbol{\nu}}_0 = 0$, $\underline{\boldsymbol{\nu}}_{t+1} = \underline{\boldsymbol{\nu}}_t + \frac{\eta/2}{\sqrt{r_s}} \boldsymbol{\nu}_{t+1,11}$.

• We define a sequence of events $\{\mathcal{E}_t\}_{t>0}$

$$\mathcal{E}_t \coloneqq \left\{ \frac{-\mathsf{rk}_{\star}^{-\frac{\alpha}{2}}}{\log^2 d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \preceq \underline{\boldsymbol{\nu}}_t \preceq \frac{\mathsf{rk}_{\star}^{-\frac{\alpha}{2}}}{\log^2 d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \right\} \cap \left\{ \frac{-\mathsf{rk}_{\star}^{-\alpha}}{\log^4 d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \preceq \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \underline{\boldsymbol{\nu}}_t \boldsymbol{\Lambda}_{11}^{\frac{1}{2}} \preceq \frac{\mathsf{rk}_{\star}^{-\alpha}}{\log^4 d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \right\},$$

We define the stopping times

$$\mathcal{T}_{\text{noise}}(\omega) \coloneqq \inf \left\{ t \geq 0 \mid \omega \not\in \mathcal{E}_t \right\} \wedge d^3, \quad \mathcal{T}_{\text{bounded}} \coloneqq \inf \left\{ t \geq 0 \mid \boldsymbol{G}_t \not\succeq \boldsymbol{I}_{\mathsf{rk}_{\star}} - \frac{2}{\log d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \right\}.$$

and $\mathcal{T}_{stable} := \mathcal{T}_{noise} \wedge \mathcal{T}_{bounded}$.

We have the following statement:

Proposition 21. Consider the parameters in Proposition 20. (SGD) guarantees that if $G_{0,11} \succeq I_{\mathsf{rk}_{\star}} - \frac{1}{\log d}$, we have $G_{t,11} \succeq I_{\mathsf{rk}_{\star}} - \frac{2}{\log d}$ for $t \leq \frac{\mathsf{rk}_{\star}^{\alpha} \log^2 d}{\eta}$ with probability $1 - d^4 \exp(-\log^2 d)$.

Proof. We define $\underline{\zeta}_t \coloneqq 2\Lambda_{\ell_2,11}^{\frac{1}{2}}\underline{\nu}_t\Lambda_{\ell_2,11}^{\frac{1}{2}}$. We make the following observations:

• Since $m{G}_{t,11}^2 + m{G}_{t,12} m{G}_{t,12}^ op m{I}_{\mathsf{rk}_\star}$ for $t \leq \mathcal{T}_{\mathsf{bounded}}$, we have

$$oldsymbol{G}_{t,12} oldsymbol{G}_{t,12}^{ op} \preceq rac{4}{\log d} oldsymbol{I}_{\mathsf{rk}_{\star}}$$

Therefore, by using (F.8), we have for $t \leq T_{\text{bounded}}$

$$egin{aligned} oldsymbol{G}_{t+1,11} \succeq oldsymbol{G}_{t,11} + \eta \Big(oldsymbol{\Lambda}_{\ell_1,11} oldsymbol{G}_{t,11} + oldsymbol{G}_{t,11} oldsymbol{\Lambda}_{\ell_1,11} - 2oldsymbol{G}_{t,11} oldsymbol{\Lambda}_{\ell_2,11} oldsymbol{G}_{t,11} \Big) - rac{4\eta}{\mathsf{rk}_{\star}^{lpha} \log d} oldsymbol{I}_{\mathsf{rk},\star} \\ & - C\eta^2 \|oldsymbol{\Lambda}\|_{\mathrm{F}}^2 r_s oldsymbol{I}_{\mathsf{rk},\star} + rac{\eta/2}{\sqrt{r_s}} oldsymbol{
u}_{t+1,11} \end{aligned}$$

Then, if we define $V_t^- \coloneqq 2\Lambda_{\ell_2,11}^{\frac{1}{2}}G_{t,11}\Lambda_{\ell_2,11}^{\frac{1}{2}} - \Lambda_{\ell_1,11}$, we have

$$\begin{split} \boldsymbol{V}_{t+1}^{-} \succeq \boldsymbol{V}_{t}^{-} - \eta (\boldsymbol{V}_{t+1}^{-})^{2} + \eta \boldsymbol{\Lambda}_{\ell_{1},11}^{2} - \left(\frac{8\eta}{\mathsf{rk}_{\star}^{\alpha} \log d} + C\eta^{2} \|\boldsymbol{\Lambda}\|_{\mathrm{F}}^{2} r_{s} \right) \boldsymbol{\Lambda}_{\ell_{2},11} \\ + \frac{\eta}{\sqrt{r_{s}}} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \\ & \succeq \boldsymbol{V}_{t}^{-} \left(\boldsymbol{I}_{r} + \frac{\eta}{1-1.1\eta} \boldsymbol{V}_{t}^{-} \right)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11}^{2} - \frac{8.1}{\mathsf{rk}_{\star}^{\alpha} \log d} \boldsymbol{\Lambda}_{\ell_{2},11} \right) + \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \end{split}$$

• To derive an upper-bound for \underline{G}_t , assuming $\underline{G}_t \preceq 1.1 I_{{\sf rk}_\star}$, we have

$$\begin{split} \underline{V}_{t+1} &= \underline{V}_{t} - \frac{\eta}{1 - 1.1\eta} \underline{V}_{t}^{2} + \frac{\eta^{2}}{(1 - 1.1\eta)^{2}} \underline{V}_{t}^{3} \left(\underline{I}_{\mathsf{rk}_{\star}} + \frac{\eta}{1 - 1.1\eta} \underline{V}_{t} \right)^{-1} \\ &+ \frac{\eta}{1 - 1.1\eta} \left(\mathbf{\Lambda}_{\ell_{1},11}^{2} - \frac{8.1}{\mathsf{rk}_{\star}^{\alpha} \log d} \mathbf{\Lambda}_{\ell_{2},11} - \frac{O(1)}{\log^{2} d} \mathbf{\Lambda}_{\ell_{2},11}^{2} \right) \\ &\leq \underline{V}_{t} - \frac{\eta}{1 - 1.1\eta} \underline{V}_{t}^{2} + \frac{\eta}{1 - 1.1\eta} \mathbf{\Lambda}_{\ell_{1},11}^{2}. \end{split}$$

Then, by Proposition 34, we have $\underline{G}_{t+1} \leq 1.1 I_{rk_{\star}}$. Since the bound holds for t = 0, it holds for all $t \in \mathbb{N}$.

• To derive a lower-bound, we first observe that by monotonicity $\underline{V}_0 \succ 0$. Therefore, assuming $V_4 \succ V_0$

$$\begin{split} \underline{\boldsymbol{V}}_{t+1} \succeq \underline{\boldsymbol{V}}_{t} - \frac{\eta/2}{1 - 1.1\eta} (2\boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \underline{\boldsymbol{G}}_{t} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} - \boldsymbol{\Lambda}_{\ell_{1},11})^{2} \\ + \frac{\eta/2}{1 - 1.1\eta} \left(\boldsymbol{\Lambda}_{\ell_{1},11}^{2} - \frac{8.1}{\mathsf{rk}_{\star}^{\alpha} \log d} \boldsymbol{\Lambda}_{\ell_{2},11} - \frac{O(1)}{\log^{2} d} \boldsymbol{\Lambda}_{\ell_{2},11}^{2} \right) \end{split}$$

$$\begin{split} & \succeq \underline{V}_t - \frac{\eta/2}{1 - 1.1\eta} \left(2 \mathbf{\Lambda}_{\ell_2, 11}^{\frac{1}{2}} \underline{G}_t \mathbf{\Lambda}_{\ell_2, 11}^{\frac{1}{2}} - \sqrt{\left(\mathbf{\Lambda}_{\ell_1, 11}^2 - \frac{8.1}{\mathsf{rk}_\star^\alpha \log d} \mathbf{\Lambda}_{\ell_2, 11} - \frac{O(1)}{\log^2 d} \mathbf{\Lambda}_{\ell_2, 11}^2 \right)} \right)^2 \\ & + \frac{\eta/2}{1 - 1.1\eta} \left(\mathbf{\Lambda}_{\ell_1, 11}^2 - \frac{8.1}{\mathsf{rk}_\star^\alpha \log d} \mathbf{\Lambda}_{\ell_2, 11} - \frac{O(1)}{\log^2 d} \mathbf{\Lambda}_{\ell_2, 11}^2 \right). \end{split}$$

Then, by Proposition 34, we have $\underline{V}_t \succeq \underline{V}_0$.

We start our proof by showing that $V_t^- \succeq \underline{V}_t + \underline{\zeta}_t$ for $t \leq \mathcal{T}_{\text{stable}}$. Assuming the statement holds for $t \in \mathbb{N}$, we have

$$\begin{split} \boldsymbol{V}_{t+1}^{-} \succeq & (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t}) \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t}) \right)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11}^{2} - \frac{8.1}{\mathsf{rk}_{\star}^{\alpha} \log d} \boldsymbol{\Lambda}_{\ell_{2},11} \right) \\ & + \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \\ & = \underline{\boldsymbol{V}}_{t} \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t} \right)^{-1} + \eta \left(\boldsymbol{\Lambda}_{\ell_{1},11}^{2} - \frac{8.1}{\mathsf{rk}_{\star}^{\alpha} \log d} \boldsymbol{\Lambda}_{\ell_{2},11} \right) \\ & + \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t} \right)^{-1} \underline{\boldsymbol{\zeta}}_{t} \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t}) \right)^{-1} + \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \boldsymbol{\nu}_{t+1,11} \boldsymbol{\Lambda}_{\ell_{2},11}^{\frac{1}{2}} \end{split}$$

We have for $t < \mathcal{T}_{\text{stable}}$

$$\begin{split} &\left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t}\right)^{-1} \underline{\boldsymbol{\zeta}}_{t} \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t})\right)^{-1} \\ &= \underline{\boldsymbol{\zeta}}_{t} - \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t} \underline{\boldsymbol{\zeta}}_{t} - \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{\zeta}}_{t} \underline{\boldsymbol{V}}_{t} - \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{\zeta}}_{t}^{2} \\ &- \frac{\eta^{2}}{(1 - 1.1\eta)^{2}} \underline{\boldsymbol{V}}_{t} \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} \underline{\boldsymbol{V}}_{t}\right)^{-1} \underline{\boldsymbol{\zeta}}_{t} \left(\boldsymbol{I}_{r} + \frac{\eta}{1 - 1.1\eta} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t})\right)^{-1} (\underline{\boldsymbol{V}}_{t} + \underline{\boldsymbol{\zeta}}_{t}) \\ &\succeq \underline{\boldsymbol{\zeta}}_{t} - \frac{\eta}{1 - 1.1\eta} \frac{1}{\log^{2} d} \underline{\boldsymbol{V}}_{t}^{2} - \frac{\eta}{1 - 1.1\eta} (1 + \log^{2} d) \underline{\boldsymbol{\zeta}}_{t}^{2} - \frac{\eta^{2}}{(1 - 1.1\eta)^{2}} \frac{O(rk_{\star}^{-\alpha})}{\log^{4} d} \boldsymbol{I}_{\mathsf{rk}_{\star}} \\ &\succeq \underline{\boldsymbol{\zeta}}_{t} - \frac{\eta}{1 - 1.1\eta} \frac{O(1)}{\log^{2} d} \boldsymbol{\Lambda}_{\ell_{2},11}^{2}. \end{split}$$

