

# Feature Learning in High-Dimensions under Structured Covariance: Scaling Laws in Quadratic Networks

**Qingyuan Yu\***

*University of Toronto and Vector Institute, Toronto, Canada*

QINGYUANYU@CS.TORONTO.EDU

**Nuri Mert Vural\***

*University of Toronto and Vector Institute, Toronto, Canada*

VURAL@CS.TORONTO.EDU

**Xin T. Tong**

*National University of Singapore, Singapore*

XIN.T.TONG@NUS.EDU.SG

**Murat A. Erdogdu**

*University of Toronto and Vector Institute, Toronto, Canada*

ERDOGDU@CS.TORONTO.COM

## Abstract

Recent theoretical work has shown that nonlinear solvable models exhibit scaling laws in the feature-learning regime. However, these results largely rely on the assumption of isotropic inputs, and understanding how these laws extend to anisotropic data remains a central open problem. In this work, we address this gap by analyzing the learning dynamics of two-layer neural networks with quadratic activations under anisotropic Gaussian inputs. We provide a sharp characterization of online stochastic gradient descent (SGD), explicitly quantifying how the covariance spectrum influences both the scaling exponent and sample complexity. Furthermore, we establish that normalization techniques overcome the intrinsic limitations of vanilla SGD, strictly improving sample efficiency in anisotropic settings. Experiments on two-layer networks with general activation functions support our theoretical predictions, suggesting that these insights extend well beyond the quadratic model.

**Keywords:** scaling laws, stochastic gradient descent, shallow neural network, multi-index model

## 1. Introduction

The predictable scaling of generalization error with respect to computational resources, dataset size, and model width—commonly referred to as scaling laws—has become a defining feature of modern machine learning [10, 11]. Because empirical losses consistently exhibit robust power-law decay, practitioners can reliably extrapolate large-scale performance from a relatively small number of measurements [2, 12]. This phenomenon has sparked significant theoretical interest in deriving scaling laws from first principles [18, 21]. While earlier work primarily focused on linear models [16, 22], recent efforts have shifted toward nonlinear solvable models, which are necessary to capture phenomena such as feature learning [7] and the heavy spectral tails observed in real-world representations [13].

Despite this progress, existing scaling analyses for nonlinear networks largely rely on the restrictive assumption of isotropic inputs [3, 23]. In this work, we study the anisotropic setting. Specifically, we analyze the learning dynamics of a two-layer neural network with quadratic activations under anisotropic Gaussian data:

$$y = \sum_{j=1}^r \lambda_j (\langle \boldsymbol{\theta}_j, \mathbf{x} \rangle^2 - \langle \boldsymbol{\theta}_j, \boldsymbol{\Sigma} \boldsymbol{\theta}_j \rangle) + \epsilon, \quad \mathbf{x} \sim \mathcal{N}(0, \boldsymbol{\Sigma}). \quad (1.1)$$

Here,  $\epsilon$  is zero-mean independent sub-Gaussian noise,  $\{\boldsymbol{\theta}_j\}_{j=1}^r \subset \mathbb{R}^d$  are unknown signal directions,  $\lambda_1 > \lambda_2 > \dots > \lambda_r > 0$  denote their associated signal strengths, and  $\boldsymbol{\Sigma}$  is the data covariance

---

\*. Equal contribution.

matrix, which may exhibit a power-law spectrum. The constant terms  $\langle \theta_j, \Sigma \theta_j \rangle$  are included to ensure the labels remain centered, i.e.,  $\mathbb{E}[y] = 0$ . Our goal is to learn this teacher network using a two-layer student network of width  $r_s$  trained via gradient-based optimization.

A central challenge in this setting is understanding the fundamental statistical and computational limits of feature learning under anisotropy. Existing work provides sample complexity guarantees bounded by the effective dimension  $r_{\text{eff}}(\Sigma) := \text{tr}(\Sigma)/\|\Sigma\|_2$ : scaling quadratically as  $(r_{\text{eff}}(\Sigma))^2$  for kernel methods [8, 24], and scaling linearly as  $r_{\text{eff}}(\Sigma)$  for infinite-width mean-field models [20]. However, the exact risk trajectory of stochastic gradient descent (SGD) applied to narrow networks—and the corresponding scaling behavior—remains poorly understood. More broadly, theoretical analyses of nonlinear multi-index models under anisotropic covariates remain limited. Motivated by these gaps, we investigate the following question:

*Can we characterize the exact risk trajectory and the resulting computational and statistical requirements of SGD for quadratic networks under anisotropic data?*

### 1.1. Our contributions

In this paper, we study the model in (1.1) under structured power-law conditions on both the signal strength and the covariance spectrum. Specifically, letting  $\sigma_i$  denote the eigenvalues of  $\Sigma$ , we assume

$$\lambda_i \asymp i^{-\alpha} \quad \text{and} \quad \sigma_i \asymp i^{-\beta}, \quad (1.2)$$

for exponents  $\alpha, \beta > 0$ . In the existing literature,  $\alpha$  and  $\beta$  is closely related to the *source* and *capacity* conditions, respectively [5, 24]. Our main contributions are summarized below:

1. **Scaling exponents:** In Section 3, we derive a sharp characterization of the excess risk under online vanilla SGD dynamics. In particular, we show that the risk exhibits an explicit power-law decay jointly governed by the source and capacity exponents  $\alpha$  and  $\beta$ . We validate these predictions empirically, observing a strong agreement between theory and simulations (Figure 1).
2. **Algorithmic lower bound:** In Section 3.1, we study the limitations of vanilla online SGD under anisotropic data. Combining theoretical results with supporting empirical evidence, we show that unconstrained SGD is restricted to conservative learning rates, which force its sample complexity to  $\Omega(d)$ . This prevents vanilla SGD from attaining the near-optimal sample complexity  $r_{\text{eff}}(\Sigma)$ .
3. **Near-optimal sample complexity:** In Section 3.2, we propose and analyze an online Stiefel SGD algorithm. We prove that it achieves a near-optimal sample complexity of  $\tilde{\mathcal{O}}(r_{\text{eff}}(\Sigma))$ , matching the effective dimension of the data up to logarithmic factors while preserving the exact scaling behavior. This establishes a sharp separation in sample efficiency between SGD learning in feature-learning regime and lazy training regimes under anisotropic data.

## 2. Problem Setting

We consider learning a two-layer teacher network with a quadratic activation when the input is anisotropic Gaussian given in (1.1), under the following assumption:

**Assumption 1 (Spectral alignment)** *Each  $\theta_j$  is the  $j$ -th eigenvector of  $\Sigma$ , such that  $\Sigma \theta_j = \sigma_j \theta_j$ . For normalization, we assume  $\sum_{i=1}^r \sigma_i^2 \lambda_i^2 = 1$  so that the variance of labels, i.e.,  $\mathbb{E}[y^2]$ , is constant.*

This assumption formalizes a classical statistical learning principle: the directions most predictive of the response  $y$  coincide with the principal components of the input covariates  $\mathbf{x}$  [9]. Combined with the power-law assumptions in (1.2), our setting extends classical three-parameter models to the multi-index setting through a shared spectral decay of the covariance and signal components [17].

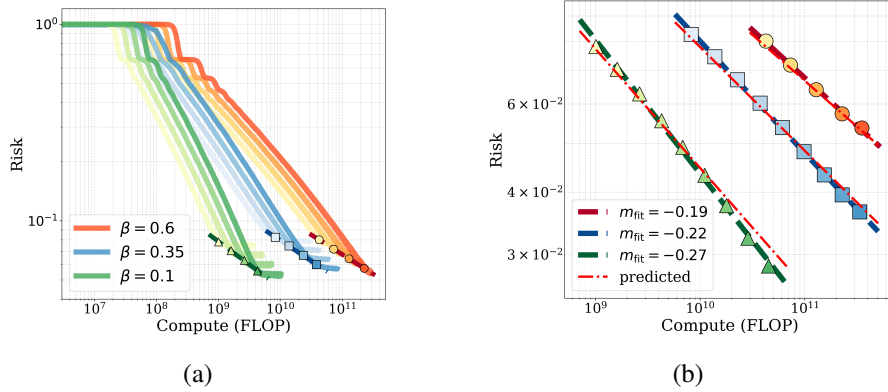


Figure 1: (a) Population risk vs. compute for two-layer quadratic networks trained via one-pass SGD for  $\beta \in \{0.1, 0.35, 0.6\}$  and  $r_s \in \{128, 152, 181, 215\}$  ( $d = 2048$ ,  $\alpha = 0.7 - \beta$ ). (b) Empirical scaling exponents across model widths, with theoretical predictions (red lines) closely matching observed slopes.

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

$$\hat{y}(\mathbf{x}; \mathbf{W}, b) = \frac{1}{r_s} \sum_{j=1}^{r_s} \langle \mathbf{w}_j, \mathbf{x} \rangle^2 - b,$$

where  $r_s$  is the width of the student network,  $b$  is the bias term, and  $\{\mathbf{w}_j\}_{j=1}^{r_s} \subset \mathbb{R}^d$  denotes the set of trainable weights. We collect these weights as the columns of the matrix  $\mathbf{W} \in \mathbb{R}^{d \times r_s}$ .

We consider minimizing the squared loss, for which the population risk is defined by

$$R(\mathbf{W}, b) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}, y)} [(y - \hat{y}(\mathbf{x}, \mathbf{W}, b))^2].$$

By substituting the teacher and student models in (1.1) and , the population risk admits the following closed-form expression:

$$R(\mathbf{W}, b) = \left\| \frac{1}{r_s} \Sigma^{\frac{1}{2}} \mathbf{W} \mathbf{W}^\top \Sigma^{\frac{1}{2}} - \Sigma^{\frac{1}{2}} \Theta \Lambda \Theta^\top \Sigma^{\frac{1}{2}} \right\|_F^2 + \frac{1}{2} \left( b - \frac{1}{r_s} \text{tr}(\Sigma \mathbf{W} \mathbf{W}^\top) \right)^2 + \frac{1}{2} \mathbb{E}[\epsilon^2]. \quad (2.1)$$

Here, we collect the teacher directions  $\{\boldsymbol{\theta}_j\}_{j=1}^r$  as the columns of the matrix  $\Theta \in \mathbb{R}^{d \times r}$ , and  $\Lambda$  is a diagonal matrix where the  $j$ -th diagonal entry is  $\lambda_j$ .

### 3. Main Result: Risk characterization under structured covariance

#### 3.1. Online vanilla SGD

We first analyze the learning dynamics of online (one-pass) SGD. Given a fresh sample  $(\mathbf{x}_{t+1}, y_{t+1})$ , we define the instantaneous loss:

$$\mathcal{L}_{t+1}(\mathbf{W}, b) := \frac{1}{2} (y_{t+1} - \hat{y}(\mathbf{x}_{t+1}; \mathbf{W}, b))^2,$$

where  $\hat{y}(\mathbf{x}_{t+1}; \mathbf{W}, b)$  denotes the network output. To learn the weights and bias simultaneously, we employ two-time-scale gradient descent, where the bias updates significantly faster than the weights. This separation of timescales yields a mathematically tractable learning trajectory and allows the algorithm to learn the covariance-dependent bias term in (2.1) that keeps the predictions centered. Specifically, the parameter updates are given by:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla \mathcal{L}_{t+1}(\mathbf{W}_t, b_t), \quad b_{t+1} = b_t - \frac{\eta}{\delta} \nabla \mathcal{L}_{t+1}(\mathbf{W}_t, b_t) \quad (\text{SGD})$$

where the parameters are initialized with  $b_0 \sim \mathcal{N}(0, 1)$  and  $\mathbf{W}_0 \in \mathbb{R}^{d \times r_s}$  where  $\mathbf{W}_{0,ij} \sim \text{i.i.d. } \mathcal{N}(0, 1/d)$ . The following theorem characterizes the resulting risk trajectory.

**Theorem 1** Let  $\alpha, \beta > 0$  denote the source and capacity exponents defined in (1.2), and assume the teacher width  $r \gg 1$ . We consider the regime  $2\alpha + 2\beta > 1$  and a mildly overparameterized student network with  $r_s \gg \text{polylog } d$ . We consider step size  $\eta$  and timescale parameter  $\delta$  such that  $\eta \asymp o(1/d)$  and  $\delta \in o(1/\log d)$ . Almost surely, we have

$$\mathcal{R}(\eta t \log d) - \frac{1}{2} \mathbb{E}[\epsilon^2] \asymp \Theta\left((\eta t)^{-\frac{2(\alpha+\beta)-1}{\alpha+2\beta}} + r_s^{-2(\alpha+\beta)+1}\right) \quad (3.1)$$

whenever the limit of  $\eta t$  is well-defined.

**Remark 2** We highlight the following:

- This asymptotic risk exhibits a decomposition similar to the neural scaling laws in [11, 12], where the error cleanly separates into optimization and approximation terms. Here, the capacity exponent  $\beta$  enters both contributions additively due to the alignment structure. While the limit  $\beta \rightarrow 0$  recovers the known isotropic scaling behavior [3], our result establishes strictly faster risk decay for any level of anisotropy ( $\beta > 0$ ).
- However, the resulting sample complexity satisfies  $T \simeq d \text{ polylog}(d)$ . Because the algorithm operates in a strict one-pass regime, runtime and sample complexity are identical. This rate is highly suboptimal: the intrinsic statistical complexity is determined by the effective dimension  $r_{\text{eff}}(\Sigma) := \text{tr}(\Sigma)/\|\Sigma\|_2$ , which scales as  $d^{1-\beta}$  for  $\beta \in (0, 1)$  and as  $\mathcal{O}(1)$  for  $\beta > 1$ .

This suboptimal rate suggests a potential algorithmic gap across all anisotropic regimes ( $\beta > 0$ ). Does the  $\tilde{\Theta}(d)$  sample requirement reflect a fundamental limitation of the dynamics in (SGD), or is it merely an artifact of the theoretical analysis? To provide evidence that vanilla SGD inherently requires  $\Omega(d)$  sample complexity, we analyze an offline empirical estimator for the teacher directions:

$$\hat{R}(\mathbf{W}) = \left\| \frac{1}{r_s} \Sigma^{\frac{1}{2}} \mathbf{W} \mathbf{W}^\top \Sigma^{\frac{1}{2}} - \hat{\mathbf{S}} \right\|_{\text{F}}^2, \quad \text{where} \quad \hat{\mathbf{S}} = \frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t \mathbf{x}_t^\top. \quad (3.2)$$

This objective coincides with the population risk in (2.1), up to the bias and the noise variance terms. The key difference is that the population covariance is replaced by its unbiased empirical estimator:  $\mathbb{E}[\hat{\mathbf{S}}] = \Sigma^{\frac{1}{2}} \Theta \Lambda \Theta^\top \Sigma^{\frac{1}{2}}$ . We observe that (SGD) is the online counterpart of minimizing this empirical objective. Since their sample complexities are known to match in the isotropic case ( $\Sigma = I_d$ ), analyzing  $\hat{R}(\mathbf{W})$  provides a reliable proxy for the online dynamics [4]. In the following proposition, we show that minimizing  $\hat{R}(\mathbf{W})$  requires  $\Omega(d)$  samples to achieve a non-vanishing distance with the true teacher directions. Consequently, when  $T \ll d$ , the estimator fundamentally fails to recover the relevant signal. This limitation stems from the  $\Sigma^{1/2}$  factors in (3.2), which implicitly offset the statistical advantages induced by anisotropy in  $\hat{\mathbf{S}}$  and forces the optimal weights  $\hat{\mathbf{W}}_T$  to align with the top eigenspace of the whitened matrix  $\Sigma^{-1/2} \hat{\mathbf{S}} \Sigma^{-1/2}$ .

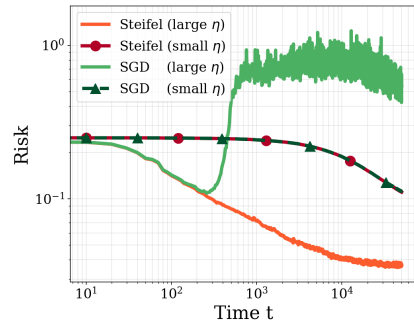


Figure 2: Population risk vs training iteration of (SGD) and Stiefel SGD. Risks are plotted after the alignment step. We set  $d = 16384$ ,  $r = 2048$ ,  $r_s = 512$ , and  $\alpha = \beta = 0.35$ . Small  $\eta$  scales as  $\tilde{\mathcal{O}}(1/d)$ , while large  $\eta$  scales as  $\tilde{\mathcal{O}}(1/r_{\text{eff}}(\Sigma))$ .

**Proposition 3 (Offline Estimator Lower Bound)** *Let  $\hat{\mathbf{W}}_T$  be the minimizer of  $\hat{R}(\mathbf{W})$  in (3.2). If  $T \ll d$ , then for a sufficiently small constant  $\varepsilon > 0$  and  $d$  sufficiently large,*

$$\mathbb{P} \left[ \left\| \frac{1}{r_s} \boldsymbol{\Sigma}^{\frac{1}{2}} \hat{\mathbf{W}}_T \hat{\mathbf{W}}_T^\top \boldsymbol{\Sigma}^{\frac{1}{2}} - \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} \right\|_F > \varepsilon \right] \geq 1 - o_d(1).$$

By showing that the population risk remains lower-bounded when  $T \ll d$ , Proposition 3 provides evidence that the sample complexity requirement in Theorem 1 cannot be substantially improved. Beyond this theoretical argument, we provide empirical observations supporting this conclusion. As illustrated in Figure 2, (SGD) becomes unstable for learning rates  $\eta \asymp \tilde{\mathcal{O}}(1/r_{\text{eff}}(\boldsymbol{\Sigma}))$ , whereas the normalized Stiefel (which is introduced in the subsequent section) remain stable in this regime. Because learning rates of this magnitude are required to achieve the statistically optimal sample complexity in one-pass settings, this instability provides an alternative argument that (SGD) is fundamentally restricted to conservative step sizes.

### 3.2. Improved sample complexity via normalized dynamics

In this section, we introduce normalized Stiefel dynamics to achieve near-optimal sample complexity under anisotropic data. The fundamental issue with the standard quadratic objective (3.2) is that it implicitly counteracts the data's anisotropic structure. Because the network variance scales as  $\mathbb{E}[\hat{y}^2] \propto \|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{W} \mathbf{W}^\top \boldsymbol{\Sigma}^{\frac{1}{2}}\|_F^2$ , unconstrained dynamics are naturally biased toward whitening the covariance spectrum, forcing the weights to align with the inverse covariance geometry.