Since $V_0^-=\underline{V}_0+\underline{\zeta}_0$, the claim follows. Then, by the third item above, we have for $t\leq \mathcal{T}_{\text{stable}}$

$$\boldsymbol{G}_t \succeq \underline{\boldsymbol{G}}_0 + \underline{\boldsymbol{\nu}}_t \succeq \boldsymbol{I}_{\mathsf{rk}_\star} - \frac{1}{\log d} \boldsymbol{I}_{\mathsf{rk}_\star} - \frac{\mathsf{rk}_\star^{\frac{-\alpha}{2}}}{\log^2 d} \boldsymbol{I}_{\mathsf{rk}_\star} \succeq \boldsymbol{I}_{\mathsf{rk}_\star} - \frac{2}{\log d} \boldsymbol{I}_{\mathsf{rk}_\star} \Rightarrow \mathcal{T}_{\mathsf{noise}} \leq \mathcal{T}_{\mathsf{bounded}}.$$

In the following, we will bound \mathcal{T}_{noise} . We have

$$\mathbb{E}_{t} \left[\underline{\boldsymbol{\nu}}_{t+1}^{2} \right] \stackrel{(a)}{\leq} \underline{\boldsymbol{\nu}}_{t}^{2} + O(\eta^{2}) \boldsymbol{I}_{\mathsf{rk}_{\star}} \Rightarrow \mathbb{E} \left[\underline{\boldsymbol{\nu}}_{t}^{2} \right] \leq O(\eta^{2} t) \boldsymbol{I}_{\mathsf{rk}_{\star}}$$

By clipping strategy we used with $L=\log^2 d$ in (F.38), and defining $\Gamma_1:=\mathbf{I}_{\mathsf{rk}_\star}, \Gamma_2:=\mathbf{\Lambda}_{11}^{\frac{1}{2}},$ and

$$\mathrm{Quad}_{k,t}^{(\ell)}(\mathsf{X}) := \sum_{j=1}^k \mathbb{E}_{j-1} \left[\left(\mathsf{\Gamma}_\ell T_t^{\frac{-1}{2}} \mathsf{X}_j T_t^{\frac{-1}{2}} \mathsf{\Gamma}_\ell \right)^2 \right], \; \ell \in \{1,2\},$$

we can show that the following events hold: For any $T \in \mathbb{N}$,

$$\begin{split} \widehat{\mathcal{E}}_{\text{ht},1} &\equiv \bigg\{ \max_{t \leq T} \lVert \mathsf{Quad}_{t,t}^{(1)} \rVert_2 \leq O(\eta^2 T) & \text{and} \ \max_{t \leq T} R_{t,t}^{(1)} \leq O(\eta \mathsf{rk}_{\star}^{\frac{1}{2}} \log^2 d) \bigg\} \\ \widehat{\mathcal{E}}_{\text{ht},2} &\equiv \bigg\{ \max_{t < T} \lVert \mathsf{Quad}_{t,t}^{(2)} \rVert_2 \leq O_{\alpha}(\eta^2 T \mathsf{rk}_{\star}^{\mathsf{p}_2 - 1}) \ \text{and} \ \max_{t < T} R_{t,t}^{(2)} \leq O_{\alpha}(\eta \mathsf{rk}_{\star}^{\mathsf{p}_2 - \frac{1}{2}} \log^2 d) \bigg\}. \end{split}$$

where p_2 is defined in Corollary 4. By using Proposition 19, we can show that with probability $d^4 \exp(-\log^2 d)$, $\mathcal{T}_{\text{noise}} \geq \frac{\mathsf{rk}_\star^\alpha \log^2 d}{\eta}$.

G Auxiliary Statements

G.1 Matrix bounds

Proposition 22. For $A, B \in \mathbb{R}^{d \times r}$, we have

$$-A^{\top}A - B^{\top}B \preceq A^{\top}B + B^{\top}A \preceq A^{\top}A + B^{\top}B$$

If r = d, then $(\mathbf{A} + \mathbf{A}^{\top})^2 \leq 2\mathbf{A}^{\top}\mathbf{A} + 2\mathbf{A}\mathbf{A}^{\top}$. Moreover, if $\mathbf{A}_1, \dots, \mathbf{A}_k$ are symmetric matrices,

$$\left(\sum_{i=1}^k \boldsymbol{A}_i\right)^2 \leq k \sum_{i=1}^k \boldsymbol{A}_i^2.$$

Proof. We have

$$(A - B)^{\top}(A - B) \succeq 0 \Rightarrow A^{\top}A + B^{\top}B \succeq A^{\top}B + B^{\top}A$$

By using $A \leftarrow -A$, we obtain the left inequality too. For the second inequality, we have

$$(\boldsymbol{A} + \boldsymbol{A}^{\top})^2 = \boldsymbol{A}^{\top} \boldsymbol{A} + \boldsymbol{A} \boldsymbol{A}^{\top} + \boldsymbol{A} \boldsymbol{A} + \boldsymbol{A}^{\top} \boldsymbol{A}^{\top}$$
$$(\boldsymbol{A} - \boldsymbol{A}^{\top})^{\top} (\boldsymbol{A} - \boldsymbol{A}^{\top}) = \boldsymbol{A}^{\top} \boldsymbol{A} + \boldsymbol{A} \boldsymbol{A}^{\top} - \boldsymbol{A} \boldsymbol{A} - \boldsymbol{A}^{\top} \boldsymbol{A}^{\top}$$

Therefore, $({\pmb A}+{\pmb A}^{\top})^2 \preceq 2\left({\pmb A}^{\top}{\pmb A}+{\pmb A}{\pmb A}^{\top}\right)$. For the last statement,

$$\left(\sum_{i=1}^k A_i\right)^2 = \sum_{i=1}^k A_i^2 + \sum_{i=1}^k \sum_{j=i+1}^k A_i A_j + \sum_{i=1}^k \sum_{j=i+1}^k A_j A_i \leq k \sum_{i=1}^k A_i^2,$$

where we use the first statement in the last inequality.

Proposition 23. Consider a symmetric square matrix with block partition

$$oldsymbol{M} = egin{bmatrix} oldsymbol{A} & oldsymbol{B} \ oldsymbol{B}^ op & oldsymbol{C} \end{bmatrix}.$$

If A is invertible, then M > 0 if and only if A > 0 and $C - B^{\top}A^{-1}B > 0$.

Proof. If A is invertible, we have

$$\begin{bmatrix} \boldsymbol{A} & \boldsymbol{B} \\ \boldsymbol{B}^\top & \boldsymbol{C} \end{bmatrix} = \begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \\ \boldsymbol{B}^\top \boldsymbol{A}^{-1} & \boldsymbol{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{A} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{C} - \boldsymbol{B}^\top \boldsymbol{A}^{-1} \boldsymbol{B} \end{bmatrix} \begin{bmatrix} \boldsymbol{I} & \boldsymbol{A}^{-1} \boldsymbol{B} \\ \boldsymbol{0} & \boldsymbol{I} \end{bmatrix}.$$

Note that

$$\begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \\ \boldsymbol{B}^{\top} \boldsymbol{A}^{-1} & \boldsymbol{I} \end{bmatrix}^{-1} = \begin{bmatrix} \boldsymbol{I} & \boldsymbol{0} \\ -\boldsymbol{B}^{\top} \boldsymbol{A}^{-1} & \boldsymbol{I} \end{bmatrix}.$$

Therefore, the statement follows.

Proposition 24. Let $r_u < r$ and $\mathbf{Z} \in \mathbb{R}^{r \times r_s}$ such that

$$m{Z} = egin{bmatrix} m{Z}_1 \ m{Z}_2 \end{bmatrix}, \ \ ext{where} \ \ m{Z}_1 \in \mathbb{R}^{r_u imes r_s}, m{Z}_2 \in \mathbb{R}^{r-r_u imes r_s}.$$

For ant $0 \le \varepsilon < 1$

$$\begin{split} \boldsymbol{Z}\boldsymbol{Z}^{\top} \succeq \varepsilon \begin{bmatrix} \boldsymbol{Z}_{1}\boldsymbol{Z}_{1}^{\top} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix} + (1-\varepsilon) \begin{bmatrix} \boldsymbol{Z}_{1}\boldsymbol{Z}_{1}^{\top} - \boldsymbol{Z}_{1}\boldsymbol{Z}_{2}^{\top}(\boldsymbol{Z}_{2}\boldsymbol{Z}_{2}^{\top})^{+}\boldsymbol{Z}_{2}\boldsymbol{Z}_{1}^{\top} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} \end{bmatrix} \\ -\frac{\varepsilon}{1-\varepsilon} \begin{bmatrix} \boldsymbol{0} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{Z}_{2}\boldsymbol{Z}_{2}^{\top} \end{bmatrix}, \end{split}$$

where $A o A^+$ denotes the pseudo inverse operator.

Proof. We will denote $x \in \mathbb{R}^r$ as

$$oldsymbol{x} = egin{bmatrix} oldsymbol{x}_1 \ oldsymbol{x}_2 \end{bmatrix} \; ext{ where } \; oldsymbol{x}_1 \in \mathbb{R}^{r_u}, oldsymbol{x}_2 \in \mathbb{R}^{r-r_u}.$$

We have

$$\boldsymbol{x}^{\top} \boldsymbol{Z} \boldsymbol{Z}^{\top} \boldsymbol{x} = \left(\boldsymbol{x}_{1}^{\top} \boldsymbol{Z}_{1} \boldsymbol{Z}_{1}^{\top} \boldsymbol{x}_{1} + 2 \boldsymbol{x}_{1}^{\top} \boldsymbol{Z}_{1} \boldsymbol{Z}_{2}^{\top} \boldsymbol{x}_{2} + \frac{1}{1 - \varepsilon} \boldsymbol{x}_{2}^{\top} \boldsymbol{Z}_{2} \boldsymbol{Z}_{2}^{\top} \boldsymbol{x}_{2} \right) - \frac{\varepsilon}{1 - \varepsilon} \boldsymbol{x}_{2}^{\top} \boldsymbol{Z}_{2} \boldsymbol{Z}_{2}^{\top} \boldsymbol{x}_{2}$$

$$\stackrel{(a)}{\geq} \left(\boldsymbol{x}_{1}^{\top} \boldsymbol{Z}_{1} \boldsymbol{Z}_{1}^{\top} \boldsymbol{x}_{1} - (1 - \varepsilon) \boldsymbol{x}_{1}^{\top} \boldsymbol{Z}_{1} \boldsymbol{Z}_{2}^{\top} (\boldsymbol{Z}_{2} \boldsymbol{Z}_{2}^{\top})^{+} \boldsymbol{Z}_{2} \boldsymbol{Z}_{1}^{\top} \boldsymbol{x}_{1} \right) - \frac{\varepsilon}{1 - \varepsilon} \boldsymbol{x}_{2}^{\top} \boldsymbol{Z}_{2} \boldsymbol{Z}_{2}^{\top} \boldsymbol{x}_{2},$$

where we minimized the first term in the first line over x_2 in (a). Since (a) holds for all x, the statement follows,

Proposition 25. Let $A \in \mathbb{R}^{r \times r}$ be a symmetric matrix. For $S \succ -A$, $S \rightarrow -(S + A)^{-1}$ is monotone.

Proof. Let $S_1 \succ S_2 \succ -A$. We have

$$-(S_1 + A)^{-1} + (S_2 + A)^{-1} = (S_2 + A)^{-1}((S_1 - S_2)^{-1} + (S_2 + A)^{-1})^{-1}(S_2 + A)^{-1} > 0.$$
(G.1)

For
$$S_1 \succeq S_2$$
, we can use $S_1 + \varepsilon I_r$ in (G.1) and take $\varepsilon \downarrow 0$

G.1.1 Additional bounds for continuous-time analysis

Proposition 26. For a symmetric positive definite D_1 , D_2 , and C > 0, we have

$$egin{split} m{D}_1 \Big(m{D}_1 + m{Z}_1 ig(Cm{I}_{r_s} + m{Z}_2^ op m{D}_2^{-1}m{Z}_2ig)^{-1}m{Z}_1^ op \Big)^{-1}m{Z}_1m{Z}_2^ op ig(m{Z}_2m{Z}_2^ op + Cm{D}_2ig)^{-1}m{D}_2 \ &= m{Z}_1 ig(Cm{I}_{r_s} + m{Z}_2^ op m{D}_2^{-1}m{Z}_2 + m{Z}_1^ op m{D}_2^{-1}m{Z}_1ig)^{-1}m{Z}_2^ op. \end{split}$$

Proof. We have

$$\begin{split} & \big(\boldsymbol{D}_1 + \boldsymbol{Z}_1 \big(\boldsymbol{C} \boldsymbol{I}_{r_s} + \boldsymbol{Z}_2^\top \boldsymbol{D}_2^{-1} \boldsymbol{Z}_2 \big)^{-1} \boldsymbol{Z}_1^\top \big)^{-1} \\ & = \boldsymbol{D}_1^{-1} - \boldsymbol{D}_1^{-1} \boldsymbol{Z}_1 \big(\boldsymbol{C} \boldsymbol{I}_{r_s} + \boldsymbol{Z}_2^\top \boldsymbol{D}_2^{-1} \boldsymbol{Z}_2 + \boldsymbol{Z}_1^\top \boldsymbol{D}_1^{-1} \boldsymbol{Z}_1 \big)^{-1} \boldsymbol{Z}_1^\top \boldsymbol{D}_1^{-1}. \end{split}$$

Therefore,

$$\begin{split} & \boldsymbol{D}_{1} \big(\boldsymbol{D}_{1} + \boldsymbol{Z}_{1} \big(\boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} \big)^{-1} \boldsymbol{Z}_{1}^{\top} \big)^{-1} \boldsymbol{Z}_{1} \\ & = \boldsymbol{Z}_{1} \Big(\boldsymbol{I}_{r_{s}} - \big(\boldsymbol{C} \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} + \boldsymbol{Z}_{1}^{\top} \boldsymbol{D}_{1}^{-1} \boldsymbol{Z}_{1} \big)^{-1} \boldsymbol{Z}_{1}^{\top} \boldsymbol{D}_{1}^{-1} \boldsymbol{Z}_{1} \Big) \\ & = \boldsymbol{Z}_{1} \big(\boldsymbol{C} \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} + \boldsymbol{Z}_{1}^{\top} \boldsymbol{D}_{1}^{-1} \boldsymbol{Z}_{1} \big)^{-1} \big(\boldsymbol{C} \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} \big) \end{split}$$

Then,

$$\begin{split} \boldsymbol{Z}_{1} \big(C \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} + \boldsymbol{Z}_{1}^{\top} \boldsymbol{D}_{1}^{-1} \boldsymbol{Z}_{1} \big)^{-1} \big(C \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} \big) \boldsymbol{Z}_{2}^{\top} \left(\boldsymbol{Z}_{2} \boldsymbol{Z}_{2}^{\top} + C \boldsymbol{D}_{2} \right)^{-1} \boldsymbol{D}_{2} \\ &= \boldsymbol{Z}_{1} \left(C \boldsymbol{I}_{r_{s}} + \boldsymbol{Z}_{2}^{\top} \boldsymbol{D}_{2}^{-1} \boldsymbol{Z}_{2} + \boldsymbol{Z}_{1}^{\top} \boldsymbol{D}_{1}^{-1} \boldsymbol{Z}_{1} \right)^{-1} \boldsymbol{Z}_{2}^{\top}. \end{split}$$

Proposition 27. For some diagonal positive definite $\mathbf{A} := \operatorname{diag}(\{a_j\}_{j=1}^{r_u})$ and $\mathbf{B} := \operatorname{diag}(\{b_j\}_{j=1}^{d-r_u})$, we let

$$oldsymbol{D}_1 \coloneqq rac{oldsymbol{A} \exp(-toldsymbol{A})}{oldsymbol{I}_{r_n} - \exp(-toldsymbol{A})}, \quad oldsymbol{D}_2 \coloneqq rac{oldsymbol{B} \exp(-toldsymbol{B})}{oldsymbol{I}_{d-r_n} - \exp(-toldsymbol{B})},$$

For some $\mathbf{Z}_1 \in \mathbb{R}^{r_u \times r_s}$, $\mathbf{Z}_2 \in \mathbb{R}^{(d-r_u) \times r_s}$, and C > 0, we define

$$M := \exp(0.5tA)Z_1 \left(CI_{r_s} + Z_2^{\top}D_2^{-1}Z_2 + Z_1^{\top}D_1^{-1}Z_1\right)^{-1}Z_2^{\top}\exp(0.5tB).$$

We have

$$\|\boldsymbol{M}\|_F^2 \leq \tilde{C} \sum_{i=1}^{r_u \wedge r_s} \left(\lambda_{\max}(\boldsymbol{Z}_1 \boldsymbol{Z}_1^\top) \exp(t(a_i + b_i)) \wedge \left(C + \frac{\lambda_{\min}(\boldsymbol{Z}_2^\top \boldsymbol{Z}_2)}{\lambda_{\max}(\boldsymbol{D}_2)} \right) \frac{a_i \exp(t(a_i + b_i))}{\exp(ta_i) - 1} \right)$$

where

$$\tilde{C} = \frac{\lambda_{\max}(\boldsymbol{Z}_2^{\top}\boldsymbol{Z}_2)}{\left(C + \frac{\lambda_{\min}(\boldsymbol{Z}_2^{\top}\boldsymbol{Z}_2)}{\lambda_{\max}(\boldsymbol{D}_2)}\right)^2}$$

Proof. For convenience, we will use

$$egin{aligned} ilde{m{D}}_1 &\coloneqq rac{m{A}}{m{I}_{r_u} - \exp(-tm{A})}, & ilde{m{D}}_2 &\equiv rac{m{B}}{m{I}_{d-r_u} - \exp(-tm{B})}, \ ilde{m{Z}}_1 &\coloneqq \exp(0.5tm{A})m{Z}_1, & ilde{m{Z}}_2 &\coloneqq \exp(0.5tm{B})m{Z}_2. \end{aligned}$$

We let

$$egin{aligned} oldsymbol{M}_1 &\coloneqq ilde{oldsymbol{Z}}_1 \left(C oldsymbol{I}_{r_s} + ilde{oldsymbol{Z}}_2^ op ilde{oldsymbol{D}}_2^{-1} ilde{oldsymbol{Z}}_2 + ilde{oldsymbol{Z}}_1^ op ilde{oldsymbol{D}}_1^{-1} ilde{oldsymbol{Z}}_1
ight)^{rac{-1}{2}} ilde{oldsymbol{Z}}_2^ op \ & oldsymbol{M}_2 centcolon = \left(C oldsymbol{I}_{r_s} + ilde{oldsymbol{Z}}_2^ op ilde{oldsymbol{D}}_2^{-1} ilde{oldsymbol{Z}}_2 + ilde{oldsymbol{Z}}_1^ op ilde{oldsymbol{D}}_1^{-1} ilde{oldsymbol{Z}}_1
ight)^{rac{-1}{2}} ilde{oldsymbol{Z}}_2^ op \ & oldsymbol{Z}_2^ op ilde{oldsymbol{D}}_1^ op ilde{$$

We observe that

$$\|oldsymbol{M}\|_F^2 = ext{Tr}(oldsymbol{M}_1^ op oldsymbol{M}_1 oldsymbol{M}_2 oldsymbol{M}_2^ op) \leq \sum_{i=1}^{r_u \wedge r_s} \lambda_i(oldsymbol{M}_1 oldsymbol{M}_1^ op) \lambda_i(oldsymbol{M}_2^ op oldsymbol{M}_2)$$

where we used that $\operatorname{rank}(\boldsymbol{M}_1\boldsymbol{M}_1^\top) \leq r_u \wedge r_s$ and Von Neumann's trace inequality in the last part. We have

$$\boldsymbol{M}_{2}^{\top} \boldsymbol{M}_{2} \leq \exp(0.5t\boldsymbol{B}) \boldsymbol{Z}_{2}^{\top} \left(C \boldsymbol{I}_{r_{s}} + \frac{1}{\lambda_{\max}(\boldsymbol{D}_{2})} \boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2} \right)^{-1} \boldsymbol{Z}_{2} \exp(0.5t\boldsymbol{B})$$
$$\leq \frac{\lambda_{\max}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{C + \frac{\lambda_{\max}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})}} \exp(t\boldsymbol{B})$$