To bypass this whitening effect and preserve the statistical advantages of anisotropy, we optimize a correlation objective while constraining the weight matrix to the Stiefel manifold  $\text{St}(r_s, d) := \{\mathbf{V} \in \mathbb{R}^{d \times r_s} \mid \mathbf{V}^\top \mathbf{V} = \mathbf{I}_{r_s}\}$ , which makes sure that the weights are unit norm and orthogonal. Given a fresh sample  $(\mathbf{x}_{t+1}, y_{t+1})$ , we define this instantaneous objective and its corresponding Riemannian gradient as:

$$\mathcal{J}_{t+1}(\mathbf{W}, b) = -y_{t+1} \hat{y}(\mathbf{x}_{t+1}; \mathbf{W}, b), \quad \nabla_{\text{St}} \mathcal{J}_{t+1}(\mathbf{W}, b) = (\mathbf{I}_d - \mathbf{W} \mathbf{W}^\top) y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{W}.$$

Following layer-wise training strategies from prior work [1, 6, 15], the training proceeds in two distinct stages:

- Phase 1 (Direction Learning): We initialize  $\mathbf{W}_0 \sim \text{Unif}(\text{St}(r_s, d))$  and  $b_0 = 0$ . For  $t = 1, 2, \dots$ , we update the weights using:

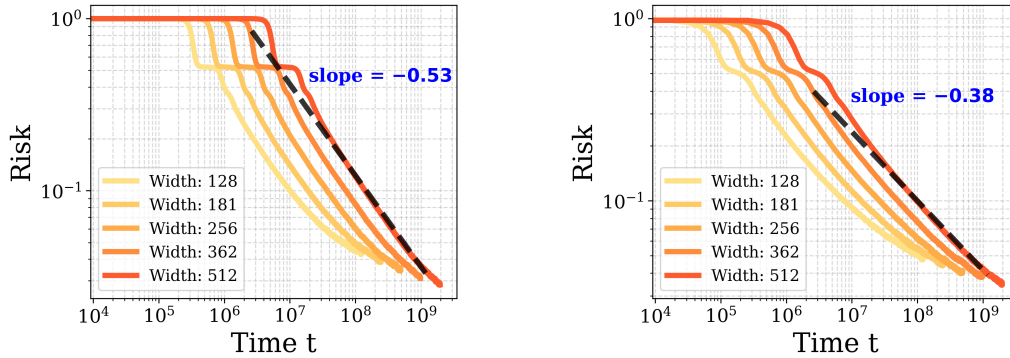
$$\tilde{\mathbf{W}}_{t+1} = \mathbf{W}_t - \eta \nabla_{\text{St}} \mathcal{J}_{t+1}(\mathbf{W}_t, b_0), \quad \mathbf{W}_{t+1} = \tilde{\mathbf{W}}_{t+1} (\tilde{\mathbf{W}}_{t+1}^\top \tilde{\mathbf{W}}_{t+1})^{-1/2}.$$

- Phase 2 (Weight Alignment): Using a set of  $T_{\text{align}}$  fresh samples  $\{(\tilde{\mathbf{x}}_j, \tilde{y}_j)\}$ , we freeze the learned features  $\mathbf{W}_t$  and optimize an alignment matrix  $\boldsymbol{\Omega}$  and scalar bias  $b$ :

$$(\boldsymbol{\Omega}_*, b_*) = \underset{\boldsymbol{\Omega} \in \mathbb{R}^{r_s \times r_s}, b \in \mathbb{R}}{\text{argmin}} \sum_{j=1}^{T_{\text{align}}} (\tilde{y}_j - \hat{y}(\tilde{\mathbf{x}}_j; \mathbf{W}_t \boldsymbol{\Omega}, b))^2 + \lambda \|\boldsymbol{\Omega}\|_F^2.$$

The final parameters for the model are then set to  $\mathbf{W}_t^{\text{final}} := \mathbf{W}_t \boldsymbol{\Omega}_*$  and  $b_t^{\text{final}} := b_*$ .

Let  $\mathcal{R}(t) := R(\mathbf{W}_t^{\text{final}}, b_t^{\text{final}})$  denote the population risk of the final predictor. The following theorem establishes that these normalized dynamics achieve near-optimal sample complexity while matching the exact asymptotic scaling trajectory established in Theorem 1.



(a)  $(\alpha, \beta) = (0.6, 0.1)$ . ideal exponent:  $-0.53$       (b)  $(\alpha, \beta) = (0.35, 0.35)$ . ideal exponent:  $-0.33$

Figure 3: Population loss versus compute for a two-layer squared ReLU network. The model is trained using (SGD) with the learning rate regime in Theorem 1. We use  $d = 2048$  and  $r = 1024$ . The resulting risk trajectories follow a power-law decay, where the empirical exponents (dashed lines) closely align with our theoretical prediction of  $\frac{1-2(\alpha+\beta)}{\alpha+2\beta}$  for the quadratic setting.

**Theorem 4** Consider the teacher model with parameters  $\alpha$ ,  $\beta$ , and  $r$  as defined in Theorem 1. Let the hyperparameters satisfy:  $\eta \ll 1/\tilde{\mathcal{O}}(r_{\text{eff}}(\Sigma))$ , and  $r_s \gg d/\text{polylog}(d)$ , and  $T_{\text{align}} \gg \text{polylog}(d)$ . Then there exists a regularization parameter  $\lambda > 0$ , such that the risk of the aligned predictor  $\mathcal{R}$  satisfies the same guarantee as in (3.1).

We conclude by emphasizing a critical distinction: unlike (SGD), Stiefel dynamics remain stable at significantly larger learning rates. Consequently, the required sample complexity improves to  $\tilde{\mathcal{O}}(r_{\text{eff}}(\Sigma))$ , which matches the effective dimension of the input distribution and is information-theoretically optimal up to logarithmic factors.

## 4. Conclusion

In this work, we characterized the feature learning dynamics of SGD under anisotropic data and specific spectral alignment conditions. We demonstrated that vanilla SGD requires a sample complexity of  $n \asymp T \simeq d \text{polylog}(d)$ , which is intrinsically suboptimal relative to the effective dimension  $r_{\text{eff}}(\Sigma)$  due to instabilities in the empirical objective. To resolve this fundamental  $\Omega(d)$  bottleneck, we introduced a normalized algorithm that optimizes a correlation objective under Stiefel constraints. This approach achieves statistical optimality by maintaining stability at larger learning rates ( $\eta \asymp 1/r_{\text{eff}}(\Sigma)$ ). Future directions include relaxing the spectral alignment assumption and more rigorously formalizing the instability regime of vanilla SGD under anisotropic data.

**Universality of exponents.** A natural next step is extending this theoretical framework beyond two-layer quadratic networks to general activation functions. Because SGD is conjectured to learn in hierarchical phases dictated by Hermite polynomials, standard nonlinearities (e.g., GeLU, tanh) with a nonzero second-order ( $\text{He}_2$ ) component should exhibit an intermediate phase governed by the exact scaling exponents derived in our analysis. Preliminary experiments with squared ReLU activations confirm the existence of this intermediate phase (Figures 3a and 3b), and we leave the formal generalization of this framework as exciting future work.

## References

- [1] 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.
- [2] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27), 2024. ISSN 1091-6490. doi: 10.1073/pnas.2311878121. URL <http://dx.doi.org/10.1073/pnas.2311878121>.
- [3] Gérard Ben Arous, Murat A. Erdogdu, Nuri Mert Vural, and Denny Wu. Learning quadratic neural networks in high dimensions: Sgd dynamics and scaling laws, 2025. URL <https://arxiv.org/abs/2508.03688>.
- [4] Gérard Ben Arous, Cédric Gerbelot, and Vanessa Piccolo. Stochastic gradient descent in high dimensions for multi-spiked tensor pca, 2025. URL <https://arxiv.org/abs/2410.18162>.
- [5] Andrea Caponnetto and Ernesto de Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7:331–368, 2007. URL <https://api.semanticscholar.org/CorpusID:207063850>.
- [6] 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.
- [7] Leonardo Defilippis, Yizhou Xu, Julius Girardin, Emanuele Troiani, Vittorio Erba, Lenka Zdeborová, Bruno Loureiro, and Florent Krzakala. Scaling laws and spectra of shallow neural networks in the feature learning regime, 2025. URL <https://arxiv.org/abs/2509.24882>.
- [8] 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.
- [9] Trevor Hastie, Robert Tibshirani, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [10] 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.
- [11] 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.
- [12] 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.

- [13] Juno Kim, Eshaan Nichani, Denny Wu, Alberto Bietti, and Jason D. Lee. Sharp capacity scaling of spectral optimizers in learning associative memory, 2026. URL <https://arxiv.org/abs/2603.26554>.
- [14] B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/2674095>.
- [15] 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.
- [16] 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.
- [17] Alexander Maloney, Daniel A Roberts, and James Sully. A solvable model of neural scaling laws. *arXiv preprint arXiv:2210.16859*, 2022.
- [18] Eric Michaud, Ziming Liu, Uzay Girit, and Max Tegmark. The quantization model of neural scaling. *Advances in Neural Information Processing Systems*, 36, 2024.
- [19] Marco Mondelli and Andrea Montanari. Fundamental limits of weak recovery with applications to phase retrieval. *Foundations of Computational Mathematics*, 19:703–773, 2017.
- [20] 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.
- [21] 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.
- [22] Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+ 3 phases of compute-optimal neural scaling laws. *arXiv preprint arXiv:2405.15074*, 2024.
- [23] 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.
- [24] Arie Wortsman and Bruno Loureiro. Kernel ridge regression under power-law data: spectrum and generalization. *ArXiv*, abs/2510.04780, 2025. URL <https://api.semanticscholar.org/CorpusID:281842750>.

**Contents**

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our contributions . . . . .	2
<b>2</b>	<b>Problem Setting</b>	<b>2</b>
<b>3</b>	<b>Main Result: Risk characterization under structured covariance</b>	<b>3</b>
3.1	Online vanilla SGD . . . . .	3
3.2	Improved sample complexity via normalized dynamics . . . . .	5
<b>4</b>	<b>Conclusion</b>	<b>6</b>
<b>A</b>	<b>Preliminaries for Proofs</b>	<b>10</b>
<b>B</b>	<b>SGD comparison proof</b>	<b>12</b>
B.1	Setup, exact recursion, and centered ideal discrete flow . . . . .	12
B.2	Scalar logistic facts . . . . .	13
B.3	Initialization, localization, and one-step bounds . . . . .	15
B.4	Discrete bias tracking . . . . .	18
B.5	Discrete Frobenius localization with explicit remainders . . . . .	19
B.6	Rest-block reduction for SGD . . . . .	20
B.7	Fast-block bootstrap for SGD . . . . .	23
<b>C</b>	<b>Proof of Proposition 3</b>	<b>30</b>
<b>D</b>	<b>Normalized SGD Dynamics</b>	<b>30</b>
D.1	Definitions and bounding systems . . . . .	34
D.2	Bounding the second order terms . . . . .	35
D.3	Noise characterization . . . . .	37
D.3.1	Application of noise bounds . . . . .	41
D.4	Alignment step . . . . .	43
<b>E</b>	<b>Some moment bounds and concentration inequalities</b>	<b>43</b>

## Appendix A. Preliminaries for Proofs

**Proof organization.** Section A introduces the additional notation and definitions used throughout, and collects several auxiliary lemma. Section B proves the discrete-time online SGD result in Theorem 1. Section C proves our lower bound argument in Proposition 3. Finally, Section D proves the online Steifel SGD result in Theorem 4.

**Additional notation and definitions.** We use a reduced time parameter  $\bar{t} > 0$  and relate it to the physical time scales by

$$T = \bar{t} T_{\text{eff}}, \quad T_{\text{eff}} := \frac{Z_r}{2} \log \frac{d}{Z_r}, \quad N = \lfloor \bar{t} N_{\text{eff}} \rfloor, \quad N_{\text{eff}} := \frac{Z_r}{2\eta} \log \frac{d}{Z_r}.$$

We work in the common eigenbasis of the covariance and the teacher:

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d), \quad \mathbf{T} = \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0),$$

where

$$\lambda_1 \geq \dots \geq \lambda_r > 0, \quad \sigma_1 \geq \dots \geq \sigma_d > 0, \quad \lambda_i \asymp i^{-\alpha}, \quad \sigma_i \asymp i^{-\beta}. \quad (\text{A.1})$$

We assume

$$2\alpha + 2\beta > 1, \quad \text{equivalently} \quad \alpha + \beta > \frac{1}{2}.$$

In particular, under the usual fixed spectral envelope in (A.1), one has  $\sigma_1 = O(1)$ .

Define

$$Z_r := \|\mathbf{T}\Sigma\|_F = \left( \sum_{i=1}^r \sigma_i^2 \lambda_i^2 \right)^{1/2}, \quad \frac{d}{Z_r} \rightarrow \infty.$$

Set

$$\mathbf{S}^* := \frac{1}{Z_r} \Sigma^{1/2} \mathbf{T} \Sigma^{1/2}, \quad \mathbf{T}_* := \Sigma \mathbf{S}^* = \mathbf{S}^* \Sigma = \frac{1}{Z_r} \Sigma \mathbf{T} \Sigma.$$

Then

$$\mathbf{T}_* = \text{diag}(\tau_1, \dots, \tau_r, 0, \dots, 0), \quad \tau_i = \frac{\sigma_i^2 \lambda_i}{Z_r}, \quad \kappa_i := \sigma_i^2 \lambda_i, \quad s_i^* := \frac{\tau_i}{\sigma_i} = \frac{\sigma_i \lambda_i}{Z_r}.$$

The normalization is

$$\|\mathbf{S}^*\|_F^2 = \sum_{i=1}^r (s_i^*)^2 = 1.$$

And

$$\|\mathbf{T}_*\|_F^2 = \frac{\sum_{i=1}^r \sigma_i^4 \lambda_i^2}{Z_r^2} \leq \sigma_1^2 \|\mathbf{S}^*\|_F^2 = O(1). \quad (\text{A.2})$$

Fix  $\bar{t} > 0$  and define

$$m_*(\bar{t}) := \max\{i \leq r : \bar{t} \kappa_i > 1\},$$

with the convention that the maximum of the empty set is 0. We write

$$k := m_*(\bar{t}).$$

If  $k < r$ , set

$$a_s := \bar{t} \kappa_{k+1} \in [0, 1),$$

and if  $k = r$ , set  $a_s := 0$ .

Fix

$$0 < \rho < \min\{\beta, 1\}.$$

When  $k < r$  and the slow block is nonempty, we additionally assume

$$\rho < a_s.$$

Define

$$m_\rho(\bar{t}) := \max\{i \leq r : \bar{t} \kappa_i > \rho\},$$

again with the convention that the maximum of the empty set is 0.

We assume an off-transition gap: there exists  $\delta_{\text{gap}} \in (0, 1)$ , independent of  $d$ , such that

$$\bar{t} \kappa_i \notin [1 - \delta_{\text{gap}}, 1 + \delta_{\text{gap}}] \quad \text{for all } i \leq m_\rho(\bar{t}).$$

Consequently,

$$k = \max\{i \leq m_\rho(\bar{t}) : \bar{t} \kappa_i \geq 1 + \delta_{\text{gap}}\},$$

and, if  $k < r$ ,

$$a_s \leq 1 - \delta_{\text{gap}}.$$

We use the diagonal projectors

$$\mathbf{P}_f := \text{diag}(\mathbb{1}\{i \leq k\})_{i=1}^d, \quad \mathbf{P}_r := \mathbf{I}_d - \mathbf{P}_f.$$

Inside the rest block, we further write

$$\mathbf{P}_s := \text{diag}(\mathbb{1}\{k < i \leq m_\rho(\bar{t})\})_{i=1}^d, \quad \mathbf{P}_t := \mathbf{I}_d - \mathbf{P}_f - \mathbf{P}_s.$$

Since  $\kappa_i \asymp i^{-(\alpha+2\beta)}$ , for every fixed  $\bar{t}$  and  $\rho$ ,

$$k = O_{\bar{t}}(1), \quad m_\rho(\bar{t}) = O_{\bar{t}, \rho}(1).$$

Unless stated otherwise, constants may depend on the fixed spectral regularity constants, on  $\bar{t}$ ,  $\rho$ ,  $\delta_{\text{gap}}$ ,  $K$ , and  $\sigma_\varepsilon$ , but not on  $d, N, \eta, \delta, r_s$ .

**Lemma 5 (Rank floor)** *For every matrix  $\mathbf{M}$  with  $r_u(\mathbf{M}) \leq r_s$ ,*

$$\|\mathbf{S}^* - \mathbf{M}\|_F^2 \geq \sum_{i > r_s} (s_i^*)^2. \quad (\text{A.3})$$

*If  $\alpha + \beta > 1/2$ , then, in the regime  $r_s \rightarrow \infty$  with  $r_s < r$ ,*

$$\sum_{i > r_s} (s_i^*)^2 \asymp r_s^{-2(\alpha+\beta)+1}. \quad (\text{A.4})$$

**Proof** The nonzero singular values of  $\mathbf{S}^*$  are precisely  $(s_i^*)_{i \leq r}$ , in nonincreasing order. The lower bound (A.3) is the Eckart–Young theorem. The tail estimate (A.4) follows from  $(s_i^*)^2 \asymp i^{-2(\alpha+\beta)}$  and the integral test.  $\blacksquare$

## Appendix B. SGD comparison proof

### B.1. Setup, exact recursion, and centered ideal discrete flow

Let

$$\mathbf{x}_n \sim \mathcal{N}(0, \Sigma) \text{ i.i.d.}, \quad \varepsilon_n \perp \mathbf{x}_n, \quad \mathbb{E}[\varepsilon_n] = 0, \quad \|\varepsilon_n\|_{\psi_2} \leq \sigma_\varepsilon.$$

Set

$$\mathbf{g}_n := \Sigma^{-1/2} \mathbf{x}_n \sim \mathcal{N}(0, \mathbf{I}_d), \quad \boldsymbol{\chi}_n := \mathbf{g}_n \mathbf{g}_n^\top, \quad \boldsymbol{\xi}_n := \boldsymbol{\chi}_n - \mathbf{I}_d.$$

The label is

$$y_n = \left\langle \frac{\mathbf{T}}{Z_r}, \mathbf{x}_n \mathbf{x}_n^\top - \Sigma \right\rangle + \varepsilon_n = \langle \mathbf{S}^*, \boldsymbol{\xi}_n \rangle + \varepsilon_n.$$

Define

$$\mathbf{G}_n := \frac{1}{r_s} \mathbf{W}_n \mathbf{W}_n^\top, \quad \mathbf{M}_n := \Sigma^{1/2} \mathbf{G}_n \Sigma^{1/2}, \quad e_n := b_n - \text{tr}(\mathbf{M}_n), \quad \boldsymbol{\Delta}_n := \mathbf{S}^* - \mathbf{M}_n.$$

The online prediction is

$$\hat{y}_{n+1} = \langle \mathbf{M}_n, \boldsymbol{\chi}_{n+1} \rangle - b_n.$$

To avoid conflict with the two-timescale parameter  $\delta$ , the prediction error is denoted by

$$\mathfrak{t}_{n+1} := y_{n+1} - \hat{y}_{n+1} = \langle \boldsymbol{\Delta}_n, \boldsymbol{\xi}_{n+1} \rangle + e_n + \varepsilon_{n+1}. \quad (\text{B.1})$$