On the other hand,

$$\begin{split} & \boldsymbol{M}_{1} \boldsymbol{M}_{1}^{\top} \preceq \tilde{\boldsymbol{Z}}_{1} \left(\left(\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})} \right) \boldsymbol{I}_{r_{s}} + \tilde{\boldsymbol{Z}}_{1}^{\top} \tilde{\boldsymbol{D}}_{1}^{-1} \tilde{\boldsymbol{Z}}_{1} \right)^{-1} \tilde{\boldsymbol{Z}}_{1}^{\top} \\ &= \frac{1}{\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})}} \tilde{\boldsymbol{Z}}_{1} \left(\boldsymbol{I}_{r_{s}} + \tilde{\boldsymbol{Z}}_{1}^{\top} \left(\left(\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})} \right) \tilde{\boldsymbol{D}}_{1} \right)^{-1} \tilde{\boldsymbol{Z}}_{1} \right)^{-1} \tilde{\boldsymbol{Z}}_{1}^{\top} \\ &= \frac{1}{\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})}} \tilde{\boldsymbol{Z}}_{1} \tilde{\boldsymbol{Z}}_{1}^{\top} \left(\left(\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})} \right) \tilde{\boldsymbol{D}}_{1} + \tilde{\boldsymbol{Z}}_{1} \tilde{\boldsymbol{Z}}_{1}^{\top} \right)^{-1} \left(\left(\boldsymbol{C} + \frac{\lambda_{\min}(\boldsymbol{Z}_{2}^{\top} \boldsymbol{Z}_{2})}{\lambda_{\max}(\boldsymbol{D}_{2})} \right) \tilde{\boldsymbol{D}}_{1} \right) \end{split}$$

We have the following at the same time:

•
$$\tilde{\mathbf{Z}}_1 \tilde{\mathbf{Z}}_1^{\top} \left(\left(C + \frac{\lambda_{\min}(\mathbf{Z}_2^{\top} \mathbf{Z}_2)}{\lambda_{\max}(\mathbf{D}_2)} \right) \tilde{\mathbf{D}}_1 + \tilde{\mathbf{Z}}_1 \tilde{\mathbf{Z}}_1^{\top} \right)^{-1} \left(C + \frac{\lambda_{\min}(\mathbf{Z}_2^{\top} \mathbf{Z}_2)}{\lambda_{\max}(\mathbf{D}_2)} \right) \tilde{\mathbf{D}}_1 \preceq \left(C + \frac{\lambda_{\min}(\mathbf{Z}_2^{\top} \mathbf{Z}_2)}{\lambda_{\max}(\mathbf{D}_2)} \right) \tilde{\mathbf{D}}_1$$

•
$$\tilde{Z}_1 \tilde{Z}_1^{\top} \left(\left(C + \frac{\lambda_{\min}(Z_2^{\top} Z_2)}{\lambda_{\max}(D_2)} \right) \tilde{D}_1 + \tilde{Z}_1 \tilde{Z}_1^{\top} \right)^{-1} \left(C + \frac{\lambda_{\min}(Z_2^{\top} Z_2)}{\lambda_{\max}(D_2)} \right) \tilde{D}_1 \leq \lambda_{\max}(Z_1 Z_1^{\top}) \exp(tA)$$

Therefore, for $i \leq r \wedge r_s$, we have

$$\lambda_i(\boldsymbol{M}_1\boldsymbol{M}_1^\top) \leq \frac{1}{C + \frac{\lambda_{\min}(\boldsymbol{Z}_2^\top \boldsymbol{Z}_2)}{\lambda_{\max}(\boldsymbol{D}_2)}} \left(\lambda_{\max}(\boldsymbol{Z}_1\boldsymbol{Z}_1^\top) \exp(ta_i) \wedge \left(C + \frac{\lambda_{\min}(\boldsymbol{Z}_2^\top \boldsymbol{Z}_2)}{\lambda_{\max}(\boldsymbol{D}_2)}\right) \frac{a_i \exp(ta_i)}{\exp(ta_i) - 1}\right)$$

Therefore,

$$\|\boldsymbol{M}\|_F^2 \leq \tilde{C} \sum_{i=1}^{r_u \wedge r_s} \left(\lambda_{\max}(\boldsymbol{Z}_1 \boldsymbol{Z}_1^\top) \exp(t(a_i + b_i)) \wedge \left(C + \frac{\lambda_{\min}(\boldsymbol{Z}_2^\top \boldsymbol{Z}_2)}{\lambda_{\max}(\boldsymbol{D}_2)} \right) \frac{a_i \exp(t(a_i + b_i))}{\exp(ta_i) - 1} \right).$$

G.1.2 Additional bounds for discrete-time analysis

Proposition 28. For some positive definite diagonal matrices $D_0, D_1 \in \mathbb{R}^{r \times r}$ and symmetric matrices $G, \nu \in \mathbb{R}^{r \times r}$, we let

$$V\coloneqq 2D_0^{\frac{1}{2}}GD_0^{\frac{1}{2}}-D_1$$
 and $\zeta\coloneqq D_0^{\frac{1}{2}}
u D_0^{\frac{1}{2}}$ and $\grave{V}\coloneqq V+\zeta,$

where

- $\|G\|_2 \le L_G$ and $\|\nu\|_2 \le L_{\nu}$ and $\|D_0\|_2 \le L_0$.
- $\|\boldsymbol{D}_0^{-1}\boldsymbol{D}_1\|_2 \leq L_{1/0}$ and $\|\boldsymbol{D}_0\boldsymbol{D}_1^{-1}\|_2 \leq L_{0/1}$.
- For notational convenience, let $L_F := 2L_G + L_{1/0}$ and $L_{\dot{\mathbf{r}}} := 2(L_G + L_{\nu}) + L_{1/0}$.

For $0 \le \eta < \frac{1}{L_{\dot{F}}L_0}$, we have that $(I_r + \eta V)$ and $(I_r + \eta \dot{V})$ are invertible and the following bounds holds:

$$-C_{1}\boldsymbol{D}_{1} \leq \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} \leq C_{1}\boldsymbol{D}_{1}, \quad -C_{2}\boldsymbol{D}_{1} \leq \boldsymbol{V} \boldsymbol{\zeta} \dot{\boldsymbol{V}}^{2} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} \leq C_{2}\boldsymbol{D}_{1}$$

$$(G.2)$$

$$-C_{3}\boldsymbol{D}_{1} \leq \boldsymbol{V}^{3} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \leq C_{3}\boldsymbol{D}_{1}, \qquad -C_{4}\boldsymbol{D}_{1} \leq \boldsymbol{\zeta} \dot{\boldsymbol{V}}^{3} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} \leq C_{4}\boldsymbol{D}_{1},$$

$$(G.3)$$

where

$$C_{1} = \frac{L_{\nu}L_{0/1}L_{F}^{2}L_{\dot{F}}L_{0}^{3}}{\left(1 - \eta L_{F}L_{0}\right)\left(1 - \eta L_{\dot{F}}L_{0}\right)}, \quad C_{2} = \frac{L_{\nu}L_{0/1}L_{F}L_{\dot{F}}^{2}L_{0}^{3}}{1 - \eta L_{\dot{F}}L_{0}}$$

$$C_{3} = \frac{L_{\nu}L_{0/1}L_{F}^{3}L_{0}^{3}}{1 - \eta L_{F}L_{0}}, \quad C_{4} = \frac{L_{\nu}L_{0/1}L_{\dot{F}}^{3}L_{0}^{3}}{1 - \eta L_{\dot{F}}L_{0}}.$$

Proof. Note that $\|V\|_2 \vee \|\dot{V}\|_2 \leq L_{\dot{F}}L_0$, therefore, if $0 \leq \eta < \frac{1}{L_{\dot{F}}L_0}$, $(I_r + \eta V)$ and $(I_r + \eta \dot{V})$ are invertible. For the following, we introduce the notation

$$\grave{\boldsymbol{G}}\coloneqq \boldsymbol{G}+\boldsymbol{\nu}$$
 and $\boldsymbol{F}=2\boldsymbol{G}+\boldsymbol{D}_0^{-1}\boldsymbol{D}_1$ and $\grave{\boldsymbol{F}}=2\grave{\boldsymbol{G}}+\boldsymbol{D}_0^{-1}\boldsymbol{D}_1.$

Note that we have $\|F\|_2 \le L_F$ and $\|\hat{F}\|_2 \le L_F$. For the left part of (G.2), we write

$$\begin{split} \boldsymbol{D}_{0}^{-\frac{1}{2}} \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} \boldsymbol{D}_{0}^{-\frac{1}{2}} \\ &= \boldsymbol{F} \boldsymbol{D}_{0} \boldsymbol{F} \boldsymbol{D}_{0} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{F} \boldsymbol{D}_{0}\right)^{-1} \boldsymbol{\nu} \boldsymbol{D}_{0} \dot{\boldsymbol{F}} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{D}_{0} \dot{\boldsymbol{F}}\right)^{-1}. \end{split}$$

Therefore, we have

$$\begin{split} \left\| \boldsymbol{F} \boldsymbol{D}_{0} \boldsymbol{F} \boldsymbol{D}_{0} \big(\boldsymbol{I}_{r} + \eta \boldsymbol{F} \boldsymbol{D}_{0} \big)^{-1} \boldsymbol{\nu} \boldsymbol{D}_{0} \dot{\boldsymbol{F}} \big(\boldsymbol{I}_{r} + \eta \boldsymbol{D}_{0} \dot{\boldsymbol{F}} \big)^{-1} \right\|_{2} \\ & \leq \frac{\| \boldsymbol{F} \|_{2}^{2} \| \dot{\boldsymbol{F}} \|_{2} \| \boldsymbol{D}_{0} \|_{2}^{3} \| \boldsymbol{\nu} \|_{2}}{\left(1 - \eta \| \boldsymbol{F} \|_{2} \| \boldsymbol{D}_{0} \|_{2} \right) \left(1 - \eta \| \dot{\boldsymbol{F}} \|_{2} \| \boldsymbol{D}_{0} \|_{2} \right)} \\ & \leq \frac{L_{F}^{2} L_{\dot{F}} L_{0}^{3} L_{\nu}}{(1 - \eta L_{F} L_{0}) (1 - \eta L_{\dot{F}} L_{0})}. \end{split}$$

Therefore, we have the bound. For the right part of (G.2), we write

$$m{D}_0^{-rac{1}{2}}m{V}m{\zeta}\hat{m{V}}^2\left(m{I}_r+\eta\hat{m{V}}
ight)^{-1}m{D}_0^{-rac{1}{2}}=m{F}m{D}_0m{
u}m{D}_0\hat{m{F}}m{D}_0\hat{m{F}}\left(m{I}_r+\etam{D}_0\hat{m{F}}
ight)^{-1}$$

Therefore, we have

$$\left\| \boldsymbol{F} \boldsymbol{D}_0 \boldsymbol{\nu} \boldsymbol{D}_0 \dot{\boldsymbol{F}} \boldsymbol{D}_0 \dot{\boldsymbol{F}} \left(\boldsymbol{I}_r + \eta \boldsymbol{D}_0 \dot{\boldsymbol{F}} \right)^{-1} \right\|_2 \leq \frac{\| \boldsymbol{F} \|_2 \| \dot{\boldsymbol{F}} \|_2^2 \| \boldsymbol{D}_0 \|_2^3 \| \boldsymbol{\nu} \|_2}{1 - \eta \| \dot{\boldsymbol{F}} \|_2 \| \boldsymbol{D}_0 \|_2} \leq \frac{L_F L_F^2 L_0^3 L_\nu}{1 - \eta L_F L_0},$$

which gives us the bound. For the left part of (G.3), we write

$$m{D}_0^{-rac{1}{2}}m{V}^3\left(m{I}_r+\etam{V}
ight)^{-1}m{\zeta}m{D}_0^{-rac{1}{2}}=\left(m{F}m{D}_0
ight)^3\left(m{I}_r+\etam{F}m{D}_0
ight)^{-1}m{
u}$$

Therefore, we have

$$\left\| (\mathbf{F} \mathbf{D}_0)^3 \left(\mathbf{I}_r + \eta \mathbf{F} \mathbf{D}_0 \right)^{-1} \nu \right\|_2 \le \frac{\|\mathbf{F}\|_2^3 \|\mathbf{D}_0\|_2^3 \|\boldsymbol{\nu}\|_2}{1 - \eta \|\mathbf{F}\|_2 \|\mathbf{D}_0\|_2} \le \frac{L_{\nu} L_F^3 L_0^3}{1 - \eta L_F L_0},$$

which gives us the bound. The the right part of (G.3) can be derived similarly.