The online SGD updates are

$$\mathbf{W}_{n+1} = \mathbf{W}_n + \frac{\eta}{2} \mathfrak{t}_{n+1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^\top \mathbf{W}_n, \quad (\text{B.2})$$

$$b_{n+1} = b_n - \frac{\eta}{4\delta} \mathfrak{t}_{n+1}. \quad (\text{B.3})$$

Let

$$N_{\text{eff}} := \frac{Z_r}{2\eta} \log \frac{d}{Z_r}, \quad N := \lfloor \bar{t} N_{\text{eff}} \rfloor, \quad \Lambda_d := \log \frac{ed}{Z_r}.$$

Let

$$\mathcal{F}_n := \sigma(\mathbf{W}_0, b_0, (\mathbf{x}_t, \varepsilon_t)_{1 \leq t \leq n}).$$

We also set

$$p_d := 40(K+2) \log(edN), \quad R_s := 1 + \sqrt{r_s}, \quad L_2 := \text{tr}(\Sigma^2).$$

**Proposition 6 (Exact sandwich recursion and first-order expansion)** *For every  $n \geq 0$ ,  $\mathbf{M}_n \succeq 0$  and*

$$\mathbf{M}_{n+1} = \left( \mathbf{I}_d + \frac{\eta}{2} \mathfrak{t}_{n+1} \Sigma \boldsymbol{\chi}_{n+1} \right) \mathbf{M}_n \left( \mathbf{I}_d + \frac{\eta}{2} \mathfrak{t}_{n+1} \boldsymbol{\chi}_{n+1} \Sigma \right). \quad (\text{B.4})$$

Equivalently,

$$\mathbf{M}_{n+1} = \mathbf{M}_n + \eta \mathcal{F}^\sharp(\mathbf{M}_n, e_n) + \eta \mathbf{Z}_{n+1} + \eta^2 \mathbf{Q}_{n+1}, \quad (\text{B.5})$$

where

$$\mathcal{F}^\sharp(\mathbf{M}, e) := \mathcal{F}(\mathbf{M}) + e\mathcal{B}(\mathbf{M}),$$

and

$$\begin{aligned} \mathbf{Z}_{n+1} &:= \frac{1}{2} \mathbf{r}_{n+1} (\boldsymbol{\Sigma} \boldsymbol{\chi}_{n+1} \mathbf{M}_n + \mathbf{M}_n \boldsymbol{\chi}_{n+1} \boldsymbol{\Sigma}) - \mathcal{F}^\sharp(\mathbf{M}_n, e_n), \\ \mathbf{Q}_{n+1} &:= \frac{1}{4} \mathbf{r}_{n+1}^2 \boldsymbol{\Sigma} \boldsymbol{\chi}_{n+1} \mathbf{M}_n \boldsymbol{\chi}_{n+1} \boldsymbol{\Sigma}. \end{aligned}$$

Moreover,

$$\mathbb{E}[\mathbf{Z}_{n+1} \mid \mathcal{F}_n] = 0. \quad (\text{B.6})$$

**Proof** Since  $\mathbf{x}_{n+1} \mathbf{x}_{n+1}^\top = \boldsymbol{\Sigma}^{1/2} \boldsymbol{\chi}_{n+1} \boldsymbol{\Sigma}^{1/2}$ , (B.2) gives (B.4); PSD is therefore preserved. Expanding the sandwich gives (B.5).

It remains to identify the conditional drift. By (B.1),

$$\mathbb{E}[\mathbf{r}_{n+1} \boldsymbol{\chi}_{n+1} \mid \mathcal{F}_n] = \mathbb{E}[\mathbf{r}_{n+1} \boldsymbol{\xi}_{n+1} \mid \mathcal{F}_n] + e_n \mathbf{I}_d.$$

We have

$$\mathbb{E}[\langle \mathbf{A}, \boldsymbol{\xi}_{n+1} \rangle \boldsymbol{\xi}_{n+1}] = 2\mathbf{A} \quad (\mathbf{A} = \mathbf{A}^\top).$$

Taking  $\mathbf{A} = \boldsymbol{\Delta}_n$  yields

$$\mathbb{E}[\mathbf{r}_{n+1} \boldsymbol{\chi}_{n+1} \mid \mathcal{F}_n] = 2\boldsymbol{\Delta}_n + e_n \mathbf{I}_d.$$

Substitution into the linear part of (B.4) gives exactly  $\mathcal{F}^\sharp(\mathbf{M}_n, e_n)$ , proving (B.6).  $\blacksquare$

The centered ideal discrete flow is the diagonal recursion

$$\mathbf{L}_{n+1} = \mathbf{L}_n + \eta (\mathbf{T}_* \mathbf{L}_n + \mathbf{L}_n \mathbf{T}_* - 2\mathbf{L}_n \boldsymbol{\Sigma} \mathbf{L}_n) + \eta^2 (\mathbf{T}_* - \boldsymbol{\Sigma} \mathbf{L}_n) \mathbf{L}_n (\mathbf{T}_* - \mathbf{L}_n \boldsymbol{\Sigma}), \quad \mathbf{L}_0 = \frac{1}{d} \boldsymbol{\Sigma}.$$

Equivalently,

$$\begin{aligned} \ell_{i,n+1} &= \ell_{i,n} (1 + \eta(\tau_i - \sigma_i \ell_{i,n}))^2, & i \leq r, \\ \ell_{i,n+1} &= \ell_{i,n} (1 - \eta \sigma_i \ell_{i,n})^2, & i > r. \end{aligned} \quad (\text{B.7})$$

## B.2. Scalar logistic facts

**Lemma 7 (Basic scalar logistic facts)** Fix  $c \in (0, 1/2]$ , and let

$$a_{n+1} = a_n (1 + c(1 - a_n))^2, \quad 0 \leq a_0 \leq 1. \quad (\text{B.8})$$

Then:

- (i)  $0 \leq a_n \leq 1$  and  $a_{n+1} \geq a_n$  for all  $n$ .
- (ii)  $a_n \leq a_0 (1 + c)^{2n}$  for all  $n$ .
- (iii) If  $a_m \leq \theta < 1$  for all  $0 \leq m \leq n$ , then

$$a_n \geq a_0 (1 + c(1 - \theta))^{2n}.$$

- (iv) If  $a_n \geq \theta > 0$ , then

$$1 - a_{n+1} \leq (1 - 2c\theta)(1 - a_n). \quad (\text{B.9})$$

**Proof** Let  $\phi(a) := a(1 + c(1 - a))^2$ . Then  $\phi([0, 1]) \subseteq [0, 1]$  and

$$\phi(a) - a = ac(1 - a)(2 + c(1 - a)) \geq 0,$$

which proves (i). Item (ii) follows from  $1 + c(1 - a_n) \leq 1 + c$ . Item (iii) follows from  $1 - a_m \geq 1 - \theta$ . Finally,

$$1 - a_{n+1} = 1 - a_n(1 + c(1 - a_n))^2 = (1 - a_n) (1 - 2ca_n - c^2a_n(1 - a_n)).$$

If  $a_n \geq \theta$ , the last factor is at most  $1 - 2c\theta$ , proving (B.9). ■

**Lemma 8 (Robust discrete threshold law)** Fix  $i \leq r$  and set

$$c_i := \eta\tau_i, \quad n(\bar{t}) := \lfloor \bar{t} N_{\text{eff}} \rfloor.$$

Assume

$$\eta\tau_1 \leq \frac{1}{2}, \quad \eta \log^2 \frac{d}{Z_r} \rightarrow 0. \quad (\text{B.10})$$

Let  $a_n$  satisfy (B.8) with  $c = c_i$  and initial data

$$c_- \frac{Z_r}{d} \leq a_0 \leq c_+ \frac{Z_r}{d},$$

where  $0 < c_- < c_+ < \infty$  are independent of  $d$ . For every fixed  $\delta_{\text{thr}} \in (0, 1)$ :

(i) If  $\bar{t} \kappa_i \leq 1 - \delta_{\text{thr}}$ , then  $a_{n(\bar{t})} \rightarrow 0$ .

(ii) If  $\bar{t} \kappa_i \geq 1 + \delta_{\text{thr}}$ , then  $a_{n(\bar{t})} \rightarrow 1$ .

**Proof** Write  $n = n(\bar{t})$ . If  $\bar{t} \kappa_i \leq 1 - \delta_{\text{thr}}$ , Lemma 7(ii) gives

$$a_n \leq c_+ \frac{Z_r}{d} \exp(2nc_i).$$

Since

$$2nc_i = \bar{t} \kappa_i \log \frac{d}{Z_r} + o(1),$$

we have

$$a_n \leq c_+ \left( \frac{d}{Z_r} \right)^{-1 + \bar{t} \kappa_i + o(1)} \leq c_+ \left( \frac{d}{Z_r} \right)^{-\delta_{\text{thr}} + o(1)} \rightarrow 0.$$

Now assume  $\bar{t} \kappa_i \geq 1 + \delta_{\text{thr}}$ . Set

$$\varepsilon := \frac{\delta_{\text{thr}}}{4}, \quad \gamma := 1 - \frac{\delta_{\text{thr}}}{2(1 + \delta_{\text{thr}})}, \quad n_- := \lfloor \gamma \bar{t} N_{\text{eff}} \rfloor.$$

Then  $(1 - \varepsilon)\gamma(1 + \delta_{\text{thr}}) > 1$ . If  $a_m \leq \varepsilon$  for all  $m \leq n_-$ , Lemma 7(iii) yields

$$a_{n_-} \geq a_0(1 + c_i(1 - \varepsilon))^{2n_-}.$$

Using  $\log(1 + x) \geq x - x^2$  for  $x \in [0, 1/2]$  and (B.10),

$$\log a_{n_-} \geq \log a_0 + ((1 - \varepsilon)\gamma \bar{t} \kappa_i + o(1)) \log \frac{d}{Z_r}.$$

The exponent is strictly larger than 1 by a fixed margin, contradicting  $a_{n_-} \leq 1$ . Hence some  $m \leq n_-$  has  $a_m \geq \varepsilon$ , and monotonicity gives  $a_{n_-} \geq \varepsilon$ . Lemma 7(iv) implies

$$1 - a_n \leq \exp(-2c_i \varepsilon (n - n_-)).$$

Finally,

$$2c_i(n - n_-) = (1 - \gamma)\bar{t}\kappa_i \log \frac{d}{Z_r} + o(1) \rightarrow \infty,$$

so  $a_n \rightarrow 1$ . ■

### B.3. Initialization, localization, and one-step bounds

**Lemma 9 (Initialization event for SGD)** *Assume*

$$(\mathbf{W}_0)_{ia} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/d), \quad b_0 = 0.$$

Then there exists an event  $\mathcal{G}_{\text{init}}$  such that

$$\mathbb{P}(\mathcal{G}_{\text{init}}) \geq 1 - Cd^{-K},$$

and on  $\mathcal{G}_{\text{init}}$ ,

$$\|\mathbf{M}_0\|_F \leq 2, \quad |e_0| \leq 2.$$

Moreover, for every deterministic diagonal projector  $\mathbf{P}$ ,

$$\text{tr}(\mathbf{P}\mathbf{M}_0) \leq \frac{C}{d} \text{tr}(\mathbf{P}\boldsymbol{\Sigma}). \quad (\text{B.11})$$

If  $k \geq 1$ , then for every  $1 \leq i \leq k$ ,

$$\left| (\mathbf{H}_0)_{ii} - \frac{\sigma_i}{d} \right| \leq \frac{C\sigma_i}{d} \sqrt{\frac{pd}{r_s}}, \quad (\mathbf{H}_0)_{ii} \geq \frac{c}{d}, \quad (\text{B.12})$$

and

$$\delta_0^{\text{coh}} := \max_{1 \leq i \neq j \leq k} \frac{|(\mathbf{H}_0)_{ij}|}{\sqrt{(\mathbf{H}_0)_{ii}(\mathbf{H}_0)_{jj}}} \leq C \sqrt{\frac{pd}{r_s}}. \quad (\text{B.13})$$

Define the localization stopping time

$$\tau_{\Delta} := \inf\{n \geq 0 : \|\boldsymbol{\Delta}_n\|_F > 4\},$$

and the one-step localized event

$$\mathcal{G}_{\text{loc},n} := \{\|\boldsymbol{\Delta}_n\|_F \leq 4, |e_n| \leq 4, \|\mathbf{M}_n\|_F \leq 5\}.$$

**Lemma 10 (Gaussian chaos bounds)** *For every  $p \geq 2$ , the following hold.*

(i) *For every symmetric  $\mathbf{A}$ ,*

$$\mathbb{E}[\langle \mathbf{A}, \boldsymbol{\xi}_n \rangle^2] = 2\|\mathbf{A}\|_F^2, \quad \|\langle \mathbf{A}, \boldsymbol{\xi}_n \rangle\|_{L^p} \leq Cp\|\mathbf{A}\|_F.$$

(ii) For symmetric  $\mathbf{A}, \mathbf{B}$ ,

$$\|\langle \mathbf{A}, \boldsymbol{\xi}_n \rangle \langle \mathbf{B}, \boldsymbol{\xi}_n \rangle\|_{L^p} \leq Cp^2 \|\mathbf{A}\|_F \|\mathbf{B}\|_F.$$

(iii) If  $\mathbf{A} \succeq 0$ , then

$$\|\mathbf{g}_n^\top \mathbf{A} \mathbf{g}_n\|_{L^p} \leq C (\text{tr}(\mathbf{A}) + p \|\mathbf{A}\|_{\text{op}}). \quad (\text{B.14})$$

(iv) If  $\mathbf{M} \succeq 0$ ,  $r_u(\mathbf{M}) \leq r_s$ , and  $\|\mathbf{M}\|_F \leq 5$ , then for every diagonal projector  $\mathbf{P}$ ,

$$\|\text{tr}(\mathbf{P} \boldsymbol{\Sigma} \boldsymbol{\chi}_n \mathbf{M} \boldsymbol{\chi}_n \boldsymbol{\Sigma})\|_{L^p} \leq Cp^2 R_s (1 + \text{tr}(\mathbf{P} \boldsymbol{\Sigma}^2)). \quad (\text{B.15})$$

**Proof** Items (i) and (ii) are standard hypercontractive estimates for Gaussian chaoses of order two and four. Item (iii) follows by diagonalization and moment bounds for weighted sums of centered  $\chi^2$  variables. For (iv),

$$\text{tr}(\mathbf{P} \boldsymbol{\Sigma} \boldsymbol{\chi}_n \mathbf{M} \boldsymbol{\chi}_n \boldsymbol{\Sigma}) = (\mathbf{g}_n^\top \mathbf{M} \mathbf{g}_n) (\mathbf{g}_n^\top \mathbf{P} \boldsymbol{\Sigma}^2 \mathbf{P} \mathbf{g}_n).$$

Apply Hölder and (B.14). Since  $\mathbf{M} \succeq 0$ ,

$$\text{tr}(\mathbf{M}) \leq \sqrt{r_u(\mathbf{M})} \|\mathbf{M}\|_F \leq CR_s.$$

This proves (B.15). ■

**Lemma 11 (Martingale maximal and stable convolution bounds)** *Let  $(X_n)_{n \geq 1}$  be a martingale difference sequence with respect to  $(\mathcal{G}_n)_{n \geq 0}$ . Then, for  $p \geq 2$ ,*

$$\left\| \sup_{1 \leq m \leq N} \left\| \sum_{n=1}^m X_n \right\|_{L^p} \right\|_{L^p} \leq C \left[ \sqrt{p} \left\| \left( \sum_{n=1}^N \mathbb{E}[X_n^2 | \mathcal{G}_{n-1}] \right)^{1/2} \right\|_{L^p} + p \left( \sum_{n=1}^N \|X_n\|_{L^p}^p \right)^{1/p} \right]. \quad (\text{B.16})$$

If  $0 < \gamma < 1$  and

$$S_n := \sum_{m=0}^{n-1} \gamma^{n-1-m} X_{m+1},$$

with

$$\mathbb{E}[X_{m+1}^2 | \mathcal{F}_m] \leq v^2, \quad \|X_{m+1}\|_{L^{p_d}(\mathbb{P}(\cdot | \mathcal{F}_m))} \leq A_{p_d},$$

then, with probability at least  $1 - Cd^{-K}$ ,

$$\sup_{0 \leq n \leq N} |S_n| \leq C \left( \sqrt{\frac{p_d v^2}{1 - \gamma}} + p_d A_{p_d} \right).$$

**Proof** The first estimate is the Burkholder–Freedman maximal inequality in  $L^p$ . For the stable convolution, apply (B.16) to the weighted martingale for fixed  $n$ , using

$$\sum_{j \geq 0} \gamma^{2j} \leq \frac{1}{1 - \gamma}, \quad \left( \sum_{j \geq 0} \gamma^{p_d j} \right)^{1/p_d} \leq C.$$

A union bound over  $0 \leq n \leq N$  is absorbed into the definition of  $p_d$ . ■

For a diagonal projector  $\mathbf{P}$ , define

$$u_n^{(\mathbf{P})} := \text{tr}(\mathbf{P}\mathbf{M}_n), \quad \tau(\mathbf{P}) := \max_{i:P_{ii}=1} \tau_i, \quad \sigma(\mathbf{P}) := \max_{i:P_{ii}=1} \sigma_i,$$

with empty maxima equal to 0. Also set

$$d_{n+1}^{(\mathbf{P})} := \text{tr}(\mathbf{P}\mathbf{Z}_{n+1}), \quad q_{n+1}^{(\mathbf{P})} := \text{tr}(\mathbf{P}\mathbf{Q}_{n+1}),$$

and

$$\bar{q}_n^{(\mathbf{P})} := \mathbb{E}[q_{n+1}^{(\mathbf{P})} \mid \mathcal{F}_n], \quad \hat{q}_{n+1}^{(\mathbf{P})} := q_{n+1}^{(\mathbf{P})} - \bar{q}_n^{(\mathbf{P})}.$$

**Lemma 12 (Localized one-step bounds)** *On  $\mathcal{G}_{\text{loc},n}$ , the following estimates hold for every  $p \geq 2$ .*

(i) *For every diagonal projector  $\mathbf{P}$ ,*

$$\mathbb{E}[d_{n+1}^{(\mathbf{P})} \mid \mathcal{F}_n] = 0, \quad \mathbb{E}[\hat{q}_{n+1}^{(\mathbf{P})} \mid \mathcal{F}_n] = 0,$$

and

$$\|d_{n+1}^{(\mathbf{P})}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq Cp^2 (1 + \text{tr}(\mathbf{P}\Sigma^2))^{1/2}, \quad (\text{B.17})$$

$$\bar{q}_n^{(\mathbf{P})} \leq CR_s (1 + \text{tr}(\mathbf{P}\Sigma^2)), \quad (\text{B.18})$$

$$\|\hat{q}_{n+1}^{(\mathbf{P})}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq CR_s p^4 (1 + \text{tr}(\mathbf{P}\Sigma^2)). \quad (\text{B.19})$$

(ii) *If  $\mathbf{P}$  has bounded support contained in  $\{1, \dots, r\}$ , then*

$$\|d_{n+1}^{(\mathbf{P})}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq Cp^2 (u_n^{(\mathbf{P})})^{1/2}, \quad (\text{B.20})$$

$$\bar{q}_n^{(\mathbf{P})} \leq CR_s, \quad \|\hat{q}_{n+1}^{(\mathbf{P})}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq CR_s p^4.$$

(iii) *Define*

$$u_{n+1} := \text{tr}(\Sigma\mathbf{X}_{n+1}\mathbf{M}_n), \quad s_n := \text{tr}(\Sigma\mathbf{M}_n),$$

and

$$\zeta_{n+1}^e := -\frac{1}{4\delta}(\mathbf{r}_{n+1} - e_n) - (\mathbf{r}_{n+1}u_{n+1} - \mathbb{E}[\mathbf{r}_{n+1}u_{n+1} \mid \mathcal{F}_n]). \quad (\text{B.21})$$

Then

$$\mathbb{E}[\zeta_{n+1}^e \mid \mathcal{F}_n] = 0, \quad \|\zeta_{n+1}^e\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq C \left( \frac{p}{\delta} + p^2 \right). \quad (\text{B.22})$$

For

$$q_{n+1}^e := q_{n+1}^{(\mathbf{I}_d)} = \text{tr}(\mathbf{Q}_{n+1}),$$

one has

$$\mathbb{E}[q_{n+1}^e \mid \mathcal{F}_n] \leq CR_s(1 + L_2), \quad \|q_{n+1}^e - \mathbb{E}[q_{n+1}^e \mid \mathcal{F}_n]\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq CR_s p^4(1 + L_2). \quad (\text{B.23})$$

**Proof** For (i), write

$$d_{n+1}^{(\mathbf{P})} = \mathbf{r}_{n+1} \langle \mathbf{K}_n^{(\mathbf{P})}, \mathbf{X}_{n+1} \rangle - \mathbb{E}[\mathbf{r}_{n+1} \langle \mathbf{K}_n^{(\mathbf{P})}, \mathbf{X}_{n+1} \rangle \mid \mathcal{F}_n],$$

where

$$\mathbf{K}_n^{(\mathbf{P})} := \frac{1}{2}(\mathbf{M}_n \mathbf{P} \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{P} \mathbf{M}_n).$$

On  $\mathcal{G}_{\text{loc},n}$ ,

$$\|\mathbf{K}_n^{(\mathbf{P})}\|_F \leq C (1 + \text{tr}(\mathbf{P} \boldsymbol{\Sigma}^2))^{1/2}.$$

Using (B.1),  $\|\boldsymbol{\Delta}_n\|_F \leq 4$ ,  $|e_n| \leq 4$ , and Lemma 10, we get (B.17). The second-order term satisfies

$$q_{n+1}^{(\mathbf{P})} = \frac{1}{4} \mathbf{r}_{n+1}^2 \text{tr}(\mathbf{P} \boldsymbol{\Sigma} \boldsymbol{\chi}_{n+1} \mathbf{M}_n \boldsymbol{\chi}_{n+1} \boldsymbol{\Sigma}),$$

so (B.18) and (B.19) follow from conditional Hölder and (B.15).

If  $\mathbf{P}$  has bounded support in  $\{1, \dots, r\}$ , PSD and bounded support give

$$\|\mathbf{K}_n^{(\mathbf{P})}\|_F \leq C \left(u_n^{(\mathbf{P})}\right)^{1/2},$$

which proves (ii).

For (iii), note that

$$u_{n+1} = s_n + \left\langle \frac{1}{2}(\boldsymbol{\Sigma} \mathbf{M}_n + \mathbf{M}_n \boldsymbol{\Sigma}), \boldsymbol{\xi}_{n+1} \right\rangle.$$

Lemma 10 gives

$$\|\mathbf{r}_{n+1} - e_n\|_{L^p(\mathbb{P}(\cdot | \mathcal{F}_n))} \leq Cp$$

and

$$\|\mathbf{r}_{n+1} u_{n+1} - \mathbb{E}[\mathbf{r}_{n+1} u_{n+1} | \mathcal{F}_n]\|_{L^p(\mathbb{P}(\cdot | \mathcal{F}_n))} \leq Cp^2.$$

This proves (B.22). The estimates for  $q_{n+1}^e$  are (B.18)–(B.19) with  $\mathbf{P} = \mathbf{I}_d$ . ■

#### B.4. Discrete bias tracking

Define

$$a_n^e := 2 \text{tr}(\boldsymbol{\Sigma} \mathbf{M}_n^2) - 2 \text{tr}(\mathbf{T}_* \mathbf{M}_n).$$

**Lemma 13 (Discrete bias recursion)** *One has*

$$e_{n+1} = \left(1 - \frac{\eta}{4\delta} - \eta s_n\right) e_n + \eta a_n^e + \eta \zeta_{n+1}^e - \eta^2 q_{n+1}^e. \quad (\text{B.24})$$

**Proof** Taking traces in (B.5),

$$\text{tr}(\mathbf{M}_{n+1}) = \text{tr}(\mathbf{M}_n) + \eta \mathbf{r}_{n+1} u_{n+1} + \eta^2 q_{n+1}^e.$$

Together with (B.3), this gives

$$e_{n+1} = e_n - \frac{\eta}{4\delta} \mathbf{r}_{n+1} - \eta \mathbf{r}_{n+1} u_{n+1} - \eta^2 q_{n+1}^e.$$

The Gaussian identity used in Proposition 6 gives

$$\mathbb{E}[\mathbf{r}_{n+1} u_{n+1} | \mathcal{F}_n] = s_n e_n - a_n^e.$$

Substituting (B.21) proves (B.24). ■

Set

$$\epsilon_N := \delta + R_s \eta \delta (1 + L_2) + \sqrt{\frac{p_d \eta}{\delta}} + \frac{\eta p_d^3}{\delta} + R_s (1 + L_2) \sqrt{p_d \eta^3 \delta} + R_s (1 + L_2) \eta^2 p_d^5. \quad (\text{B.25})$$

**Proposition 14 (Discrete bias tracking)** *Assume*

$$\frac{\eta}{4\delta} \leq \frac{1}{4}. \quad (\text{B.26})$$

There exists an event  $\mathcal{G}_e$  with

$$\mathbb{P}(\mathcal{G}_e) \geq 1 - Cd^{-K}$$

such that on  $\mathcal{G}_{\text{init}} \cap \mathcal{G}_e$ , up to  $\tau_\Delta$ ,

$$|e_n| \leq \exp\left(-\frac{n\eta}{8\delta}\right) |e_0| + C\epsilon_N, \quad 0 \leq n \leq N \wedge \tau_\Delta, \quad (\text{B.27})$$

and

$$\eta \sum_{m=0}^{n-1} |e_m| \leq C\delta |e_0| + C\Lambda_d \epsilon_N, \quad 0 \leq n \leq N \wedge \tau_\Delta. \quad (\text{B.28})$$

**Proof** On  $\{n < \tau_\Delta\}$ ,

$$\|\mathbf{M}_n\|_F \leq 5, \quad |a_n^e| + s_n \leq C,$$

where the bound on  $a_n^e$  uses  $\|\mathbf{T}_\star\|_F = O(1)$  from (A.2). For sufficiently large  $d$ , (B.26) and  $\delta \rightarrow 0$  imply

$$0 \leq 1 - \frac{\eta}{4\delta} - \eta s_n \leq 1 - \frac{\eta}{8\delta}.$$

Iterating (B.24) with this stable factor gives a deterministic contribution  $O(\delta)$  from  $a_n^e$  and  $O(R_s \eta \delta (1 + L_2))$  from the predictable part of  $q_{n+1}^e$ . Lemma 11 applied to the stable convolutions of  $\zeta_{n+1}^e$  and of the centered part of  $q_{n+1}^e$ , using (B.22) and (B.23), gives the remaining terms in (B.25). This proves (B.27). Summing the same bound over  $m$  and using  $N\eta \leq C\Lambda_d$  gives (B.28). ■

## B.5. Discrete Frobenius localization with explicit remainders

Define the localization remainder scale

$$I_N := \sqrt{p_d \eta \Lambda_d} + \eta p_d^3 + R_s (1 + L_2) \left( \eta \Lambda_d + \sqrt{p_d \eta^3 \Lambda_d} + \eta^2 p_d^5 \right).$$

**Proposition 15 (Discrete Frobenius localization)** *There exists an event  $\mathcal{G}_{\text{en}}$  with*

$$\mathbb{P}(\mathcal{G}_{\text{en}}) \geq 1 - Cd^{-K}$$

such that on  $\mathcal{G}_{\text{init}} \cap \mathcal{G}_e \cap \mathcal{G}_{\text{en}}$ , if

$$\Lambda_d \epsilon_N \rightarrow 0, \quad I_N \rightarrow 0, \quad (\text{B.29})$$

then

$$\tau_\Delta > N, \quad \sup_{0 \leq n \leq N} \|\mathbf{M}_n\|_F \leq 5, \quad I_N := \eta \sum_{n=0}^{N-1} |e_n| = o(1).$$

**Proof** By Proposition 14,

$$I_N \leq C\delta|e_0| + C\Lambda_d\epsilon_N = o(1).$$

Let

$$\mathcal{E}_n := \|\Delta_n\|_F^2.$$

On  $\{n < \tau_\Delta\}$ , write

$$\mathcal{D}_n := \mathcal{F}(\mathbf{M}_n) + e_n\mathcal{B}(\mathbf{M}_n).$$

Then (B.5) gives

$$\Delta_{n+1} = \Delta_n - \eta\mathcal{D}_n - \eta\mathbf{Z}_{n+1} - \eta^2\mathbf{Q}_{n+1}.$$

Expanding the square,

$$\mathcal{E}_{n+1} - \mathcal{E}_n = -2\eta\langle \Delta_n, \mathcal{F}(\mathbf{M}_n) \rangle - 2\eta e_n \langle \Delta_n, \mathcal{B}(\mathbf{M}_n) \rangle - 2\eta\langle \Delta_n, \mathbf{Z}_{n+1} \rangle + \mathcal{R}_{n+1},$$

where the full second-order remainder is

$$\begin{aligned} \mathcal{R}_{n+1} := & \eta^2\|\mathcal{D}_n\|_F^2 + 2\eta^2\langle \mathcal{D}_n, \mathbf{Z}_{n+1} \rangle + \eta^2\|\mathbf{Z}_{n+1}\|_F^2 - 2\eta^2\langle \Delta_n, \mathbf{Q}_{n+1} \rangle \\ & + 2\eta^3\langle \mathcal{D}_n + \mathbf{Z}_{n+1}, \mathbf{Q}_{n+1} \rangle + \eta^4\|\mathbf{Q}_{n+1}\|_F^2. \end{aligned}$$

The monotonicity identity (??) gives

$$\langle \Delta_n, \mathcal{F}(\mathbf{M}_n) \rangle = 2\|\Sigma^{1/2}\Delta_n\mathbf{M}_n^{1/2}\|_F^2 \geq 0.$$

Moreover,

$$|e_n\langle \Delta_n, \mathcal{B}(\mathbf{M}_n) \rangle| \leq C|e_n|.$$

The remainder terms are controlled by Lemmas 10, 12, and 11, and the scale conditions (B.29) ensure that their cumulative contribution is  $o(1)$  uniformly over  $n \leq N$  on  $\mathcal{G}_{\text{en}}$ . Hence

$$\sup_{m \leq N \wedge \tau_\Delta} \mathcal{E}_m \leq \mathcal{E}_0 + CI_N + o(1).$$

On  $\mathcal{G}_{\text{init}}$ ,  $\mathcal{E}_0 \leq 9$ . Thus the right-hand side is strictly less than 16 for all large  $d$ , so  $\tau_\Delta > N$ . Finally,

$$\|\mathbf{M}_n\|_F \leq \|\mathbf{S}^*\|_F + \|\Delta_n\|_F \leq 5$$

for all  $0 \leq n \leq N$ . ■

## B.6. Rest-block reduction for SGD

Write

$$\mathbf{M}_n = \begin{pmatrix} \mathbf{H}_n & \mathbf{R}_n \\ \mathbf{R}_n^\top & \mathbf{J}_n \end{pmatrix}, \quad \theta_n := \text{tr}(\mathbf{J}_n).$$

Inside the rest block, define

$$a_n := \text{tr}(\mathbf{P}_s\mathbf{M}_n), \quad c_n := \text{tr}(\mathbf{P}_t\mathbf{M}_n), \quad \theta_n = a_n + c_n, \quad \Theta_N := \eta \sum_{n=0}^{N-1} \theta_n.$$

**Lemma 16 (Projector trace recursion)** *For every diagonal projector  $\mathbf{P}$ ,*

$$u_{n+1}^{(\mathbf{P})} \leq (1 + 2\eta\tau(\mathbf{P}) + \eta\sigma(\mathbf{P})|e_n|) u_n^{(\mathbf{P})} + \eta d_{n+1}^{(\mathbf{P})} + \eta^2 q_{n+1}^{(\mathbf{P})}. \quad (\text{B.30})$$

**Proof** Apply  $\text{tr}(\mathbf{P}\cdot)$  to (B.5). Since  $\mathbf{P}$ ,  $\Sigma$ , and  $\mathbf{T}_*$  are diagonal and  $\mathbf{M}_n \succeq 0$ ,

$$\text{tr}(\mathbf{P}\mathcal{F}^\sharp(\mathbf{M}_n, e_n)) \leq (2\tau(\mathbf{P}) + \sigma(\mathbf{P})|e_n|) u_n^{(\mathbf{P})}.$$

This proves (B.30). ■

Set

$$\tau_t := \tau(\mathbf{P}_t), \quad \sigma_t := \sigma(\mathbf{P}_t), \quad \gamma_t := 1 + 2\eta\tau_t.$$

**Proposition 17 (Tail trace bound)** *There exists an event  $\mathcal{G}_t$  with*

$$\mathbb{P}(\mathcal{G}_t) \geq 1 - Cd^{-K}$$

*such that on  $\mathcal{G}_{\text{init}} \cap \mathcal{G}_e \cap \mathcal{G}_{\text{en}} \cap \mathcal{G}_t$ ,*

$$\sup_{0 \leq n \leq N} c_n \leq e^{\sigma_t I_N} [\gamma_t^N c_0 + CR_s \gamma_t^N \eta^2 N(1 + L_2) + \mathfrak{M}_t], \quad (\text{B.31})$$

*where*

$$\begin{aligned} \mathfrak{M}_t := & C\gamma_t^N \left[ \sqrt{p_d \eta(1 + L_2) \Lambda_d} + \eta p_d^3 (1 + L_2)^{1/2} \right. \\ & \left. + R_s (1 + L_2) \left( \sqrt{p_d \eta^3 \Lambda_d} + \eta^2 p_d^5 \right) \right]. \end{aligned}$$

*If*

$$\left( \frac{d}{Z_r} \right)^\rho \frac{1}{d} \sum_{i > m_\rho(\bar{t})} \sigma_i \rightarrow 0, \quad \left( \frac{d}{Z_r} \right)^{2\rho} R_s \eta (1 + L_2) \Lambda_d^8 \rightarrow 0, \quad (\text{B.32})$$

*then*

$$\sup_{0 \leq n \leq N} c_n = o(1), \quad \eta \sum_{n=0}^{N-1} c_n = o(1).$$

**Proof** Apply Lemma 16 with  $\mathbf{P} = \mathbf{P}_t$ :

$$c_{n+1} \leq (\gamma_t + \eta\sigma_t|e_n|)c_n + \eta d_{n+1}^{(\mathbf{P}_t)} + \eta^2 q_{n+1}^{(\mathbf{P}_t)}.$$

Iterating and using

$$\prod_{\ell=m+1}^{n-1} (\gamma_t + \eta\sigma_t|e_\ell|) \leq \gamma_t^{n-1-m} e^{\sigma_t I_N}$$

gives (B.31), after applying Lemmas 12 and 11 to the martingale terms.

Since

$$\gamma_t^N \leq \left( \frac{d}{Z_r} \right)^{\rho+o(1)}$$

by the definition of  $\mathbf{P}_t$ , and

$$c_0 \leq \frac{C}{d} \sum_{i > m_\rho(\bar{t})} \sigma_i$$

by (B.11), the initial term is  $o(1)$ . The martingale and predictable second-order terms are  $o(1)$  under (B.32). Summing the same recursion over  $n$  gives the integrated bound.  $\blacksquare$

Set

$$\tau_s := \tau(\mathbf{P}_s), \quad \sigma_s := \sigma(\mathbf{P}_s), \quad \gamma_s := 1 + 2\eta\tau_s.$$

**Proposition 18 (Slow trace bound)** *If  $k < r$ , there exists an event  $\mathcal{G}_s$  with*

$$\mathbb{P}(\mathcal{G}_s) \geq 1 - Cd^{-K}$$

*such that on  $\mathcal{G}_{\text{init}} \cap \mathcal{G}_e \cap \mathcal{G}_{\text{en}} \cap \mathcal{G}_s$ ,*

$$\sup_{0 \leq n \leq N} a_n \leq Ce^{\sigma_s I_N} \gamma_s^N \left( \frac{1}{d} + R_s \eta p_d \right), \quad (\text{B.33})$$

$$\eta \sum_{n=0}^{N-1} a_n \leq Ce^{\sigma_s I_N} \gamma_s^N \left( \frac{1}{d} + R_s \eta p_d \right). \quad (\text{B.34})$$

If

$$R_s \eta p_d \left( \frac{d}{Z_r} \right)^{a_s} \rightarrow 0,$$

then both quantities in (B.33)–(B.34) are  $o(1)$ . If  $k = r$ , then  $a_n \equiv 0$ .