Proposition 29. Let $V, \dot{V} \in \mathbb{R}^{r \times r}$ be symmetric matrices such that $\dot{V} = V + \zeta$. We have

$$\dot{\boldsymbol{V}} \left(\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}}\right)^{-1} - \boldsymbol{V} \left(\boldsymbol{I}_r + \eta \boldsymbol{V}\right)^{-1} = \left(\boldsymbol{I}_r + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \left(\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}}\right)^{-1}.$$

Moreover, given that

$$M := \zeta - \eta V \zeta - \eta \zeta V - \eta \zeta^2 + \eta^2 \zeta V \zeta + \eta^2 \zeta^3 + \eta^2 V \zeta V$$

under the conditions of Proposition 28, we have for any $\kappa_d > 0$,

$$\begin{aligned} -\frac{2}{\kappa_d^2} \eta^2 \boldsymbol{\zeta}^2 - \eta^2 \kappa_d^2 \boldsymbol{V}^4 - \eta^2 \kappa_d^2 \boldsymbol{V} \boldsymbol{\zeta}^2 \boldsymbol{V} - C \eta^3 \boldsymbol{D}_1 &\leq \left(\boldsymbol{I}_r + \eta \boldsymbol{V} \right)^{-1} \boldsymbol{\zeta} \left(\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}} \right)^{-1} - \boldsymbol{M} \\ &\leq \frac{2}{\kappa_d^2} \eta^2 \boldsymbol{\zeta}^2 + \eta^2 \kappa_d^2 \boldsymbol{V}^4 + \eta^2 \kappa_d^2 \boldsymbol{V} \boldsymbol{\zeta}^2 \boldsymbol{V} + C \eta^3 \boldsymbol{D}_1 \end{aligned}$$

where $C = C_1 + C_2 + C_3 + C_4$, i.e., the sum of the constants given in Proposition 28.

Proof. We write

$$\dot{\boldsymbol{V}} (\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}})^{-1} - (\boldsymbol{I}_r + \eta \boldsymbol{V})^{-1} \boldsymbol{V}
= (\boldsymbol{I}_r + \eta \boldsymbol{V})^{-1} \Big((\boldsymbol{I}_r + \eta \boldsymbol{V}) (\boldsymbol{V} + \boldsymbol{\zeta}) - \boldsymbol{V} (\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}}) \Big) (\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}})^{-1}
= (\boldsymbol{I}_r + \eta \boldsymbol{V})^{-1} \boldsymbol{\zeta} (\boldsymbol{I}_r + \eta \dot{\boldsymbol{V}})^{-1}.$$

For the second part, we write

$$(\boldsymbol{I}_{r} + \eta \boldsymbol{V})^{-1} \zeta \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

$$= \left(\boldsymbol{I}_{r} - \eta \boldsymbol{V} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1}\right) \zeta \left(\boldsymbol{I}_{r} - \eta \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}\right)$$

$$= \zeta - \eta \boldsymbol{V} \left(\boldsymbol{I}_{r} - \eta \boldsymbol{V} + \eta^{2} \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1}\right) \zeta$$

$$- \eta \zeta \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} - \eta \dot{\boldsymbol{V}} + \eta^{2} \dot{\boldsymbol{V}}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1}\right) + \eta^{2} \boldsymbol{V} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \zeta \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

$$= \zeta - \eta \boldsymbol{V} \zeta - \eta \zeta \boldsymbol{V} - \eta \zeta^{2} + \eta^{2} \boldsymbol{V}^{2} \zeta + \eta^{2} \zeta \dot{\boldsymbol{V}}^{2} - \eta^{3} \boldsymbol{V}^{3} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \zeta$$

$$- \eta^{3} \zeta \dot{\boldsymbol{V}}^{3} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} + \eta^{2} \boldsymbol{V} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \zeta \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

We have

$$\eta^2 \boldsymbol{V}^2 \boldsymbol{\zeta} + \eta^2 \boldsymbol{\zeta} \dot{\boldsymbol{V}}^2 = \eta^2 \boldsymbol{V}^2 \boldsymbol{\zeta} + \eta^2 \boldsymbol{\zeta} (\boldsymbol{V} + \boldsymbol{\zeta})^2 = \underbrace{\eta^2 \boldsymbol{V}^2 \boldsymbol{\zeta} + \eta^2 \boldsymbol{\zeta} \boldsymbol{V}^2 + \eta^2 \boldsymbol{\zeta}^2 \boldsymbol{V}}_{:= \boldsymbol{M}_1} + \eta^2 \boldsymbol{\zeta} \boldsymbol{V} \boldsymbol{\zeta} + \eta^2 \boldsymbol{\zeta}^3.$$

Moreover,

$$\eta^{2} \boldsymbol{V} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

$$= \eta^{2} \boldsymbol{V} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} - \eta^{3} \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

$$= \eta^{2} \boldsymbol{V} \boldsymbol{\zeta} \dot{\boldsymbol{V}} - \eta^{3} \boldsymbol{V} \boldsymbol{\zeta} \dot{\boldsymbol{V}}^{2} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} - \eta^{3} \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

$$= \eta^{2} \boldsymbol{V} \boldsymbol{\zeta} \boldsymbol{V} + \underbrace{\eta^{2} \boldsymbol{V} \boldsymbol{\zeta}^{2}}_{:=\boldsymbol{M}_{2}} - \eta^{3} \boldsymbol{V} \boldsymbol{\zeta} \dot{\boldsymbol{V}}^{2} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1} - \eta^{3} \boldsymbol{V}^{2} \left(\boldsymbol{I}_{r} + \eta \boldsymbol{V}\right)^{-1} \boldsymbol{\zeta} \dot{\boldsymbol{V}} \left(\boldsymbol{I}_{r} + \eta \dot{\boldsymbol{V}}\right)^{-1}$$

By Proposition 22, we have

$$-2\eta^{2}\zeta^{2} - \eta^{2}V^{4} - \eta^{2}V\zeta^{2}V \leq M_{1} + M_{2} \leq 2\eta^{2}\zeta^{2} + \eta^{2}V^{4} + \eta^{2}V\zeta^{2}V.$$

Therefore by Proposition 28, we have

$$-\frac{2}{\kappa_d^2}\eta^2\boldsymbol{\zeta}^2 - \eta^2\kappa_d^2\boldsymbol{V}^4 - \eta^2\kappa_d^2\boldsymbol{V}\boldsymbol{\zeta}^2\boldsymbol{V} - C\eta^3\boldsymbol{D}_1 \leq (\boldsymbol{I}_r + \eta\boldsymbol{V})^{-1}\boldsymbol{\zeta}\left(\boldsymbol{I}_r + \eta\dot{\boldsymbol{V}}\right)^{-1} - \boldsymbol{M}$$
$$\leq \frac{2}{\kappa_d^2}\eta^2\boldsymbol{\zeta}^2 + \eta^2\kappa_d^2\boldsymbol{V}^4 + \eta^2\kappa_d^2\boldsymbol{V}\boldsymbol{\zeta}^2\boldsymbol{V} + C\eta^3\boldsymbol{D}_1.$$

Proposition 30. By using the notation in Proposition 28, we consider

$$\eta < rac{1}{L_F L_0}$$
 and $0 < arepsilon < rac{0.5/\eta}{L_F L_0} - 1$

Then,

$$V (I_r + \eta V)^{-1} - \varepsilon \eta V^2 \succeq V (I_r + \eta (1 + \varepsilon) V)^{-1} - 2.5\varepsilon \eta^2 C D_1$$
$$V (I_r + \eta V)^{-1} + \varepsilon \eta V^2 \preceq V (I_r + \eta (1 - \varepsilon) V)^{-1} + 1.5\varepsilon \eta^2 C D_1,$$

where $C = \frac{L_{0/1}L_F^3L_0^3}{1-\eta L_FL_0}$.

Proof. For the lower bound, we have

$$V (I_r + \eta V)^{-1} - \varepsilon \eta V^2$$

$$= V - (1 + \varepsilon) \eta V^2 + \eta^2 V^3 (I_r + \eta V)^{-1}$$

$$= V - (1 + \varepsilon) \eta V^2 + (1 + \varepsilon)^2 \eta^2 V^3 (I_r + (1 + \varepsilon) \eta V)^{-1}$$

$$- (2\varepsilon + \varepsilon^2) \eta^2 V^3 (I_r + \eta V)^{-1} + (1 + \varepsilon)^2 \varepsilon \eta^3 V^4 (I_r + (1 + \varepsilon) \eta V)^{-1} (I_r + \eta V)^{-1}$$

$$\succeq V (I_r + \eta (1 + \varepsilon) V)^{-1} - 2.5\varepsilon \eta^2 C D_1,$$

where we used C_3 with $L_{\nu}=1$ in Proposition 28 in the last step. For the upper bound,

$$V (I_r + \eta V)^{-1} + \varepsilon \eta V^2$$

$$= V - (1 - \varepsilon)\eta V^2 + \eta^2 V^3 (I_r + \eta V)^{-1}$$

$$= V - (1 - \varepsilon)\eta V^2 + (1 - \varepsilon)^2 \eta^2 V^3 (I_r + (1 - \varepsilon)\eta V)^{-1}$$

$$+ (2\varepsilon - \varepsilon^2)\eta^2 V^3 (I_r + \eta V)^{-1} - (1 - \varepsilon)^2 \varepsilon \eta^3 V^4 (I_r + (1 - \varepsilon)\eta V)^{-1} (I_r + \eta V)^{-1}$$

$$\leq V (I_r + \eta (1 + \varepsilon)V)^{-1} + 1.5\varepsilon \eta^2 C D_1.$$

Lemma 7. For any $\eta \in \mathbb{R}$ and $t \in \mathbb{N}$, we have

$$\begin{bmatrix} \boldsymbol{I}_r & \eta \boldsymbol{I}_r \\ \eta \boldsymbol{\Lambda}^2 & \boldsymbol{I}_r \end{bmatrix}^t = \begin{bmatrix} \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{2} & \boldsymbol{\Lambda}^{-1} \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{2} \\ \boldsymbol{\Lambda} \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t - (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{2} & \frac{(\boldsymbol{I}_r + \eta \boldsymbol{\Lambda})^t + (\boldsymbol{I}_r - \eta \boldsymbol{\Lambda})^t}{2} \end{bmatrix}.$$
(G.4)