**Proof** Assume  $k < r$ . Apply Lemma 16 with  $\mathbf{P} = \mathbf{P}_s$ . Let

$$R_\sharp := K_0 \left( \frac{1}{d} + R_s \eta p_d \right),$$

where  $K_0$  is a sufficiently large constant, and stop the normalized process  $\gamma_s^{-n} a_n$  at its first exit above  $R_\sharp$ . On the stopped interval,

$$a_m \leq R_\sharp e^{\sigma_s I_N} \gamma_s^m.$$

The bounded-support estimate (B.20) and Lemma 11 give

$$\left\| \sup_{n \leq N} \left| \eta \sum_{m=0}^{n-1} \gamma_s^{-m-1} d_{m+1}^{(\mathbf{P}_s)} \right| \right\|_{L^p d} \leq C \left( \sqrt{p_d \eta R_\sharp} + \eta p_d^3 R_\sharp^{1/2} \right),$$

and the centered second-order part is bounded by

$$CR_s \left( \sqrt{p_d \eta^3 \Lambda_d} + \eta^2 p_d^5 \right).$$

The predictable second-order contribution is  $O(R_s \eta)$ . These terms are absorbed into  $R_\sharp$  with probability at least  $1 - Cd^{-K}$  if  $K_0$  is large. Since  $a_0 \leq C/d$ , the stopped process does not exit before  $N$ . Finally,

$$\gamma_s^N \leq \left( \frac{d}{Z_r} \right)^{a_s + o(1)},$$

which proves the result. If  $k = r$ , then  $\mathbf{P}_s = 0$ .  $\blacksquare$

**Corollary 19 (Rest block and cross block for SGD)** *On the intersection of the good events constructed above,*

$$\sup_{0 \leq n \leq N} \theta_n = o(1), \quad \Theta_N = o(1), \quad \sup_{0 \leq n \leq N} \|\mathbf{J}_n\|_F = o(1), \quad \sup_{0 \leq n \leq N} \|\mathbf{R}_n\|_F = o(1).$$

### B.7. Fast-block bootstrap for SGD

Assume  $k \geq 1$ . Write

$$\mathbf{H}_n = \mathbf{D}_n + \mathbf{E}_n, \quad \mathbf{D}_n := \text{diag}(h_{1,n}, \dots, h_{k,n}), \quad (\mathbf{E}_n)_{ii} = 0, \quad \mathbf{\Gamma}_n := \mathbf{R}_n \mathbf{R}_n^\top.$$

For  $1 \leq i \leq k$ , define the matched logistic sequence

$$\tilde{\ell}_{i,n+1} = \tilde{\ell}_{i,n} \left( 1 + \eta(\tau_i - \sigma_i \tilde{\ell}_{i,n}) \right)^2, \quad \tilde{\ell}_{i,0} = h_{i,0}, \quad (\text{B.35})$$

and

$$\tilde{\mathbf{L}}_n^f := \text{diag}(\tilde{\ell}_{1,n}, \dots, \tilde{\ell}_{k,n}).$$

For  $i \neq j$ , define

$$Y_{ij,n} := \frac{(\mathbf{E}_n)_{ij}}{\sqrt{\tilde{\ell}_{i,n} \tilde{\ell}_{j,n}}}, \quad X_n := \max_{i \neq j} \sup_{0 \leq m \leq n} |Y_{ij,m}|,$$

and

$$W_n := \max_{1 \leq i \leq k} \sup_{0 \leq m \leq n} \left| \frac{h_{i,m}}{\tilde{\ell}_{i,m}} - 1 \right|.$$

**Lemma 20 (Reciprocal sums for the matched logistics)** *On  $\mathcal{G}_{\text{init}}$ , for every  $1 \leq i \leq k$ ,*

$$\eta \sum_{n=0}^{N-1} \tilde{\ell}_{i,n}^{-1} \leq Cd, \quad (\text{B.36})$$

and

$$\eta \sum_{n=0}^{N-1} \tilde{\ell}_{i,n}^{-2} \leq Cd^2. \quad (\text{B.37})$$

**Proof** Fix  $i \leq k$  and set  $s^* := \tau_i / \sigma_i = s_i^* = \sigma_i \lambda_i / Z_r$ . Define

$$z_n := \frac{\tilde{\ell}_{i,n}}{s^*} \in (0, 1].$$

Then

$$z_{n+1} = z_n \left( 1 + \eta \tau_i (1 - z_n) \right)^2, \quad 0 < z_0 = \frac{h_{i,0}}{s^*}.$$

On  $\mathcal{G}_{\text{init}}$ ,  $h_{i,0} \asymp \sigma_i / d$ , hence for fixed  $i$ ,

$$z_0 \asymp \frac{\sigma_i / d}{\sigma_i \lambda_i / Z_r} = \frac{Z_r}{d \lambda_i} \asymp \frac{Z_r}{d}.$$

Let  $n_{1/2} := \inf\{n \geq 0 : z_n \geq 1/2\}$ . For  $n < n_{1/2}$ ,  $1 - z_n \geq 1/2$  and hence

$$z_{n+1} \geq z_n \left(1 + \frac{1}{2}\eta\tau_i\right)^2.$$

Therefore

$$\eta \sum_{n=0}^{(n_{1/2}-1) \wedge (N-1)} z_n^{-1} \leq \eta z_0^{-1} \sum_{n \geq 0} \left(1 + \frac{1}{2}\eta\tau_i\right)^{-2n} \leq C z_0^{-1} \asymp C \frac{d}{Z_r}.$$

For  $n \geq n_{1/2}$ ,  $z_n^{-1} \leq 2$ , hence

$$\eta \sum_{n=n_{1/2} \wedge N}^{N-1} z_n^{-1} \leq 2\eta N \leq C \log \frac{d}{Z_r}.$$

Combining and using  $\tilde{\ell}_{i,n}^{-1} = (s^*)^{-1} z_n^{-1}$  with  $(s^*)^{-1} = Z_r / (\sigma_i \lambda_i) = O(Z_r)$  (for fixed  $i$ ), we obtain (B.36).

The bound (B.37) is analogous, using  $z_n^{-2}$  and the same splitting at  $n_{1/2}$ , which yields  $\eta \sum z_n^{-2} \lesssim z_0^{-2} + \eta N \asymp d^2/Z_r^2 + \log(d/Z_r)$  and multiplying by  $(s^*)^{-2} \asymp Z_r^2$ . ■

**Lemma 21 (Fast-block recursion)** *On  $\mathcal{G}_{\text{loc},n}$ , for  $1 \leq i \leq k$ ,*

$$\begin{aligned} h_{i,n+1} &= h_{i,n} + 2\eta\tau_i h_{i,n} - 2\eta\sigma_i h_{i,n}^2 - \eta q_{i,n}^f + \eta\sigma_i e_n h_{i,n} + \eta z_{i,n+1}^d + \eta^2 b_{i,n+1}^d, \quad (\text{B.38}) \\ q_{i,n}^f &:= 2\sigma_i \sum_{j \neq i} (\mathbf{E}_n)_{ij}^2 + 2\sigma_i (\mathbf{\Gamma}_n)_{ii} \geq 0. \end{aligned}$$

For  $i \neq j$ ,

$$\begin{aligned} (\mathbf{E}_{n+1})_{ij} &= (\mathbf{E}_n)_{ij} + \eta c_{ij,n} (\mathbf{E}_n)_{ij} - \eta(\sigma_i + \sigma_j) (\mathbf{E}_n^2)_{ij} - \eta(\sigma_i + \sigma_j) (\mathbf{\Gamma}_n)_{ij} \\ &\quad + \frac{\eta}{2} (\sigma_i + \sigma_j) e_n (\mathbf{E}_n)_{ij} + \eta z_{ij,n+1}^{\text{od}} + \eta^2 b_{ij,n+1}^{\text{od}}, \quad (\text{B.39}) \\ c_{ij,n} &:= \tau_i + \tau_j - (\sigma_i + \sigma_j) (h_{i,n} + h_{j,n}). \end{aligned}$$

Here

$$z_{i,n+1}^d := (\mathbf{Z}_{n+1})_{ii}, \quad b_{i,n+1}^d := (\mathbf{Q}_{n+1})_{ii}, \quad z_{ij,n+1}^{\text{od}} := (\mathbf{Z}_{n+1})_{ij}, \quad b_{ij,n+1}^{\text{od}} := (\mathbf{Q}_{n+1})_{ij}.$$

**Proof** Take the  $(i, i)$  and  $(i, j)$  components of (B.5). The algebra is left to the reader. ■

Fix

$$0 < \varepsilon_\star < \frac{\sigma_{\min}}{2\sigma_{\max}},$$

and define

$$\nu_\star := \inf\{m \geq 0 : W_m > \varepsilon_\star/2\}, \quad \mu_\theta := \inf\{m \geq 0 : \theta_m > 1\}.$$

**Lemma 22 (Fast noise bounds)** *On  $\{n < \nu_* \wedge \mu_\theta\} \cap \mathcal{G}_{\text{loc},n}$ , for all  $p \geq 2$ ,*

$$\|z_{i,n+1}^{\text{d}}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq Cp^2 \tilde{\ell}_{i,n}^{1/2}, \quad (\text{B.40})$$

$$\|z_{ij,n+1}^{\text{od}}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq Cp^2 \left( \tilde{\ell}_{i,n}^{1/2} + \tilde{\ell}_{j,n}^{1/2} \right). \quad (\text{B.41})$$

*Writing*

$$b_{i,n+1}^{\text{d}} = \bar{b}_{i,n}^{\text{d}} + \widehat{b}_{i,n+1}^{\text{d}}, \quad b_{ij,n+1}^{\text{od}} = \bar{b}_{ij,n}^{\text{od}} + \widehat{b}_{ij,n+1}^{\text{od}},$$

*with conditional means as predictable parts,*

$$|\bar{b}_{i,n}^{\text{d}}| + |\bar{b}_{ij,n}^{\text{od}}| \leq C, \quad \|\widehat{b}_{i,n+1}^{\text{d}}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} + \|\widehat{b}_{ij,n+1}^{\text{od}}\|_{L^p(\mathbb{P}(\cdot|\mathcal{F}_n))} \leq Cp^4.$$

**Proof** For the diagonal first-order term,

$$z_{i,n+1}^{\text{d}} = \mathbf{r}_{n+1} \langle \mathbf{K}_{i,n}^{\text{d}}, \boldsymbol{\chi}_{n+1} \rangle - \mathbb{E}[\mathbf{r}_{n+1} \langle \mathbf{K}_{i,n}^{\text{d}}, \boldsymbol{\chi}_{n+1} \rangle | \mathcal{F}_n],$$

where

$$\mathbf{K}_{i,n}^{\text{d}} = \frac{1}{2} (\mathbf{M}_n \mathbf{e}_i \mathbf{e}_i^\top \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{M}_n).$$

On  $\{n < \nu_*\}$ ,  $h_{i,n} \asymp \tilde{\ell}_{i,n}$ ; PSD gives

$$\|\mathbf{K}_{i,n}^{\text{d}}\|_F \leq C \tilde{\ell}_{i,n}^{1/2}.$$

The off-diagonal case is identical. Lemma 10 proves (B.40)–(B.41). The second-order bounds follow from conditional Hölder, using  $k = O_{\tilde{\ell}}(1)$ ,  $\theta_n \leq 1$ , and  $\|\mathbf{M}_n\|_F \leq 5$ .  $\blacksquare$

**Lemma 23 (Fast-block algebra)** *On  $\{W_n \leq \varepsilon_*/2\}$ ,*

$$\left(1 - \frac{\varepsilon_*}{2}\right) \tilde{\ell}_{i,n} \leq h_{i,n} \leq \left(1 + \frac{\varepsilon_*}{2}\right) \tilde{\ell}_{i,n}. \quad (\text{B.42})$$

*Moreover,*

$$\frac{q_{i,n}^{\text{f}}}{\tilde{\ell}_{i,n}} \leq C(X_n^2 + \theta_n), \quad (\text{B.43})$$

$$\frac{|(\mathbf{E}_n^2)_{ij}|}{\sqrt{\tilde{\ell}_{i,n} \tilde{\ell}_{j,n}}} \leq CX_n^2, \quad (\text{B.44})$$

$$\frac{|(\boldsymbol{\Gamma}_n)_{ij}|}{\sqrt{\tilde{\ell}_{i,n} \tilde{\ell}_{j,n}}} + \frac{(\boldsymbol{\Gamma}_n)_{ii}}{\tilde{\ell}_{i,n}} \leq C\theta_n. \quad (\text{B.45})$$

*Finally, for all sufficiently large  $d$ ,*

$$0 \leq \frac{1 + \eta c_{ij,n}}{\left(1 + \eta(\tau_i - \sigma_i \tilde{\ell}_{i,n})\right) \left(1 + \eta(\tau_j - \sigma_j \tilde{\ell}_{j,n})\right)} \leq 1. \quad (\text{B.46})$$

**Proof** The first display is the definition of  $W_n$ . The estimate (B.44) follows from

$$(\mathbf{E}_n)_{ab} = Y_{ab,n} \sqrt{\tilde{\ell}_{a,n} \tilde{\ell}_{b,n}}$$

and  $k = O_{\tilde{\tau}}(1)$ . We have

$$\|\text{row}_i(\mathbf{R}_n)\|_2^2 \leq h_{i,n} \theta_n,$$

which gives (B.45); (B.43) follows from the definition of  $q_{i,n}^f$ . The last claim follows by expanding numerator and denominator in (B.46), using (B.42),  $\varepsilon_\star < \sigma_{\min}/(2\sigma_{\max})$ , and  $\eta\tau_1 \leq 1/4$ .  $\blacksquare$

Define

$$\mathfrak{M}_h := C \left( \sqrt{dp_d\eta} + \sqrt{d}\eta p_d^3 \right), \quad \mathfrak{B}_h := C \left( d\sqrt{p_d\eta^3} + d\eta^2 p_d^5 \right).$$

**Proposition 24 (Fast martingale event)** *There exists an event  $\mathcal{G}_h$  with*

$$\mathbb{P}(\mathcal{G}_h) \geq 1 - Cd^{-K}$$

such that on  $\mathcal{G}_h$ , for all  $1 \leq i \leq k$  and  $i \neq j$ ,

$$\sup_{0 \leq n \leq N} \left| \sum_{\ell=0}^{(n-1) \wedge (\nu_\star \wedge \mu_\theta - 1)} \eta \frac{z_{i,\ell+1}^d}{\tilde{\ell}_{i,\ell}} \right| + \sup_{0 \leq n \leq N} \left| \sum_{\ell=0}^{(n-1) \wedge (\nu_\star \wedge \mu_\theta - 1)} \eta \frac{z_{ij,\ell+1}^{\text{od}}}{\sqrt{\tilde{\ell}_{i,\ell} \tilde{\ell}_{j,\ell}}} \right| \leq \mathfrak{M}_h, \quad (\text{B.47})$$

$$\sup_{0 \leq n \leq N} \left| \sum_{\ell=0}^{(n-1) \wedge (\nu_\star \wedge \mu_\theta - 1)} \eta^2 \frac{\widehat{b}_{i,\ell+1}^d}{\tilde{\ell}_{i,\ell}} \right| + \sup_{0 \leq n \leq N} \left| \sum_{\ell=0}^{(n-1) \wedge (\nu_\star \wedge \mu_\theta - 1)} \eta^2 \frac{\widehat{b}_{ij,\ell+1}^{\text{od}}}{\sqrt{\tilde{\ell}_{i,\ell} \tilde{\ell}_{j,\ell}}} \right| \leq \mathfrak{B}_h. \quad (\text{B.48})$$

**Proof** Apply Lemma 11 to the martingale sums in (B.47)–(B.48). Lemma 22 gives the one-step moments, and Lemma 20 gives the reciprocal variance sums. A union bound over the fixed fast block completes the proof.  $\blacksquare$

**Proposition 25 (Fast bootstrap for SGD)** *Work on the intersection of all good events needed above. Set*

$$\begin{aligned} \Xi_N &:= e^{\sigma_{\max} I_N} \left( \delta_0^{\text{coh}} + C\Theta_N + Cd\eta + C\mathfrak{M}_h + C\mathfrak{B}_h \right), \\ R_W &:= C \left( N\eta \Xi_N^2 + \Theta_N + d\eta + \mathfrak{M}_h + \mathfrak{B}_h + I_N \right). \end{aligned}$$

If

$$4CN\eta\Xi_N \leq 1, \quad R_W \leq \frac{\varepsilon_\star}{2}, \quad (\text{B.49})$$

then, for all  $0 \leq n \leq N$ ,

$$X_n \leq 2\Xi_N, \quad W_n \leq R_W.$$

Consequently,

$$\sup_{0 \leq n \leq N} \|\mathbf{E}_n\|_F \leq C\Xi_N, \quad \sup_{0 \leq n \leq N} \|\mathbf{D}_n - \tilde{\mathbf{L}}_n^f\|_F \leq CR_W.$$

**Proof** Since  $\Theta_N = o(1)$  on the good event,  $\mu_\theta > N$  for all large  $d$ . Let

$$\nu := \inf\{n \geq 0 : W_n > R_W\}.$$

On  $n < \nu$ , Lemma 23 applies. Divide (B.39) by  $\sqrt{\tilde{\ell}_{i,n+1}\tilde{\ell}_{j,n+1}}$ . Using (B.46),

$$|Y_{ij,n+1}| \leq (1 + C\eta|e_n|)|Y_{ij,n}| + C\eta X_n^2 + C\eta\theta_n + \xi_{ij,n+1},$$

where the cumulative normalized martingale and centered second-order contributions are bounded by

$$\sup_{0 \leq n \leq N} \left| \sum_{m=0}^{n-1} \xi_{ij,m+1} \right| \leq C(d\eta + \mathfrak{M}_h + \mathfrak{B}_h)$$

thanks to Proposition 24; the predictable second-order contribution is  $O(d\eta)$  by Lemma 20. Hence

$$X_n \leq \Xi_N + C\eta \sum_{m=0}^{n-1} X_m^2 \quad (n \leq \nu).$$

The quadratic induction under (B.49) gives  $X_n \leq 2\Xi_N$  for  $n \leq \nu$ .

For the diagonal part, set

$$w_{i,n} := \frac{h_{i,n}}{\tilde{\ell}_{i,n}} - 1.$$

Subtract (B.35) from (B.38), divide by  $\tilde{\ell}_{i,n+1}$ , and use Lemma 23. The deterministic logistic term is contracting up to harmless  $O(\eta|e_n|)$  factors, while

$$\frac{q_{i,n}^f}{\tilde{\ell}_{i,n}} \leq C(X_n^2 + \theta_n)$$

by (B.43). Lemma 20 and Proposition 24 control the normalized stochastic and second-order terms. Consequently,

$$W_n \leq C(N\eta\Xi_N^2 + \Theta_N + d\eta + \mathfrak{M}_h + \mathfrak{B}_h + I_N) = R_W \quad (n \leq \nu).$$

Thus the bootstrap improves strictly, so  $\nu > N$ . The matrix bounds follow from the definitions of  $X_n, W_n$  and  $k = O_{\tilde{\ell}}(1)$ .  $\blacksquare$

**Lemma 26 (Matched and canonical logistics in discrete time)** *Assume*

$$\frac{pd}{r_s} \rightarrow 0.$$

*Then*

$$\|\tilde{\mathbf{L}}_N^f - \mathbf{L}_N^f\|_F = o(1) \tag{B.50}$$

*with probability*  $1 - o(1)$ .

**Proof** Fix  $i \leq k$  and define

$$a_{i,n} := \frac{\ell_{i,n}}{s_i^*}, \quad \tilde{a}_{i,n} := \frac{\tilde{\ell}_{i,n}}{s_i^*}.$$

Both sequences satisfy

$$u_{n+1} = u_n(1 + \eta\tau_i(1 - u_n))^2.$$

By (B.12),

$$\frac{d h_{i,0}}{\sigma_i} = 1 + O_{\mathbb{P}}\left(\sqrt{\frac{pd}{r_s}}\right),$$

so  $a_{i,0}$  and  $\tilde{a}_{i,0}$  have the same  $Z_r/d$  scale. Since  $i \leq k$  implies  $\bar{\ell}\kappa_i \geq 1 + \delta_{\text{gap}}$ , Lemma 8 yields

$$a_{i,N} \rightarrow 1, \quad \tilde{a}_{i,N} \rightarrow 1.$$

Therefore  $|\tilde{\ell}_{i,N} - \ell_{i,N}| = o(1)$ . Summing over  $O_{\bar{\ell}}(1)$  fast indices proves (B.50).  $\blacksquare$

**Theorem 27 (Comparison with the centered ideal discrete flow)** *Assume*

$$\eta\tau_1 \leq \frac{1}{4}, \quad \frac{\eta}{4\delta} \leq \frac{1}{4}, \quad r_s \gg pd\Lambda_d^2.$$

*Assume also*

$$\Lambda_d \epsilon_N \rightarrow 0, \quad d\eta\Lambda_d^8 \rightarrow 0, \tag{B.51}$$

*and*

$$\left(\frac{d}{Z_r}\right)^{2\rho} R_s \eta (1 + L_2) \Lambda_d^8 \rightarrow 0, \quad \left(\frac{d}{Z_r}\right)^{\rho} \frac{1}{d} \sum_{i > m_{\rho}(\bar{\ell})} \sigma_i \rightarrow 0. \tag{B.52}$$

*If  $k < r$ , assume further*

$$R_s \eta p d \left(\frac{d}{Z_r}\right)^{a_s} \rightarrow 0.$$

*Then there exists an event  $\mathcal{G}$  with*

$$\mathbb{P}(\mathcal{G}) \geq 1 - Cd^{-K} \tag{B.53}$$

*such that on  $\mathcal{G}$ ,*

$$\|\mathbf{M}_N - \mathbf{L}_N\|_F = o(1), \quad |e_N| = o(1). \tag{B.54}$$

**Proof** Let  $\mathcal{G}$  be the intersection of  $\mathcal{G}_{\text{init}}$ ,  $\mathcal{G}_e$ ,  $\mathcal{G}_{\text{en}}$ ,  $\mathcal{G}_t$ , together with  $\mathcal{G}_s$  when  $k < r$ , and  $\mathcal{G}_h$  when  $k \geq 1$ . Then (B.53) holds.

The assumptions (B.51)–(B.52) imply  $\iota_N = o(1)$ , so Proposition 15 gives

$$I_N = o(1), \quad e_N = o(1), \quad \sup_{0 \leq n \leq N} \|\mathbf{M}_n\|_F \leq 5.$$

Corollary 19 gives

$$\sup_{0 \leq n \leq N} \|\mathbf{J}_n\|_F = o(1), \quad \sup_{0 \leq n \leq N} \|\mathbf{R}_n\|_F = o(1), \quad \Theta_N = o(1).$$

If  $k = 0$ , the fast block is empty. By the off-transition gap, choose  $\delta_{\text{thr}} > 0$  such that  $\bar{t}\kappa_i \leq 1 - \delta_{\text{thr}}$  for all  $i \leq r$  in the relevant finite set. Lemma 8 gives  $\ell_{i,N} \rightarrow 0$  for all  $i \leq r$ , and (B.7) is nonincreasing for  $i > r$ . Hence  $\|\mathbf{L}_N\|_F = o(1)$ , and Lemma ?? gives  $\|\mathbf{M}_N - \mathbf{L}_N\|_F = o(1)$ .

Assume  $k \geq 1$ . By (B.13) and  $r_s \gg p_d \Lambda_d^2$ ,

$$\delta_0^{\text{coh}} = O\left(\sqrt{\frac{p_d}{r_s}}\right) = o(\Lambda_d^{-1}).$$

The scale conditions imply

$$\mathfrak{M}_h = o(1), \quad \mathfrak{B}_h = o(1), \quad d\eta = o(1), \quad \Theta_N = o(1), \quad I_N = o(1).$$

Therefore

$$\Xi_N = o(1), \quad R_W = o(1), \quad N\eta\Xi_N = o(1).$$

Proposition 25 gives

$$\|\mathbf{H}_N - \tilde{\mathbf{L}}_N^f\|_F = o(1),$$

and Lemma 26 gives

$$\|\tilde{\mathbf{L}}_N^f - \mathbf{L}_N^f\|_F = o(1).$$

The rest-block estimates and Lemma ?? complete the proof of (B.54).  $\blacksquare$

For the discrete dynamics, define the one-step conditional excess risk

$$R_n := \frac{1}{8}\mathbb{E}[(y_{n+1} - \hat{y}_{n+1})^2 \mid \mathcal{F}_n] - \frac{1}{8}\mathbb{E}[\varepsilon_1^2].$$

Then

$$R_n = \frac{1}{4}\|\mathbf{S}^* - \mathbf{M}_n\|_F^2 + \frac{1}{8}e_n^2. \quad (\text{B.55})$$

We also write  $R(\eta n) := R_n$ .

**Corollary 28 (Discrete risk asymptotics and main-text form)** *Under the assumptions of Theorem 27, the following hold.*

(i) For fixed reduced time  $\bar{t}$ ,

$$\|\mathbf{S}^* - \mathbf{M}_N\|_F^2 = \sum_{i > m_*(\bar{t})} (s_i^*)^2 + o_{\mathbb{P}}(1), \quad (\text{B.56})$$

and therefore

$$R_N = \frac{1}{4} \sum_{i > m_*(\bar{t})} (s_i^*)^2 + o_{\mathbb{P}}(1). \quad (\text{B.57})$$

(ii) Away from the transition set,

$$R_N = \Theta\left(\bar{t}^{-\frac{2(\alpha+\beta)-1}{\alpha+2\beta}}\right). \quad (\text{B.58})$$

(iii) If  $Z_r = \Theta(1)$ , write  $n = \lfloor t \log d \rfloor$  and suppose the corresponding reduced time satisfies  $\bar{t} \asymp \eta t$  and remains off the transition set. In the mildly overparameterized regime  $r_s \gg \text{polylog } d$ ,

$$R(\eta t \log d) \asymp \Theta\left((\eta t)^{-\frac{2(\alpha+\beta)-1}{\alpha+2\beta}} + r_s^{-2(\alpha+\beta)+1}\right) \quad (\text{B.59})$$

with high probability, and almost surely after Borel–Cantelli.

**Proof** Theorem 27 gives

$$\|\mathbf{M}_N - \mathbf{L}_N\|_F = o_{\mathbb{P}}(1), \quad e_N = o_{\mathbb{P}}(1).$$

Thus

$$\|\mathbf{S}^* - \mathbf{M}_N\|_F^2 = \|\mathbf{S}^* - \mathbf{L}_N\|_F^2 + o_{\mathbb{P}}(1).$$

By Lemma 8 and the off-transition gap, the coordinates with  $\bar{t}\kappa_i > 1$  converge to  $s_i^*$ , while the subcritical coordinates are negligible. Hence (B.56). Equation (B.57) follows from (B.55).

The power law (B.58) is the same integral-test estimate as in Corollary ???. Finally, when  $Z_r = \Theta(1)$ ,  $N_{\text{eff}} \asymp \eta^{-1} \log d$ , so  $n = t \log d$  corresponds to reduced time  $\bar{t} \asymp \eta t$ . Combining the dynamical term with the rank floor in Lemma 5 gives (B.59). The almost-sure statement follows from Borel–Cantelli because all high-probability estimates are available with probability at least  $1 - C_K d^{-K}$  for arbitrary fixed  $K$ .  $\blacksquare$

### Appendix C. Proof of Proposition 3

Let  $\mathbf{G} := \frac{1}{r_s} \mathbf{W} \mathbf{W}^\top$  and redefine  $\hat{R}$  in (3.2) in terms of  $\mathbf{G}$ , i.e.,  $\hat{R}(\mathbf{G}) = \|\Sigma^{\frac{1}{2}} \mathbf{G} \Sigma^{\frac{1}{2}} - \hat{\mathbf{S}}\|_F^2$ . We observe that the minimizer of  $\hat{R}(\mathbf{G})$  is low-rank approximation of  $\Sigma^{-\frac{1}{2}} \hat{\mathbf{S}} \Sigma^{-\frac{1}{2}}$ . For strong anisotropy ( $\beta \gg 1/2$ ), driving  $\hat{R}(\mathbf{G})$  to zero strictly requires recovering the principal eigenvector of  $\Sigma$ , which aligns with  $\theta_1$ . Because estimating this principal direction reduces to a standard phase retrieval problem via second-order Stein-based spectral estimators, the desired sample complexity lower bound follows immediately from [19, Theorem 3].

### Appendix D. Normalized SGD Dynamics

We consider

$$y_{t+1} = \langle \mathbf{T}, \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top - \Sigma \rangle + \epsilon_{t+1} \quad \text{and} \quad \hat{y}(\mathbf{W}_t; \mathbf{x}_{t+1}) = \left\langle \frac{1}{r_s} \mathbf{W}_t \mathbf{W}_t^\top, \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right\rangle$$

where  $\|\epsilon_{t+1}\|_{\psi_2} \leq \sigma$ ,  $\mathbf{W}_t \in \mathbb{R}^{d \times r_s}$  and

$$r_s \gg \frac{d}{\text{polylog } d}.$$

We use  $\hat{y}_{t+1} := \hat{y}(\mathbf{W}_t; \mathbf{x}_{t+1})$  and consider

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

$$\nabla_{\text{St}} \mathcal{L}(\mathbf{W}_t) = \frac{-1}{r_s} \left( \mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^\top \right) y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{W}_t.$$

- We recall that  $\mathbf{G}_t = \mathbf{W}_t \mathbf{W}_t^\top$ .

Then, (SGD) reads

$$\begin{aligned}\hat{\mathbf{W}}_{t+1} &= \mathbf{W}_t + \frac{\eta}{2r_s} \underbrace{\left(\mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^\top\right) y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{W}_t}_{:= \nabla_{\text{St}} \mathbf{L}_{t+1}} \\ \mathbf{W}_{t+1} &= \hat{\mathbf{W}}_{t+1} \left( \mathbf{I}_{r_s} + \frac{\eta^2}{4r_s^2} \underbrace{\nabla_{\text{St}} \mathbf{L}_{t+1}^\top \nabla_{\text{St}} \mathbf{L}_{t+1}}_{:= \mathcal{P}_{t+1}} \right)^{-1/2}.\end{aligned}\quad (\text{SGD})$$

We observe that

$$\begin{aligned}\frac{\eta^2}{4r_s^2} \mathcal{P}_{t+1} &= \frac{\eta^2}{r_s^2} y_{t+1}^2 \mathbf{W}_t^\top \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \left(\mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^\top\right) \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{W}_t \\ &= \frac{\eta^2}{4r_s^2} y_{t+1}^2 \left\| \left(\mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^\top\right) \mathbf{x}_{t+1} \right\|_2^2 \mathbf{W}_t^\top \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{W}_t.\end{aligned}$$

Let

$$c_{t+1}^2 := \frac{\eta^2}{4r_s^2} \|\mathcal{P}_{t+1}\|_2 = \frac{\eta^2}{4r_s^2} y_{t+1}^2 \left\| \left(\mathbf{I}_d - \mathbf{W}_t \mathbf{W}_t^\top\right) \mathbf{x}_{t+1} \right\|_2^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2.$$

We define  $\mathbf{P}_{t+1} := \left(\mathbf{I}_{r_s} + \frac{\eta^2}{4r_s^2} \mathcal{P}_{t+1}\right)^{-1/2}$  and since  $\mathcal{P}_{t+1}$  is 1-rank, we have

$$\mathbf{P}_{t+1}^2 = \mathbf{I}_{r_s} - \frac{\eta^2}{4r_s^2} \frac{\mathcal{P}_{t+1}}{1 + c_{t+1}^2}.$$

By recalling that  $\mathbf{G}_t = \mathbf{W}_t \mathbf{W}_t^\top$ , we have

$$\begin{aligned}\mathbf{G}_{t+1} &= \hat{\mathbf{W}}_{t+1} \hat{\mathbf{W}}_{t+1}^\top + \hat{\mathbf{W}}_{t+1} (\mathbf{P}_{t+1}^2 - \mathbf{I}_{r_s}) \hat{\mathbf{W}}_{t+1}^\top \\ &= \mathbf{G}_t + \frac{\eta}{2r_s} (\mathbf{I}_d - \mathbf{G}_t) y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{G}_t + \frac{\eta}{2r_s} \mathbf{G}_t y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top (\mathbf{I}_d - \mathbf{G}_t) \\ &\quad + \frac{\eta^2}{4r_s^2} (\mathbf{I}_d - \mathbf{G}_t) y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{G}_t y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top (\mathbf{I}_d - \mathbf{G}_t) - \frac{\eta^2}{4r_s^2} \frac{\hat{\mathbf{W}}_{t+1} \mathcal{P}_{t+1} \hat{\mathbf{W}}_{t+1}^\top}{1 + c_{t+1}^2}.\end{aligned}$$

For  $\mathbf{T}_* := \Sigma \mathbf{T} \Sigma$  and  $\mathbf{N}_{t+1} := \frac{1}{2} y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top - \mathbf{T}_*$ , we write

$$\begin{aligned}\mathbf{G}_{t+1} &= \mathbf{G}_t + \frac{\eta}{r_s} (\mathbf{I}_d - \mathbf{G}_t) (\mathbf{T}_* + \mathbf{N}_{t+1}) \mathbf{G}_t + \frac{\eta}{r_s} \mathbf{G}_t (\mathbf{T}_* + \mathbf{N}_{t+1}) (\mathbf{I}_d - \mathbf{G}_t) \\ &\quad + \frac{\eta^2}{4r_s^2} \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top - \frac{\eta^2}{4r_s^2} \frac{\hat{\mathbf{W}}_{t+1} \mathcal{P}_{t+1} \hat{\mathbf{W}}_{t+1}^\top}{1 + c_{t+1}^2}.\end{aligned}$$

On the other hand,

$$\begin{aligned}\hat{\mathbf{W}}_{t+1} \mathcal{P}_{t+1} \hat{\mathbf{W}}_{t+1}^\top &= \left( \mathbf{W}_t + \frac{\eta}{2r_s} \nabla_{\text{St}} \mathbf{L}_{t+1} \right) \mathcal{P}_{t+1} \left( \mathbf{W}_t + \frac{\eta}{2r_s} \nabla_{\text{St}} \mathbf{L}_{t+1} \right)^\top \\ &= \mathbf{W}_t \mathcal{P}_{t+1} \mathbf{W}_t^\top + \frac{\eta}{r_s} \text{sym} \left( \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbf{W}_t^\top \right) + \frac{\eta^2}{4r_s^2} \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top.\end{aligned}$$

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

$$\begin{aligned}
 R_{\text{so}}[\mathbf{G}_t] &:= \frac{\eta^2}{4r_s^2} \mathbb{E}_t \left[ \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \right] - \frac{\eta^2}{4r_s^2} \mathbf{W}_t \mathbb{E}_t \left[ \frac{\mathcal{P}_{t+1}}{1+c_{t+1}^2} \right] \mathbf{W}_t^\top \\
 &\quad - \frac{\eta^3}{4r_s^3} \text{sym} \left( \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1+c_{t+1}^2} \right] \mathbf{W}_t^\top \right) - \frac{\eta^4}{16r_s^4} \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top}{1+c_{t+1}^2} \right]. \quad (\text{D.2})
 \end{aligned}$$

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

$$\begin{aligned}
 \frac{\eta}{r_s} \boldsymbol{\xi}_{t+1} &:= \frac{\eta}{r_s} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{N}_{t+1} \mathbf{G}_t + \frac{\eta}{r_s} \mathbf{G}_t \mathbf{N}_{t+1} (\mathbf{I}_d - \mathbf{G}_t) \\
 &\quad - \frac{\eta^2}{4r_s^2} \mathbf{W}_t \left( \frac{\mathcal{P}_{t+1}}{1+c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\mathcal{P}_{t+1}}{1+c_{t+1}^2} \right] \right) \mathbf{W}_t^\top \\
 &\quad + \frac{\eta^2}{4r_s^2} \left( \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top - \mathbb{E}_t \left[ \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \right] \right) \\
 &\quad - \frac{\eta^3}{4r_s^3} \text{sym} \left( \left( \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1+c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1+c_{t+1}^2} \right] \right) \mathbf{W}_t^\top \right) \\
 &\quad - \frac{\eta^4}{16r_s^4} \left( \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top}{1+c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top}{1+c_{t+1}^2} \right] \right).
 \end{aligned}$$

With these definitions in hand, we have

$$\mathbf{G}_{t+1} = \mathbf{G}_t + \frac{\eta}{r_s} (\mathbb{T}_\star \mathbf{G}_t + \mathbf{G}_t \mathbb{T}_\star - 2\mathbf{G}_t \mathbb{T}_\star \mathbf{G}_t) + R_{\text{so}}[\mathbf{G}_t] + \frac{\eta}{2r_s} \boldsymbol{\xi}_{t+1}.$$

**Proposition 29 (Including second order terms)** For  $\frac{\eta}{r_s} \ll \frac{1}{\sqrt{\text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})}}$ , by including the bounds in (D.5), we get

$$\mathbf{G}_{t+1} \preceq \mathbf{G}_t + \frac{\eta}{r_s} (\mathbb{T}_\star \mathbf{G}_t + \mathbf{G}_t \mathbb{T}_\star - 2\mathbf{G}_t \mathbb{T}_\star \mathbf{G}_t) + \frac{C\eta^2 \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})}{2r_s^2} (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) + \frac{\eta}{r_s} \boldsymbol{\xi}_{t+1} \quad (\text{D.3})$$

$$\begin{aligned}
 \mathbf{G}_{t+1} \succeq \mathbf{G}_t + \frac{\eta}{r_s} (\mathbb{T}_\star \mathbf{G}_t + \mathbf{G}_t \mathbb{T}_\star - 2\mathbf{G}_t (\mathbb{T}_\star + \frac{C\eta \text{tr}(\boldsymbol{\Sigma})}{r_s} \boldsymbol{\Sigma}) \mathbf{G}_t) \\
 + \frac{c\eta^2 \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})}{2r_s^2} (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) + \frac{\eta}{r_s} \boldsymbol{\xi}_{t+1} \quad (\text{D.4})
 \end{aligned}$$

**Proof** By using  $\text{tr}(\mathbf{G}_t \boldsymbol{\Sigma}) \leq r_s$ , the result is immediate use of the bound (D.5).  $\blacksquare$

**Proposition 30 (Sufficient statistics)** Let

$$\mathbb{T} =: \boldsymbol{\Theta} \boldsymbol{\Lambda} \boldsymbol{\Theta}^\top, \quad \boldsymbol{\Sigma}_1 = \boldsymbol{\Theta}^\top \boldsymbol{\Sigma} \boldsymbol{\Theta}, \quad \mathbf{M}_t := \boldsymbol{\Sigma}_1^{\frac{1}{2}} \mathbf{G}_t \boldsymbol{\Sigma}_1^{\frac{1}{2}}, \quad \boldsymbol{\zeta}_t := \boldsymbol{\Sigma}_1^{\frac{1}{2}} \boldsymbol{\xi}_t \boldsymbol{\Sigma}_1^{\frac{1}{2}}.$$

We have

$$\mathbf{M}_{t+1} \preceq \mathbf{M}_t + \frac{\eta}{r_s} (\boldsymbol{\Lambda} \boldsymbol{\Sigma}_1 \mathbf{M}_t + \mathbf{M}_t \boldsymbol{\Sigma}_1 \boldsymbol{\Lambda} - 2\mathbf{M}_t \boldsymbol{\Lambda} \mathbf{M}_t) + \frac{C\eta^2 \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})}{2r_s^2} (\boldsymbol{\Sigma}_1 - \mathbf{M}_t)^2 + \frac{\eta}{r_s} \boldsymbol{\zeta}_{t+1}$$

$$M_{t+1} \succeq M_t + \frac{\eta}{r_s} \left( \Lambda \Sigma_1 M_t + M_t \Sigma_1 \Lambda - 2M_t \Lambda M_t \right) + \frac{c\eta^2 \text{tr}(\mathbf{G}_t \Sigma)}{2r_s^2} (\Sigma_1 - M_t)^2 + \frac{\eta}{r_s} \zeta_{t+1}$$

**Proof** The proof immediately follows from (D.3) and (D.4). The additional term in the quadratic term of lower bound can be ignored since it does not contribute in asymptotic limit.  $\blacksquare$

We introduce block matrix notation:

$$M_t =: \begin{bmatrix} M_{t,11} & M_{t,12} \\ M_{t,12}^\top & M_{t,22} \end{bmatrix}, \quad \Lambda =: \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix}, \quad \zeta_t =: \begin{bmatrix} \zeta_{t,11} & \zeta_{t,12} \\ \zeta_{t,12}^\top & \zeta_{t,22} \end{bmatrix}, \quad \Sigma_1 =: \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{12} \end{bmatrix}$$

where  $M_{t,11}$ ,  $\Lambda_1$ ,  $\zeta_{t,11}$ ,  $\Sigma_{11} \in \mathbb{R}^{r_u \times r_u}$ .