Proof. We observe that

$$\mathbf{A} \coloneqq \begin{bmatrix} 0 & \mathbf{I}_r \\ \mathbf{\Lambda}^2 & 0 \end{bmatrix} \Rightarrow (\mathbf{G.4}) = \sum_{k=0}^t \binom{t}{k} \eta^k \mathbf{A}^k.$$

Note that

$$\boldsymbol{A}^{2k} = \begin{bmatrix} \boldsymbol{\Lambda}^{2k} & 0 \\ 0 & \boldsymbol{\Lambda}^{2k} \end{bmatrix} \ \text{ and } \ \boldsymbol{A}^{2k+1} = \begin{bmatrix} 0 & \boldsymbol{\Lambda}^{2k} \\ \boldsymbol{\Lambda}^{2k+2} & 0 \end{bmatrix}.$$

Therefore,

$$\begin{split} \sum_{k=0}^{t} \binom{t}{k} \eta^k \mathbf{A}^k &= \begin{bmatrix} \sum_{\substack{k=0 \\ k \text{ even}}}^{t} \binom{t}{k} \eta^k \mathbf{\Lambda}^k & \sum_{\substack{k=0 \\ k \text{ odd}}}^{t} \binom{t}{k} \eta^k \mathbf{\Lambda}^{k-1} \\ \sum_{\substack{k=0 \\ k \text{ odd}}}^{t} \binom{t}{k} \eta^k \mathbf{\Lambda}^{k+1} & \sum_{\substack{k=0 \\ k \text{ even}}}^{t} \binom{t}{k} \eta^k \mathbf{\Lambda}^k \end{bmatrix} \\ &= \begin{bmatrix} \frac{(\mathbf{I}_r + \eta \mathbf{\Lambda})^t + (\mathbf{I}_r - \eta \mathbf{\Lambda})^t}{2} & \mathbf{\Lambda}^{-1} \frac{(\mathbf{I}_r + \eta \mathbf{\Lambda})^t - (\mathbf{I}_r - \eta \mathbf{\Lambda})^t}{2} \\ \mathbf{\Lambda} \frac{(\mathbf{I}_r + \eta \mathbf{\Lambda})^t - (\mathbf{I}_r - \eta \mathbf{\Lambda})^t}{2} & \frac{(\mathbf{I}_r + \eta \mathbf{\Lambda})^t + (\mathbf{I}_r - \eta \mathbf{\Lambda})^t}{2} \end{bmatrix}. \end{split}$$

G.2 Some moment bounds and concentration inequalities

Lemma 8 (Hypercontractivity). Let $P_k : \mathbb{R}^d \to \mathbb{R}$ be a polynomial of degree-k and $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. For $q \geq 2$, we have $\mathbb{E}\left[P_k(\mathbf{x})^q\right]^{1/q} \leq (q-1)^{k/2}\mathbb{E}\left[P_k(\mathbf{x})^2\right]^{1/2}$.

Lemma 9. Let $x \sim \mathcal{N}(0, I_d)$ and $S \in \mathbb{R}^{d \times d}$ be a symmetric matrix. For u > 0,

$$\mathbb{P}\left[|\boldsymbol{x}^{\top}\boldsymbol{S}\boldsymbol{x} - Tr(\boldsymbol{S})| \ge 2\|\boldsymbol{S}\|_{F}u + 2\|\boldsymbol{S}\|_{2}u^{2}\right] \le 2e^{-u^{2}}.$$

Proof. We note that $\boldsymbol{x}^{\top} \boldsymbol{S} \boldsymbol{x} - \text{Tr}(\boldsymbol{S})$ has the same distribution with $\sum_{i=1}^{d} \lambda_i(\boldsymbol{S})(Z_i^2 - 1)$, where $Z_i \sim_{iid} \mathcal{N}(0,1)$. By using the Laurent-Massart lemma [LM00], we have the result.

Corollary 6. Let $y = \mathbf{x}^{\top} \mathbf{S} \mathbf{x} - Tr(\mathbf{S})$ and $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. For $p \geq 2$, we have $\mathbb{E}[|y|^p]^{\frac{1}{p}} \leq (p-1)\sqrt{2}\|\mathbf{S}\|_F$.

Proof. By observing that $\mathbb{E}[|y|^2] = 2\|S\|_F^2$, we have the result.

Corollary 7. For $\mathbf{A} \in \mathbb{R}^{d \times r}$, $p \geq 2$ and $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, we have $\mathbb{E}[\|\mathbf{A}^{\top}\mathbf{x}\|_2^{2p}]^{\frac{1}{p}} \leq \sqrt{3}(p-1)Tr(\mathbf{A}^{\top}\mathbf{A})$.

Proof. By Lemma 8, we have $\mathbb{E}[\|\boldsymbol{A}^{\top}\boldsymbol{x}\|_{2}^{2p}]^{\frac{1}{p}} \leq (p-1)\mathbb{E}[\|\boldsymbol{A}^{\top}\boldsymbol{x}\|_{2}^{4}]^{\frac{1}{2}}$. For $\boldsymbol{S} = \boldsymbol{A}\boldsymbol{A}^{\top}$, we have $\mathbb{E}[\|\boldsymbol{A}^{\top}\boldsymbol{x}\|_{2}^{4}] = \mathbb{E}[(\boldsymbol{x}^{\top}\boldsymbol{S}\boldsymbol{x})^{2}] = \operatorname{Tr}(\mathbb{E}[(\boldsymbol{x}^{\top}\boldsymbol{S}\boldsymbol{x})\boldsymbol{x}\boldsymbol{x}^{\top}]\boldsymbol{S})$.

We have

$$\mathbb{E}[(\boldsymbol{x}^{\top}\boldsymbol{S}\boldsymbol{x})\boldsymbol{x}\boldsymbol{x}^{\top}] = \text{Tr}(\boldsymbol{S})\boldsymbol{I}_d + 2\boldsymbol{S} \Rightarrow \mathbb{E}[\|\boldsymbol{A}^{\top}\boldsymbol{x}\|_2^4] = \text{Tr}(\boldsymbol{S})^2 + 2\|\boldsymbol{S}\|_F^2 \stackrel{(a)}{\leq} 3\text{Tr}(\boldsymbol{S})^2,$$

where (a) follows that S is positive semi-definite. Since $Tr(S) = Tr(A^{T}A)$, we have the statement.

Proposition 31. Let $x_j \sim_{i.i.d} \mathcal{N}(0, I_r)$, for $j \in [N]$. There exists a constant c > 0 such that for $\delta = \frac{u(r + \sqrt{Cr \log d} + C \log d)}{\sqrt{N}}$, we have

$$\mathbb{P}\left[\sup_{\substack{\boldsymbol{S} \in \mathbb{R}^{r \times r} \\ \|\boldsymbol{S}\|_{F} = 1}} \left| \frac{1}{N} \sum_{j=1}^{N} \frac{1}{2} Tr \left(\boldsymbol{S}(\boldsymbol{x}_{j} \boldsymbol{x}_{j}^{\top} - \boldsymbol{I}_{r})\right)^{2} - 1 \right| \ge \max\{2\delta, \delta^{2}\} + 10d^{-C/2}\right] \\
\le d^{2} \exp(-cu^{2}) + 2Nd^{-C}.$$

Proof. We observe that

$$\frac{1}{2} \| \boldsymbol{x}_j \boldsymbol{x}_j^{\top} - \boldsymbol{I}_r \|_F \le \frac{1}{2} (\| \boldsymbol{x}_j \|_2^2 + r).$$

By using Lemma 9, we can derive

$$\mathbb{P}\Big[\underbrace{\|\boldsymbol{x}_j\|_2^2 \leq r + 2\sqrt{r}\sqrt{C\log d} + 2C\log d}_{=:\mathcal{E}_j}\Big] \geq 1 - 2d^{-C}.$$

We have

$$\begin{split} \left| \mathbb{E} \left[\frac{1}{2} \text{Tr} \left(\boldsymbol{S} (\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_r) \right)^2 \mathbb{1}_{\mathcal{E}_j} \right] - 1 \right| &= \frac{1}{2} \mathbb{E} \left[\text{Tr} \left(\boldsymbol{S} (\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_r) \right)^2 \mathbb{1}_{\mathcal{E}_j^c} \right] \\ &\leq \frac{1}{2} \mathbb{E} \left[\text{Tr} \left(\boldsymbol{S} (\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_r) \right)^4 \right]^{1/2} \sqrt{2} d^{-C/2} \\ &< 9 \sqrt{2} d^{-C/2}. \end{split}$$

By using [Ver10, Theorem 5.41], for $\delta = \frac{u(r + \sqrt{Cr \log d} + C \log d)}{\sqrt{N}}$, we have

$$\mathbb{P}\left[\sup_{\substack{\boldsymbol{S} \in \mathbb{R}^{r \times r} \\ \|\boldsymbol{S}\|_{F} = 1}} \left| \frac{1}{N} \sum_{j=1}^{N} \frac{1}{2} \text{Tr} \left(\boldsymbol{S}(\boldsymbol{x}_{j} \boldsymbol{x}_{j}^{\top} - \boldsymbol{I}_{r})\right)^{2} - 1 \right| \ge \max\{2\delta, \delta^{2}\} + 10d^{-C/2} \right] \\
\le \mathbb{P}\left[\sup_{\substack{\boldsymbol{S} \in \mathbb{R}^{r \times r} \\ \|\boldsymbol{S}\|_{F} = 1}} \left| \frac{1}{N} \sum_{j=1}^{N} \frac{1}{2} \text{Tr} \left(\boldsymbol{S}(\boldsymbol{x}_{j} \boldsymbol{x}_{j}^{\top} - \boldsymbol{I}_{r})\right)^{2} \mathbb{1}_{\mathcal{E}_{j}} - \mathbb{E}\left[\frac{1}{2} \text{Tr} \left(\boldsymbol{S}(\boldsymbol{x}_{j} \boldsymbol{x}_{j}^{\top} - \boldsymbol{I}_{r})\right)^{2} \mathbb{1}_{\mathcal{E}_{j}}\right] \right| \ge \max\{2\delta, \delta^{2}\} \\
+ 2Nd^{-C} \\
\le d^{2} \exp(-cu^{2}) + 2Nd^{-C}.$$