**Proposition 31 (Sparsification)** *We define*

$$\begin{aligned} \Lambda_{u_1} &:= \Lambda_1 \Sigma_{11} - C\eta \text{tr}(\Sigma) (\Sigma_{11} - \|\Sigma_{12}\| \mathbf{I}_{r_u}), \quad \text{and} \quad \Lambda_{u_2} := \Lambda_1 - C\eta \text{tr}(\Sigma) \mathbf{I}_{r_u} \\ \Lambda_{\ell_1} &:= \Lambda_1 \Sigma_{11} - \|\Sigma_{12}\|_2 \|\Lambda_2\|_2 \mathbf{I}_{r_u}, \quad \text{and} \quad \Lambda_{\ell_2} := \Lambda_1, \end{aligned}$$

and

$$V_t^+ := 2\Lambda_{u_2}^{\frac{1}{2}} M_{t,11} \Lambda_{u_2}^{\frac{1}{2}} - \Lambda_{u_1} \quad \text{and} \quad V_t^- := 2\Lambda_{\ell_2}^{\frac{1}{2}} M_{t,11} \Lambda_{\ell_2}^{\frac{1}{2}} - \Lambda.$$

For  $\eta = \frac{\eta}{r_s}$ , we have

$$\begin{aligned} V_{t+1}^+ &\preceq V_t^+ \left( \mathbf{I}_{r_u} + \frac{\eta}{1+1.1\eta} V_t^+ \right)^{-1} + \eta \Lambda_{u_1}^2 + \eta^2 \text{tr}(\Sigma) \Lambda_{u_2} \Sigma_{11}^2 + 2\eta \Lambda_{u_2}^{\frac{1}{2}} \zeta_{t+1,11} \Lambda_{u_2}^{\frac{1}{2}} \\ V_{t+1}^- &\succeq V_t^- \left( \mathbf{I}_{r_u} + \frac{\eta}{1-1.1\eta} V_t^- \right)^{-1} + \eta \Lambda_{\ell_2}^2 + 2\eta \Lambda_{\ell_2}^{\frac{1}{2}} \zeta_{t+1,11} \Lambda_{\ell_2}^{\frac{1}{2}} \end{aligned}$$

**Proof** For the upper bound, we observe that

$$(\Sigma_1 - M_t)^2 = \Sigma_1^2 - \Sigma_1 M_t - M_t \Sigma_1 + M_t^2$$

Note that  $M_{t,12} = \Sigma_{11}^{\frac{1}{2}} \mathbf{G}_{t,12} \Sigma_{12}^{\frac{1}{2}}$ . Therefore,

$$M_{t,12} M_{t,12}^\top = \Sigma_{11}^{\frac{1}{2}} \mathbf{G}_{t,12} \Sigma_{12}^{\frac{1}{2}} \Sigma_{12}^{\frac{1}{2}} \mathbf{G}_{t,12}^\top \Sigma_{11}^{\frac{1}{2}} \preceq \|\Sigma_{12}\|_2 \Sigma_{11}^{\frac{1}{2}} \mathbf{G}_{t,12} \mathbf{G}_{t,12}^\top \Sigma_{11}^{\frac{1}{2}} \preceq \|\Sigma_{12}\|_2 M_{t,11}.$$

For the lower bound,

$$M_{t,12} \Lambda_2 M_{t,12}^\top = \Sigma_{11}^{\frac{1}{2}} \mathbf{G}_{t,12} \Sigma_{12}^{\frac{1}{2}} \Lambda_2 \Sigma_{12}^{\frac{1}{2}} \mathbf{G}_{t,12}^\top \Sigma_{11}^{\frac{1}{2}} \preceq \|\Sigma_{12}\|_2 \|\Lambda_2\|_2 M_{t,11}.$$

Therefore, the coefficients  $\Lambda_{u_1}$ ,  $\Lambda_{u_2}$ ,  $\Lambda_{\ell_1}$ ,  $\Lambda_{\ell_2}$  follow. The bounds can be derived by ordering the higher order terms in  $V \rightarrow V(\mathbf{I}_{r_u} + \eta V)^{-1}$   $\blacksquare$

### D.1. Definitions and bounding systems

Throughout the proof, we will also use a constant  $\kappa_d \in o_d(1)$  that will be specified later. Moreover, we make the following definitions:

- **Noise sequence.** For  $\nu_0 = 0$ , we define the noise sequence  $\nu_{t+1} := \nu_t + \eta \zeta_{t+1,11}$ .
- **Reference sequence.** For  $\mathbf{R}_0 = \frac{\kappa_d r_u}{d} \Sigma_{11}$ , we define the reference sequence

$$\mathbf{R}_{t+1} = \mathbf{R}_t + 2(1 - 2\kappa_d)\eta \left( \Lambda_{\ell_1} \mathbf{R}_t - \frac{3\kappa_d + 1}{\kappa_d(1 - 2\kappa_d)} \Lambda_{u_1} \mathbf{R}_t^2 \right).$$

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

$$\begin{aligned} \underline{\mathbf{V}}_{t+1} &= \underline{\mathbf{V}}_t \left( \mathbf{I}_{r_u} + \frac{\eta(1 + 2\kappa_d)}{1 - 1.2\eta} \underline{\mathbf{V}}_t \right)^{-1} + \frac{\eta(1 + 2\kappa_d)}{1 - 1.2\eta} \left( \frac{\Lambda_{\ell_1}^2}{(1 + 2\kappa_d)^2} - C(\eta \Lambda_{\ell_1}^2 + \eta^2 \Lambda_{\ell_1}) \right), \\ \bar{\mathbf{V}}_{t+1} &= \bar{\mathbf{V}}_t \left( \mathbf{I}_{r_u} + \frac{\eta(1 - 2\kappa_d)}{1 + 1.2\eta} \bar{\mathbf{V}}_t \right)^{-1} + \frac{\eta(1 - 2\kappa_d)}{1 + 1.2\eta} \left( \frac{\Lambda_{u_1}^2}{(1 - 2\kappa_d)^2} + C(\eta \text{tr}(\Sigma) \Lambda_{u_2} \Sigma_{11}^2 + \eta \Lambda_{u_2}^2 + \eta^2 \Lambda_{u_2}) \right). \end{aligned}$$

where the iterates  $\{\underline{\mathbf{V}}_{t+1}\}_{t \in \mathbb{N}}$  and  $\{\bar{\mathbf{V}}_{t+1}\}_{t \in \mathbb{N}}$  are functions of the bounding sequences  $\{\underline{\mathbf{M}}_t\}_{t \in \mathbb{N}}$  and  $\{\bar{\mathbf{M}}_t\}_{t \in \mathbb{N}}$  as following:

$$\begin{aligned} \underline{\mathbf{V}}_t &= 2\Lambda_{u_2}^{\frac{1}{2}} \underline{\mathbf{M}}_t \Lambda_{u_2}^{\frac{1}{2}} - \frac{\Lambda_{u_1}}{1 + 2\kappa_d} \quad \text{and} \quad \underline{\mathbf{M}}_0 \preceq \mathbf{M}_{0,11} - \mathbf{R}_0, \\ \bar{\mathbf{V}}_t &= 2\Lambda_{\ell_2}^{\frac{1}{2}} \bar{\mathbf{M}}_t \Lambda_{\ell_2}^{\frac{1}{2}} - \frac{\Lambda_{\ell_1}}{1 - 2\kappa_d} \quad \text{and} \quad \bar{\mathbf{M}}_0 \succeq \mathbf{M}_{0,11} + \mathbf{R}_0. \end{aligned}$$

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

$$\mathcal{E}_t := \left\{ -\kappa_d r_u^{\frac{-\alpha-\beta}{2}} \mathbf{R}_t \preceq \nu_t \preceq \kappa_d r_u^{\frac{-\alpha-\beta}{2}} \mathbf{R}_t \right\} \cap \left\{ -\frac{\kappa_d^2}{4} r_u^{\frac{-\alpha-\beta}{2}} \mathbf{R}_t \preceq \Sigma_{11}^{\frac{1}{2}} \Lambda_1^{\frac{1}{2}} \nu_t \Lambda_1^{\frac{1}{2}} \Sigma_{11}^{\frac{1}{2}} \preceq \frac{\kappa_d^2}{4} r_u^{\frac{-\alpha-\beta}{2}} \mathbf{R}_t \right\}.$$

We define the stopping times

$$\mathcal{T}_{\text{noise}}(\omega) := \inf \{t \geq 0 \mid \omega \notin \mathcal{E}_t\} \wedge d^3 \quad \text{and} \quad \mathcal{T}_{\text{bounded}} := \inf \{t \geq 0 \mid \|\Sigma_{11}^{-1} \underline{\mathbf{M}}_t\|_2 \vee \|\Sigma_{11}^{-1} \bar{\mathbf{M}}_t\|_2 > 1.5\},$$

and

$$\mathcal{T}_{\text{bad}} := \mathcal{T}_{\text{noise}} \wedge \mathcal{T}_{\text{bounded}} \wedge \{t \geq 0 : \|\Sigma_{11}^{-1} \mathbf{R}_t\|_2 > 1.2\kappa_d\}.$$

We start with the following lemma:

**Lemma 32** *We consider  $\kappa_d \ll \frac{r_u^{-1}}{\log d}$ . The event  $\mathcal{E}_t$  implies for  $d \geq \Omega(1)$  and  $t \leq \mathcal{T}_{\text{bounded}} \wedge \{t : \|\Sigma_{11}^{-1} \mathbf{R}_t\|_2 > 1.2\kappa_d\}$  that*

$$1. \quad -3\kappa_d \Lambda_{\ell_1}^{\frac{1}{2}} \mathbf{R}_t \Lambda_{\ell_1}^{\frac{1}{2}} \preceq \Lambda_{\ell_1} \nu_t + \nu_t \Lambda_{\ell_1}$$

2.  $\Lambda_{u_1} \boldsymbol{\nu}_t + \boldsymbol{\nu}_t \Lambda_{u_1} \preceq 3\kappa_d \Lambda_{u_1}^{\frac{1}{2}} \mathbf{R}_t \Lambda_{u_1}^{\frac{1}{2}}$
3.  $(\Lambda_{\ell_2}^{\frac{1}{2}} \boldsymbol{\nu}_t \Lambda_{\ell_2}^{\frac{1}{2}})^2 \preceq \frac{\kappa_d^2}{4} \Lambda_{\ell_1} \mathbf{R}_t \Lambda_{\ell_1}$
4.  $(\Lambda_{u_2}^{\frac{1}{2}} \boldsymbol{\nu}_t \Lambda_{u_2}^{\frac{1}{2}})^2 \preceq \frac{\kappa_d^2}{4} \Lambda_{u_1} \mathbf{R}_t \Lambda_{u_1}$

**Proof** For notational convenience, we define  $\tilde{\boldsymbol{\nu}}_t := \mathbf{R}_t^{-\frac{1}{2}} \boldsymbol{\nu}_t \mathbf{R}_t^{-\frac{1}{2}}$ . We observe that

$$\begin{aligned} \Sigma_{11} \Lambda_1 \tilde{\boldsymbol{\nu}}_t + \tilde{\boldsymbol{\nu}}_t \Lambda_1 \Sigma_{11} &\preceq \frac{\kappa_d^2}{4} \Lambda_1 \Sigma_{11} + \frac{4}{\kappa_d^2} \Sigma_{11}^{-\frac{1}{2}} \Lambda_1^{-\frac{1}{2}} \left( \Sigma_{11}^{\frac{1}{2}} \Lambda_1^{\frac{1}{2}} \tilde{\boldsymbol{\nu}}_t \Lambda_1^{\frac{1}{2}} \Sigma_{11}^{\frac{1}{2}} \right)^2 \Lambda_1^{-\frac{1}{2}} \Sigma_{11}^{-\frac{1}{2}} \\ &\preceq \frac{\kappa_d^2}{4} \Lambda_1 \Sigma_{11} + \frac{\kappa_d^2}{4} r_u^{-\alpha-\beta} \mathbf{I}_{r_u} \\ \Sigma_{11} \Lambda_1 \tilde{\boldsymbol{\nu}}_t + \tilde{\boldsymbol{\nu}}_t \Lambda_1 \Sigma_{11} &\succeq -\frac{\kappa_d^2}{4} \Lambda_1 \Sigma_{11} - \frac{4}{\kappa_d^2} \Sigma_{11}^{-\frac{1}{2}} \Lambda_1^{-\frac{1}{2}} \left( \Sigma_{11}^{\frac{1}{2}} \Lambda_1^{\frac{1}{2}} \tilde{\boldsymbol{\nu}}_t \Lambda_1^{\frac{1}{2}} \Sigma_{11}^{\frac{1}{2}} \right)^2 \Lambda_1^{-\frac{1}{2}} \Sigma_{11}^{-\frac{1}{2}} \\ &\succeq -\frac{\kappa_d^2}{4} \Lambda_1 \Sigma_{11} - \frac{\kappa_d^2}{4} r_u^{-\alpha-\beta} \mathbf{I}_{r_u} \end{aligned}$$

The items follows the use of matrix Young's inequality with the given bounds.  $\blacksquare$

The main result of this section is the following:

**Proposition 33** Consider  $d$  large enough so that  $\kappa_d \leq \frac{1}{50}$  and  $\eta$  small enough to make sure that

$$\left(1 - \frac{C}{\log d}\right) \mathbf{I}_{r_u} \preceq \Lambda_{u_1} \Lambda_{u_2}^{-1}, \Lambda_{\ell_1} \Lambda_{\ell_2}^{-1} \preceq \left(1 + \frac{C}{\log d}\right) \mathbf{I}_{r_u}.$$

We have

$$\underline{\mathbf{M}}_{t \wedge \mathcal{T}_{bad}} + \mathbf{R}_{t \wedge \mathcal{T}_{bad}} + \boldsymbol{\nu}_{t \wedge \mathcal{T}_{bad}} \preceq \mathbf{M}_{t \wedge \mathcal{T}_{bad}, 11} \preceq \bar{\mathbf{M}}_{t \wedge \mathcal{T}_{bad}} - \mathbf{R}_{t \wedge \mathcal{T}_{bad}} + \boldsymbol{\nu}_{t \wedge \mathcal{T}_{bad}}.$$

**Proof** The proof is identical with [3, Proposition 15].  $\blacksquare$

## D.2. Bounding the second order terms

**Proposition 34** There exists a universal constant  $C > 0$  such that

$$R_{so}[\mathbf{G}_t] \sim \frac{\eta^2 \text{tr}(\boldsymbol{\Sigma})}{r_s^2} \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t + \left( \frac{\eta^2 \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})}{r_s^2} + \frac{\eta^4 \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})^2}{r_s^4} \right) (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t). \quad (\text{D.5})$$

**Proof** We will bound each term in (D.2) separately. For the first term,

$$\mathbb{E}_t \left[ \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \right] = (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t).$$

Therefore,

$$0 \preceq \frac{\eta^2}{4r_s^2} \mathbb{E}_t \left[ \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \right] \preceq \frac{C\eta^2}{r_s^2} \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma}) (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t).$$

For the second term, we have

$$\mathbf{W}_t \mathbb{E}_t \left[ \frac{\mathcal{P}_{t+1}}{1 + c_{t+1}^2} \right] \mathbf{W}_t^\top = \mathbf{G}_t \mathbb{E}_t \left[ \frac{y_{t+1}^2 \|( \mathbf{I}_d - \mathbf{G}_t ) \mathbf{x}_{t+1} \|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top}{1 + c_{t+1}^2} \right] \mathbf{G}_t.$$

Therefore,

$$0 \preceq \frac{\eta^2}{4r_s^2} \mathbf{W}_t \mathbb{E}_t \left[ \frac{\mathcal{P}_{t+1}}{1 + c_{t+1}^2} \right] \mathbf{W}_t^\top \preceq C \frac{\eta^2 \text{tr}(\boldsymbol{\Sigma})}{r_s^2} \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t.$$

For the third term, we have

$$\mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1 + c_{t+1}^2} \right] \mathbf{W}_t^\top = (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ \frac{y_{t+1}^3}{1 + c_{t+1}^2} \|( \mathbf{I}_d - \mathbf{G}_t ) \mathbf{x}_{t+1} \|_2^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] \mathbf{G}_t.$$

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

$$\left\| \mathbb{E}_t \left[ \frac{y_{t+1}^3}{1 + c_{t+1}^2} \|( \mathbf{I}_d - \mathbf{G}_t ) \mathbf{x}_{t+1} \|_2^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2 \boldsymbol{\Sigma}^{-\frac{1}{2}} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \boldsymbol{\Sigma}^{-\frac{1}{2}} \right] \right\|_2 \leq C \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma}).$$