Proposition 32. Let $\mathbf{x}_j \sim_{i.i.d} \mathcal{N}(0, \mathbf{I}_d)$, for $j \in [N]$, and $\mathbf{W} \in \mathbb{R}^{d \times r}$ be an orthonormal matrix. For a fixed $\mathbf{S} \in \mathbb{R}^{d \times d}$, $C \geq 16$ and $N \geq Cr \log d$, we have

$$\mathbb{P}\Big[\Big\|\frac{1}{N}\sum_{j=1}^{N}\frac{1}{2}Tr\big(\boldsymbol{S}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\big)\boldsymbol{W}^{\top}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\boldsymbol{W}-\boldsymbol{W}^{\top}\boldsymbol{S}\boldsymbol{W}\Big\|_{2}\geq 24e\|\boldsymbol{S}\|_{F}\Big(\sqrt{\frac{Cr}{N}}+d^{\frac{-C}{2}}\Big)\Big]$$

$$\leq 2e^{\frac{-Cr}{8}}+2Nd^{-C}.$$

Proof. Without loss of generality, we assume $||S||_F = 1$. By using Lemma 9, we have

$$\mathbb{P}\big[\underbrace{|\text{Tr}\big(\boldsymbol{S}(\boldsymbol{x}_j\boldsymbol{x}_j^\top - \boldsymbol{I}_d)\big)| \leq 4\sqrt{C\log d}}_{=:\mathcal{E}_j}\big] \geq 1 - 2d^{-C}.$$

For the following, we fix a $v \in S^{d-1}$. First, to bound the bias due to clipping, we write:

$$\begin{split} & \left| \mathbb{E} \left[\text{Tr} \left(\boldsymbol{S} (\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_d) \right) (\langle \boldsymbol{v}, \boldsymbol{x} \rangle^2 - 1) \mathbb{1}_{\mathcal{E}_j^c} \right] \right| \\ & \leq \mathbb{E} \left[\text{Tr} \left(\boldsymbol{S} (\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_d) \right)^4 \right]^{\frac{1}{4}} \mathbb{E} \left[(\langle \boldsymbol{v}, \boldsymbol{x} \rangle^2 - 1)^4 \right]^{\frac{1}{4}} \sqrt{2} d^{-C/2} \leq 18\sqrt{2} d^{-C/2} . \end{split}$$

On the other hand, to bound the moments of the clipped random variable, we have for $p \geq 2$,

$$\mathbb{E}\left[|\operatorname{Tr}\left(\boldsymbol{S}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\right)(\langle\boldsymbol{v},\boldsymbol{x}\rangle^{2}-1)|^{p}\mathbb{1}_{\mathcal{E}_{j}}\right]$$

$$\leq (4\sqrt{C\log d})^{p-2}\mathbb{E}\left[\operatorname{Tr}\left(\boldsymbol{S}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\right)^{2}|(\langle\boldsymbol{v},\boldsymbol{x}\rangle^{2}-1)|^{p}\right]\leq (12e)^{2}(8\sqrt{2}e\sqrt{C\log d})^{p-2}\frac{p!}{2}$$

By using ε -cover argument, we can derive

$$\mathbb{P}\left[\left\|\frac{1}{N}\sum_{j=1}^{N}\frac{1}{2}\mathrm{Tr}\left(\boldsymbol{S}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\right)\boldsymbol{W}^{\top}(\boldsymbol{x}_{j}\boldsymbol{x}_{j}^{\top}-\boldsymbol{I}_{d})\boldsymbol{W}-\boldsymbol{W}^{\top}\boldsymbol{S}\boldsymbol{W}\right\|_{2} \geq 24eu+18\sqrt{2}d^{-C/2}\right]$$

$$\leq 2\cdot 9^{r}\exp\left(\frac{-Nu^{2}/2}{1+u\sqrt{C\log d}}\right)+2Nd^{-C}.$$

By using $u = \sqrt{Cr/N}$, we have the result.

Proof. Without loss of generality, we assume $\|S\|_F = 1$. We have

$$\left\| \sum_{j=1}^N y_j(\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_d) - \boldsymbol{S} \right\|_F^2 \le \sup_{\boldsymbol{S} \in \mathbb{R}^{r \times r} \|\boldsymbol{S}\|_F = 1} \left| \frac{1}{N} \sum_{j=1}^N \frac{1}{2} \text{Tr} \left(\boldsymbol{S}(\boldsymbol{x}_j \boldsymbol{x}_j^\top - \boldsymbol{I}_r) \right)^2 - 1 \right|^2.$$

Hence, by considering the event in Proposition 31, we have the statement.

Proposition 33. Let $X \in \mathbb{R}$ be a random variable such that for some K, C > 0, $\mathbb{E}[|X|^p] \leq CK^pp^{pc}$ for some c > 0 and $p \geq k$. Then, $\mathbb{P}[|X| \geq Ku] \leq Ce^{-\frac{u^{1/c}}{c}}$ for $u \geq (ke)^c$.

Proof. Use Markov inequality with $p = \frac{u^{1/c}}{e}$.

G.3 Miscellaneous

Proposition 34. We consider $\eta \leq \frac{1}{10}$. The following statements holds:

• For $0.2 \ge \delta > 0$, let

$$u_{t+1} = u_t + \eta u_t (1 - u_t), \ 1 + \delta \ge u_0 \ge 0.$$

We have $1+\left(\delta\vee\frac{\eta^2}{4}\right)\geq \sup_t u_t\geq 0$. Moreover, $t^*=\inf\{t:u_t\geq 1\}$, we have $u_{t+1}\geq u_t$ for $t< t^*$ and $u_{t^*}\geq u_t\geq 1$ for $t\geq t^*$.

• For $0.5 > \varepsilon > 0$ and $1.1 > \overline{u}_0 \ge u_0 \ge \underline{u}_0 > 0$, let

$$\overline{u}_{t+1} = \overline{u}_t + \eta(1+\varepsilon)\overline{u}_t(1-\overline{u}_t)$$
 and $\underline{u}_{t+1} = \underline{u}_t + \eta\underline{u}_t(1-\underline{u}_t)$.

and

$$u_t + \eta u_t (1 - u_t) \le u_{t+1} \le u_t + \eta (1 + \varepsilon) u_t (1 - u_t).$$

We have

$$\frac{1}{2}\left(1 \wedge \underline{u}_0 e^{\frac{\eta t}{1+\eta}}\right) \leq \underline{u}_t \leq u_t \leq \overline{u}_t \leq \left(1.1 \wedge \overline{u}_0 e^{\eta(1+\varepsilon)t}\right).$$

Proof. If $t \geq t^*$, by monotonicity of the update, we have $1 \leq u_{t+1} \leq u_t \leq u_{t^*}$. If $t^* > 0$, then for $t < t^*$, we have $1 \geq u_t \geq 0$ and $u_t(1-u_t) \geq 0$, and thus, we have $u_{t+1} \geq u_t \geq 0$. Next, we observe that $u_{t^*} \leq 1 + 0.25\eta$ and by monotonicity of the update for $t \geq t^*$, we have $1 \leq u_t \leq u_{t^*}$. Hence, it is sufficient to bound u_{t^*} to bound $\sup_t u_t$. Note that, we have $1 \geq u_{t^*-1} \geq 1 - 0.25\eta$, and thus,

$$\frac{u_{t^*}}{u_{t^*-1}} = 1 + \eta(1 - u_{t^*-1}) \le 1 + \eta^2 \Rightarrow u_{t^*} \le 1 + \frac{\eta^2}{4}.$$

For the second item, by monotonicity, we have $0 < \underline{u}_t \le u_t \le \overline{u}_t < 1.1$. Moreover, by [AGP24, Lemma A.2], we have for $t < t_u := \inf\{t : \underline{u}_t \ge 0.5\}$

$$\frac{\underline{u}_0 e^{\frac{\eta t}{1+\eta}}}{1+\underline{u}_0 e^{\frac{\eta t}{1+\eta}}} \leq \underline{u}_t \Rightarrow \frac{\underline{u}_0}{2} e^{\frac{\eta t}{1+\eta}} \leq \underline{u}_0.$$

For $t \ge t_u$, by the first item, we have $\underline{u}_t \ge 0.5$. Therefore, we have

$$\frac{1}{2} \left(\underline{u}_0 e^{\frac{\eta t}{1+\eta}} \wedge 1 \right) \leq \underline{u}_t.$$

On the other hand, for all $t \in \mathbb{N}$, we have $\overline{u}_t \leq \overline{u}_0 e^{\eta(1+\varepsilon)t}$. By the first item, we have $\overline{u}_t \leq \overline{u}_0 e^{\eta(1+\varepsilon)t} \wedge 1.1$.

Proposition 35. For $t, \lambda > 0$, we have

$$\frac{1}{t \exp(t\lambda)} \le \frac{\lambda}{\exp(t\lambda) - 1} \le \frac{1}{t}.$$

Proof. The upper bound follows $\exp(t\lambda) - 1 \ge t\lambda$. For the lower bound,

$$\frac{1}{t} - \frac{\lambda}{\exp(t\lambda) - 1} = \frac{\exp(t\lambda) - t\lambda - 1}{t\left(\exp(t\lambda) - 1\right)}.$$
 (G.5)

 \Box

We have

$$\exp(t\lambda) - t\lambda - 1 \le \sum_{k=2}^{\infty} \frac{(t\lambda)^k}{k!} = t\lambda \sum_{k=1}^{\infty} \frac{(t\lambda)^k}{(k+1)!} \le t\lambda \sum_{k=1}^{\infty} \frac{(t\lambda)^k}{k!} = t\lambda \left(\exp(t\lambda) - 1\right).$$

Therefore.

$$(\mathbf{G.5}) \le \lambda \Rightarrow \frac{1}{t} \le \frac{\lambda}{\exp(t\lambda) - 1} + \lambda \Rightarrow \frac{1}{t \exp(t\lambda)} \le \frac{\lambda}{\exp(t\lambda) - 1}.$$

Lemma 10. Let $r_s \simeq d^{\gamma}$, $\gamma \in [0,1)$, and $\log^{-1} d \ll C_d \ll \log^{10} d$. We define F_d, G_d, H_d as

$$F_d(u) \coloneqq \left(1 - \frac{1}{1 + \left(\frac{dC_d}{r_s} \frac{1}{u} - 1\right) \left(\frac{d}{r_s}\right)^{-\frac{1}{u}}}\right)^2, \qquad G_d(u) \coloneqq 1 - \frac{1}{1 + \left(\frac{dC_d}{r_s} - 1\right) \left(\frac{d}{r_s}\right)^{-\frac{1}{u}}},$$

$$H_d(u) := \left(1 - C_d \left(\frac{d}{r_s}\right)^{\frac{1}{u} - 1}\right)_+,$$

We have

- For any C > 0, $\sup_{u < \log^C d} |F_d(u)| \le 1$ for $d \ge \Omega_C(1)$.
- $\sup_{u} |G_d(u)| \vee |H_d(u)| \leq 1$.
- For any $\delta \in (0,0.5)$, let $C_{\delta} := \{u \geq 0 : |u-1| < \delta\}$. For any compact $\mathcal{K} \subset (0,\infty] \setminus C_{\delta}$, we have $F_d(u) \xrightarrow{d \to \infty} \mathbb{1}\{u > 1\}$ uniformly on \mathcal{K} .
- For any compact $\mathcal{K} \subseteq [0, \infty] \setminus \mathcal{C}_{\delta}$, we have

$$G_d(u), H_d(u) \xrightarrow{d \to \infty} \mathbb{1}\{u > 1\}, \quad G_d^2(u) \xrightarrow{d \to \infty} \mathbb{1}\{u > 1\}$$

all uniformly on K.

Proof. For the first item, if $u \leq \log^C d$, for $d \geq \Omega_C(1)$

$$\frac{dC_d u}{r_s} - 1 \ge \frac{d}{r_s} \frac{t}{\log^{C+1} d} - 1 > 0.$$

Therefore, $|F_d(u)| \leq 1$. For the second item, since $\frac{dC_d}{r_s} > 1$ for $d \geq \Omega(1)$, the item follows.

For the third item, since $E \coloneqq [0,\infty) \setminus \mathcal{C}_\delta$ is closed in $[0,\infty)$, it suffices to establish the result on small open intervals around each point of of E within $[0,\infty)$. Fix $u_0 \in E$ and choose $\epsilon \in (0,\delta/2)$. Since $B(u_0,\epsilon) \coloneqq (u_0-\epsilon,u_0+\epsilon) \cap [0,\infty)$ is convex it can be either in $P_> \coloneqq \{u: u>1+\delta/2\}$ or $P_< \coloneqq \{u: u<1-\delta/2\}$. Without loss of generality let us assume it is in $P_<$. Then,

$$\sup_{u \in B(u_0, \epsilon) \subset P_{\leq}} |F_d(u)| \le 1 - \frac{1}{1 + \left(O_{\delta}(C_d) - \left(\frac{d}{r_s}\right)^{-1}\right) \left(\frac{d}{r_s}\right)^{-O_{\delta}(1)}} \to 0.$$

A similar step can be repeated if $B_{\epsilon} \subset P_{<}$.

For the last item, we first observe that uniform convergence of $G_d(u)$ implies the uniform convergence of $G_d^2(u)$. Therefore, we will only prove the first result. Since $E \coloneqq [0,\infty] \setminus \mathcal{C}_{\delta}$, is compact, and thus, $P_> \cap E$ and $P_< \cap E$ are also compact, we can directly use these sets. Without loss of generality let us use $P_< \cap E$. Then,

$$\sup_{u \in P_{\leq} \cap E} |G_d(u)| \vee |H_d(u)| \le \left(1 - C_d \frac{d}{r_s} C_{\delta}(1)\right)_+ \to 0.$$

A similar step can be repeated if $P_> \cap E$. Therefore, the statement follows.

Proposition 36. Let $r_u \leq r$ and

$$t \in \begin{cases} (0,\infty), & \alpha \in [0,0.5) \\ (0,\infty) \setminus \{j^\alpha : j \in \mathbb{N}\}, & \alpha > 0.5, \end{cases} \quad \kappa_{\text{eff}} \coloneqq \begin{cases} r^\alpha, & \alpha \in [0,0.5) \\ 1, & \alpha > 0.5. \end{cases}$$

We have

• For $K \in \{G, H\}$ and $t \neq \lim_{d \to \infty} \frac{1}{\lambda_i \kappa_{\text{eff}}}$, we have

$$\mathsf{K}_d(\frac{1}{\lambda_i t \kappa_{\text{eff}}}) - \mathbb{1}\{\frac{1}{\lambda_i} > t \kappa_{\text{eff}}\} = o_d(1).$$

• *For* $K \in \{F, G, H\}$,

$$\frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{i=1}^{r_u} \lambda_j^2 \left(\mathsf{K}_d(\frac{1}{\lambda_j t \kappa_{\mathrm{eff}}}) - \mathbb{1} \left\{ \frac{1}{\lambda_j} > t \kappa_{\mathrm{eff}} \right\} \right) = o_d(1). \tag{G.6}$$

Proof. The first item immediately follows Lemma 10. In the following, we will prove the second item for the heavy and light tailed cases separately.

For $\alpha \in [0, 0.5)$: We define a sequence of measures $\mu_d\{j/r\} \propto j^{-2\alpha}, j \leq [r]$. We observe that

- We have $\mu_d \to \mu$ weakly such that μ is supported on [0,1] and $\mu([0,\tau]) = \tau^{1-2\alpha}$ for $\tau \in [0,1]$.
- Moreover, (G.6) = $\mathbb{E}_{X \sim \mu_d} \left[(\mathsf{K}_d(X^\alpha/t) \mathbb{1}\{X^\alpha > t\}) \mathbb{1}\{X \leq \frac{r_u}{r}\} \right]$.

By using the C_{δ} definition in Lemma 10:

$$\begin{split} & \left| \mathbb{E}_{X \sim \mu_d} \left[(\mathsf{K}_d(X^\alpha/t) - \mathbb{1}\{X^\alpha > t\}) \mathbb{1}\{X \leq \frac{r_u}{r}\}\} \right] \right| \\ & \leq \mathbb{E}_{X \sim \mu_d} \left[|\mathsf{K}_d(X^\alpha/t) - \mathbb{1}\{X^\alpha > t\}| \mathbb{1}\{X^\alpha \in [0, 1] \setminus \mathcal{C}_\delta\} \right] \\ & + \mathbb{E}_{X \sim \mu_d} \left[|\mathsf{K}_d(X^\alpha/t) - \mathbb{1}\{X^\alpha > t\}| \mathbb{1}\{X^\alpha \in \mathcal{C}_\delta\} \right] \\ & \leq o_d(1) + \mathbb{P}_{X \sim \mu} [X^\alpha \in \mathcal{C}_\delta], \end{split}$$

where we used the second item in Lemma 10 for (a). Since $\mathbb{P}_{X \sim \mu}[X^{\alpha} \in C_{\delta}] \xrightarrow{\delta \to 0} 0$, we have the first result.

For $\alpha > 0.5$: We define a sequence of measures $\mu_d\{j\} \propto j^{-2\alpha}, j \leq [r]$. We observe that

- We have $\mu_d \to \mu$ weakly such that $\mu\{j\} \propto j^{-2\alpha}$ for $j \in \mathbb{N}$.
- Moreover, $(G.6) = \mathbb{E}_{X \sim \mu_d} \left[(\mathsf{K}_d(X^\alpha/t) \mathbb{1}\{X^\alpha > t\}) \mathbb{1}\{X \le r_u\} \right].$

Let $t \in ((j-1)^{\alpha}, j^{\alpha})$ for some $j \in \mathbb{N}$. For small enough $\delta > 0$, we have

$$\begin{split} & \left| \mathbb{E}_{X \sim \mu_d} \big[(\mathsf{K}_d(X^\alpha/t) - \mathbbm{1}\{X^\alpha > t\}) \mathbbm{1}\{X \leq r_u\} \big] \right| \\ & = \mathbb{E}_{X \sim \mu_d} [|\mathsf{K}_d(X^\alpha/t) - \mathbbm{1}\{X^\alpha > t\}| \mathbbm{1}\{X \in [0, r_u]\} \mathbbm{1}\{X^\alpha \not\in \mathcal{C}_\delta\}] \stackrel{(b)}{=} o_d(1), \end{split}$$

where we used both items in Lemma 10 for (b).

Corollary 8. For $1 \ge c_d \gg \log^{-5} d$, we define

$$g_d(\lambda, t) := \frac{-\lambda \exp(-t\lambda)}{1 - \exp(-t\lambda)} + \frac{\lambda^2 \exp(-t\lambda)}{(1 - \exp(-t\lambda))^2} \left(\frac{c_d}{t} \frac{r_s}{d} + \frac{\lambda \exp(-t\lambda)}{1 - \exp(-t\lambda)}\right)^{-1}$$

Let $r_u \leq r$ and

$$\kappa_{\text{eff}} \coloneqq \begin{cases} r^{\alpha}, & \alpha \in [0, 0.5) \\ 1, & \alpha > 0.5 \end{cases}, \qquad \mathsf{T}_{\text{eff}} \coloneqq \kappa_{\text{eff}} \log d / r_s.$$

We have

$$\frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \left(\sum_{j=1}^{r_u} g_d^2(\lambda_j; t\mathsf{T}_{\mathrm{eff}}) - \sum_{j=1}^{r_u} \lambda_j^2 \mathbb{1}\left\{\frac{1}{\lambda_j} > t\kappa_{\mathrm{eff}}\right\} \right) = o_d(1)$$

for any fixed

$$t \in \begin{cases} (0, \infty), & \alpha \in [0, 0.5) \\ (0, \infty) \setminus \{j^{\alpha} : j \in \mathbb{N}\}, & \alpha > 0.5. \end{cases}$$

Proof. We observe that

$$g_d(\lambda;t) = \lambda \left(1 - \frac{1}{1 - \exp(-t\lambda) + \frac{d}{r_s} \frac{\lambda t}{c_d} \exp(-t\lambda)} \right)$$

Therefore, we have $g_d^2(\lambda; t\mathsf{T}_{\mathrm{eff}}) = \lambda^2 F_d(\frac{1}{\lambda_j t \kappa_{\mathrm{eff}}})$. Then, by Proposition 36

$$\begin{split} \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \left(\sum_{j=1}^{r_u} g_d^2(\lambda_j; t\mathsf{T}_{\mathrm{eff}}) - \sum_{j=1}^{r_u} \lambda_j^2 \mathbb{1} \left\{ \frac{1}{\lambda_j} \ge t\kappa_{\mathrm{eff}} \right\} \right) \\ &= \frac{1}{\|\mathbf{\Lambda}\|_{\mathrm{F}}^2} \sum_{j=1}^{r_u} \lambda_j^2 \left(F_d(\frac{1}{\lambda_j t \kappa_{\mathrm{eff}}}) - \mathbb{1} \left\{ \frac{1}{\lambda_j} > t\kappa_{\mathrm{eff}} \right\} \right) = o_d(1). \end{split}$$

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: : The main claims made in the abstract and introduction are supported by our main theorems, Theorem 1 and 2, and their corollaries, Corollary 1 and 2.

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: Yes, we discuss limitations in the Conclusion part (see Section 6)

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: The details, including assumptions, for the theoretical setup are detailed in Section 2. The proofs are presented in appendix (see the Supplementary file)

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 crossreferenced.
- 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: Details to reproduce the simulations in Figure 1b are provided in its caption. 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: [NA]

Justification: This is a theory paper; toy experiments are conducted on Gaussian data.

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: The details of our toy experiments are provided in the caption of Figure 1.

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: Error bars are provided in Figure 1b; however, they are not visible.

- 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: [NA]

Justification: This is a theory paper; toy experiments are conducted on Gaussian data.

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: Yes, this submission follows the NeurIPS Code of Ethics

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: [NA]

: [NA]

Justification: The goal of this paper is to advance the theoretical understanding of training dynamics of two-layer neural networks. There are no direct societal impacts of our work that should be highlighted here

- 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: This paper poses no such risks.

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: This paper does not use existing assets.

- 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: This paper does not release new assets.

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: This paper does not involve human subjects.

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: This paper does not involve human subjects.

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: We do not use LLMs.

- 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.