Therefore,

$$\frac{\eta^3}{4r_s^3} \text{sym} \left( \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1 + c_{t+1}^2} \right] \mathbf{W}_t^\top \right) \preceq C \left( \frac{\eta^4 \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})^2}{r_s^4} (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) + \frac{\eta^2 \text{tr}(\boldsymbol{\Sigma})}{r_s^2} \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \right)$$

By repeating the argument for the lower bound,

$$\frac{\eta^3}{4r_s^3} \text{sym} \left( \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1}}{1 + c_{t+1}^2} \right] \mathbf{W}_t^\top \right) \succeq -C \left( \frac{\eta^4 \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})^2}{r_s^4} (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) + \frac{\eta^2 \text{tr}(\boldsymbol{\Sigma})}{r_s^2} \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \right).$$

For the last term, we write

$$\mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] = (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ \frac{y_{t+1}^4}{1 + c_{t+1}^2} \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|( \mathbf{I}_d - \mathbf{G}_t ) \mathbf{x}_{t+1} \|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t).$$

We have

$$\left\| \mathbb{E}_t \left[ \frac{y_{t+1}^4}{1 + c_{t+1}^2} \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|( \mathbf{I}_d - \mathbf{G}_t ) \mathbf{x}_{t+1} \|_2^2 \boldsymbol{\Sigma}^{-\frac{1}{2}} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \boldsymbol{\Sigma}^{-\frac{1}{2}} \right] \right\|_2 \leq C \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})^2.$$

Therefore,

$$0 \preceq \frac{\eta^4}{16r_s^4} \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top}{1 + c_{t+1}^2} \right] \preceq C \frac{\eta^4 \text{tr}(\boldsymbol{\Sigma}) \text{tr}(\mathbf{G}_t \boldsymbol{\Sigma})^2}{r_s^4} (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t).$$

■

### D.3. Noise characterization

For  $\mathbf{K}_1, \mathbf{K}_2 \in \mathbb{R}^{d \times m}$ , we define

$$\mathcal{A}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \equiv \left\{ \left\| \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{N}_{t+1} \mathbf{G}_t \mathbf{K}_2 \right\|_2 \leq \frac{L^2}{2} \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right\}$$

$$\mathcal{B}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \equiv \left\{ \left\| \mathbf{K}_1^\top \mathbf{W}_t \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 \right\|_2 \leq \frac{L^4 \text{tr}(\boldsymbol{\Sigma})}{2} \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right\}$$

$$\mathcal{C}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \equiv \left\{ \left\| \mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbf{K}_2 \right\|_2 \leq \frac{L^4 \text{tr}(\boldsymbol{\Sigma})}{2} \sqrt{\prod_{i=1}^2 \text{tr}(\mathbf{K}_i \mathbf{K}_i^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right\}$$

$$\mathcal{D}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \equiv \left\{ \left\| \mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 \right\|_2 \leq \frac{L^6 \text{tr}(\boldsymbol{\Sigma})^2}{2} \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right\}$$

$$\mathcal{F}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \equiv \left\{ \left\| \mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbf{K}_2 \right\|_2 \leq \frac{L^8 \text{tr}(\boldsymbol{\Sigma})^3}{2} \sqrt{\prod_{i=1}^2 \text{tr}(\mathbf{K}_i \mathbf{K}_i^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right\}.$$

We start with the following statement:

**Proposition 35** *There exists a universal constant  $C > 0$  such that for  $L \geq e(\sqrt{8} + \sigma)$ , the following holds:*

1. We have  $\mathbb{P}_t[\mathcal{A}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \cap \{|y_{t+1}| \leq L\}] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8} + \sigma)}}$ . Moreover,

$$\begin{aligned} & \mathbb{E}_t \left[ \left( \text{sym} \left( \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{N}_{t+1} \mathbf{G}_t \mathbf{K}_2 \right) \right)^2 \right] \\ & \leq C \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\ & + C \text{tr} \left( \mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \right) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2. \end{aligned}$$

2. We have  $\mathbb{P}_t[\mathcal{B}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \cap \{|y_{t+1}| \leq L\}] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8} + \sigma)}}$ . Moreover,

$$\begin{aligned} & \mathbb{E}_t \left[ \left( \text{sym} \left( \mathbf{K}_1^\top \mathbf{W}_t \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 \right) \right)^2 \right] \leq C \text{tr}(\boldsymbol{\Sigma})^2 (\text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_1 \\ & + C \text{tr}(\boldsymbol{\Sigma})^2 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2. \end{aligned}$$

3. We have  $\mathbb{P}_t[\mathcal{C}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \cap \{|y_{t+1}| \leq L\}] \geq 1 - 2e^{\frac{-L/e}{(\sqrt{8} + \sigma)}}$ . Moreover,

$$\begin{aligned} & \mathbb{E}_t \left[ \left( \text{sym} \left( \mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbf{K}_2 \right) \right)^2 \right] \\ & \leq C \text{tr}(\boldsymbol{\Sigma})^2 \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\ & + C \text{tr}(\boldsymbol{\Sigma})^2 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2 \end{aligned}$$

4. We have  $\mathbb{P}_t[\mathcal{D}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \cap \{|y_{t+1}| \leq L\}] \geq 1 - 2e^{\frac{-L/e}{(7/2 + \sigma)}}$ . Moreover,

$$\begin{aligned} & \mathbb{E}_t \left[ \left( \text{sym} \left( \mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 \right) \right)^2 \right] \leq C \text{tr}(\boldsymbol{\Sigma})^4 \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\ & + C \text{tr}(\boldsymbol{\Sigma})^4 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2. \end{aligned}$$

5. We have  $\mathbb{P}_t[\mathcal{F}_{t+1}(\mathbf{K}_1, \mathbf{K}_2) \cap \{|y_{t+1}| \leq L\}] \geq 1 - 2e^{\frac{-L/e}{(4\sqrt{2}+\sigma)}}$ . Moreover,

$$\begin{aligned} & \mathbb{E}_t \left[ \left( \text{sym}(\mathbf{K}_1^\top \nabla_{S_t} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{S_t} \mathbf{L}_{t+1}^\top \mathbf{K}_2) \right)^2 \right] \\ & \leq C \text{tr}(\boldsymbol{\Sigma})^6 \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\ & \quad + C \text{tr}(\boldsymbol{\Sigma})^6 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2. \end{aligned}$$

**Proof** First, we derive a concentration bound for  $|y_{t+1}|$ . We have

$$\mathbb{E}_t[|y_{t+1}|^p] \leq (\sqrt{2} + \sigma/2)^p p^p \text{ for } p \geq 2,$$

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

**First item.** We define

$$\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{N}_{t+1} \mathbf{G}_t \mathbf{K}_2 = \underbrace{\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) y_{t+1} \mathbf{x}_{t+1}}_{:= \mathbf{u}_{t+1}} \underbrace{\mathbf{x}_{t+1}^\top \mathbf{G}_t \mathbf{K}_2}_{:= \mathbf{v}_{t+1}^\top}.$$

For  $u, L > 0$

$$\begin{aligned} & \mathbb{P}_t \left[ \left\| \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq uL \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \text{ or } |y_{t+1}| \geq L \right] \\ & \leq \mathbb{P}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq uL \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right] \\ & \quad + \mathbb{P}_t[|y_{t+1}| \geq L]. \end{aligned}$$

We have for  $p \geq 2$

$$\begin{aligned} \mathbb{E}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2^p \right] & \leq L^p \mathbb{E}_t \left[ \left\| \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^p \right] \mathbb{E}_t \left[ \left\| \mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1} \right\|_2^p \right] \\ & \stackrel{(a)}{\leq} L^p \left( \frac{p}{2} \right)^p \left( 3 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \right)^{\frac{p}{2}}. \end{aligned}$$

We have for  $u \geq 2e$

$$\mathbb{P}_t \left[ \left\| \mathbb{1}\{|e_{t+1}| \leq t\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq uL \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right] \leq e^{-\frac{u}{e}}.$$

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

For the variance bound, we have

$$\begin{aligned} \mathbb{E}_t \left[ \text{sym} \left( \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right)^2 \right] & \leq \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^2 \left\| \mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1} \right\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1^\top \\ & \quad + \mathbf{K}_2^\top \mathbf{G}_t \mathbb{E}_t \left[ y_{t+1}^2 \left\| \mathbf{K}_1 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] \mathbf{G}_t \mathbf{K}_2^\top \\ & \stackrel{(b)}{\leq} C \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\ & \quad + C \text{tr} \left( \mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \right) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2, \end{aligned}$$

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

**Second item.** We define

$$\mathbf{K}_1^\top \mathbf{W}_t \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 = \underbrace{y_{t+1}^2 \|(\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2}_{:=\mathbf{u}_{t+1}} \underbrace{\mathbf{K}_1^\top \mathbf{G}_t \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{G}_t \mathbf{K}_2}_{:=\mathbf{v}_{t+1}^\top}.$$

We have for  $p \geq 2$

$$\begin{aligned} & \mathbb{E}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2^p \right] \\ & \leq L^{2p} \mathbb{E}_t \left[ \|(\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^{2p} \right] \mathbb{E}_t \left[ \|\mathbf{K}_1^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^{2p} \right]^{\frac{1}{2}} \mathbb{E}_t \left[ \|\mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^{2p} \right]^{\frac{1}{2}} \\ & \stackrel{(c)}{\leq} L^{2p} p^{2p} \left( 3 \operatorname{tr}(\boldsymbol{\Sigma}) \sqrt{\operatorname{tr}(\mathbf{K}_1 \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \operatorname{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right)^p. \end{aligned}$$

We have for  $u \geq (2e)^2$

$$\mathbb{P}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq u L^2 \operatorname{tr}(\boldsymbol{\Sigma}) \sqrt{\operatorname{tr}(\mathbf{K}_1 \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \operatorname{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right] \leq e^{-\frac{u^{1/2}}{e}}.$$

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

$$\begin{aligned} \mathbb{E}_t \left[ \operatorname{sym} \left( \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right)^2 \right] & \preceq \mathbf{K}_1^\top \mathbf{G}_t \mathbb{E}_t \left[ y_{t+1}^4 \|(\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] \mathbf{G}_t \mathbf{K}_1 \\ & \quad + \mathbf{K}_2^\top \mathbf{G}_t \mathbb{E}_t \left[ y_{t+1}^4 \|(\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_1^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] \mathbf{G}_t \mathbf{K}_2 \\ & \stackrel{(d)}{\preceq} C \operatorname{tr}(\boldsymbol{\Sigma})^2 \left( \operatorname{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_1 + \operatorname{tr}(\mathbf{K}_1 \mathbf{K}_1^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2 \right), \end{aligned}$$

where we use the Cauchy-Schwartz inequality in (d).

**Third item.** We define

$$\mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbf{K}_2 = \underbrace{y_{t+1}^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2}_{:=\mathbf{u}_{t+1}} \underbrace{\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top (\mathbf{I}_d - \mathbf{G}_t)}_{:=\mathbf{v}_{t+1}^\top} \mathbf{K}_2.$$

We have for  $p \geq 2$

$$\begin{aligned} & \mathbb{E}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2^p \right] \\ & \leq L^{2p} \mathbb{E}_t \left[ \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^{2p} \right] \mathbb{E}_t \left[ \|\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^{2p} \right]^{\frac{1}{2}} \mathbb{E}_t \left[ \|\mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^{2p} \right]^{\frac{1}{2}} \\ & \stackrel{(e)}{\leq} L^{2p} p^{2p} \left( 3r_s \sqrt{\operatorname{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \operatorname{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right)^p, \end{aligned}$$

We have for  $u \geq (2e)^2$

$$\begin{aligned} & \mathbb{P}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \right. \\ & \quad \left. \geq u L^2 3r_s \sqrt{\operatorname{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \operatorname{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right] \leq e^{-\frac{u^{1/2}}{e}}. \end{aligned}$$

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

$$\begin{aligned}
 & \mathbb{E}_t \left[ \text{sym} \left( \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right)^2 \right] \\
 & \preceq \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^4 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & + \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^4 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_1 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2 \\
 & \stackrel{(f)}{\preceq} C r_s^2 \text{tr} \left( \mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \right) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & + C r_s^2 \text{tr} \left( \mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \right) \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2,
 \end{aligned}$$

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

**Fourth item.** We define

$$\mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbf{W}_t^\top \mathbf{K}_2 = \underbrace{y_{t+1}^3 \|\mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2}_{:= \mathbf{u}_{t+1}} \underbrace{\|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^2 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{G}_t \mathbf{K}_2}_{:= \mathbf{v}_{t+1}^\top}.$$

We have for  $p \geq 2$

$$\begin{aligned}
 & \mathbb{E}_t \left[ \left\| \mathbb{1}_{\{|y_{t+1}| \leq L\}} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2^p \right] \\
 & \leq L^{3p} \mathbb{E}_t \left[ \left\| (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^{2p} \|\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^p \right] \mathbb{E}_t \left[ \left\| \mathbf{W}_t^\top \mathbf{x}_{t+1} \right\|_2^{2p} \|\mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^p \right] \\
 & \leq L^{3p} (12\sqrt{3})^p p^{3p} \left( \text{tr}(\boldsymbol{\Sigma}) r_s \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right)^p.
 \end{aligned}$$

We have for  $u \geq (2e)^3$

$$\begin{aligned}
 \mathbb{P}_t \left[ \left\| \mathbb{1}_{\{|t_{t+1}| \leq L\}} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq u L^3 12\sqrt{3} \text{tr}(\boldsymbol{\Sigma}) r_s \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t)} \right] \\
 \leq e^{-\frac{u^{1/3}}{e}}.
 \end{aligned}$$

By choosing  $u = \frac{L^3}{24\sqrt{3}}$ , we have the probability bound. For the variance bound, we have

$$\begin{aligned}
 & \mathbb{E}_t \left[ \text{sym} \left( \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right)^2 \right] \\
 & \preceq \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^6 \|\mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^4 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_2^\top \mathbf{G}_t \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & + \mathbf{K}_2^\top \mathbf{G}_t \mathbb{E}_t \left[ y_{t+1}^6 \|\mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^4 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4 \|\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] \mathbf{G}_t \mathbf{K}_2 \\
 & \preceq C \text{tr}(\boldsymbol{\Sigma})^2 r_s^2 \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & + C \text{tr}(\boldsymbol{\Sigma})^2 r_s^2 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_2^\top \mathbf{G}_t \boldsymbol{\Sigma} \mathbf{G}_t \mathbf{K}_2.
 \end{aligned}$$

**Fifth item.** We define

$$\mathbf{K}_1^\top \nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbf{K}_2 = \underbrace{y_{t+1}^4 \|\mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1}\|_2^2 \|\mathbf{W}_t^\top \mathbf{x}_{t+1}\|_2^4}_{:= \mathbf{u}_{t+1}} \underbrace{\mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2}_{:= \mathbf{v}_{t+1}^\top}.$$

We have for  $p \geq 2$

$$\begin{aligned}
 & \mathbb{E}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2^p \right] \\
 & \leq L^{4p} \mathbb{E}_t \left[ \left\| (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^{2p} \left\| \mathbf{W}_t^\top \mathbf{x}_{t+1} \right\|_2^{4p} \left\| \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^p \left\| \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^p \right] \\
 & \leq L^{4p} \mathbb{E}_t \left[ \left\| (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^{4p} \right]^{\frac{1}{2}} \mathbb{E}_t \left[ \left\| \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^{4p} \right]^{\frac{1}{4}} \mathbb{E}_t \left[ \left\| \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^{4p} \right]^{\frac{1}{4}} \mathbb{E}_t \left[ \left\| \mathbf{W}_t^\top \mathbf{x}_{t+1} \right\|_2^{4p} \right] \\
 & \leq L^{4p} (16p^4)^p (3\sqrt{3}r_s^2 \text{tr}(\boldsymbol{\Sigma}))^p \left( 3 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \right)^{\frac{p}{2}} \\
 & = L^{4p} (2\sqrt{3})^{4p} p^{4p} \left( \text{tr}(\boldsymbol{\Sigma}) r_s^2 \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right)^p.
 \end{aligned}$$

We have for  $u \geq (2e)^4$

$$\begin{aligned}
 & \mathbb{P}_t \left[ \left\| \mathbb{1}\{|y_{t+1}| \leq L\} \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right\|_2 \geq \right. \\
 & \quad \left. u (2\sqrt{3}L)^4 \text{tr}(\boldsymbol{\Sigma}) r_s^2 \sqrt{\text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t))} \right] \leq e^{-\frac{u^{1/4}}{e}}.
 \end{aligned}$$

By choosing  $u = \frac{L^4}{2(2\sqrt{3})^4}$ , we have the probability bound. For the variance bound, we have

$$\begin{aligned}
 & \mathbb{E}_t \left[ \text{sym} \left( \mathbf{u}_{t+1} \mathbf{v}_{t+1}^\top \right)^2 \right] \\
 & \leq \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^8 \left\| (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^4 \left\| \mathbf{W}_t^\top \mathbf{x}_{t+1} \right\|_2^8 \left\| \mathbf{K}_2 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & \quad + \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \mathbb{E}_t \left[ y_{t+1}^8 \left\| (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^4 \left\| \mathbf{W}_t^\top \mathbf{x}_{t+1} \right\|_2^8 \left\| \mathbf{K}_1 (\mathbf{I}_d - \mathbf{G}_t) \mathbf{x}_{t+1} \right\|_2^2 \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \right] (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2 \\
 & \leq C \text{tr}(\boldsymbol{\Sigma})^2 r_s^4 \text{tr}(\mathbf{K}_2 \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_1 \\
 & \quad + C \text{tr}(\boldsymbol{\Sigma})^2 r_s^4 \text{tr}(\mathbf{K}_1 \mathbf{K}_1^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t)) \mathbf{K}_2^\top (\mathbf{I}_d - \mathbf{G}_t) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{G}_t) \mathbf{K}_2.
 \end{aligned}$$

■

### D.3.1. APPLICATION OF NOISE BOUNDS

We will use  $r_u = \log^{\frac{1}{\alpha+\beta}} d$ . For  $\eta \ll \frac{r_s}{d \text{tr}(\boldsymbol{\Sigma})}$  and  $r_s \gg \frac{d}{\text{polylog } d}$ , we can show that  $\underline{\mathbf{M}}_t$  and  $\bar{\mathbf{M}}_t$  have asymptotically same dynamics with

$$\mathbf{M}_{t+1} = \mathbf{M}_t + 2\eta(\boldsymbol{\Lambda}_1 \boldsymbol{\Sigma}_{11} \mathbf{M}_t - \boldsymbol{\Lambda}_1 \mathbf{M}_t^2), \quad \text{with } \mathbf{M}_0 = \frac{r_s}{d} \boldsymbol{\Sigma}_{11}.$$

We start with the following proposition:

**Proposition 36** *Under  $r_u = \log^{\frac{1}{\alpha+\beta}} d$  and  $r_s \gg \frac{d}{\text{polylog } d}$ , we have*

- $\mathbf{M}_t \mathbf{R}_t^{-1} = (1 + o_d(1)) \kappa_d$
- $(\eta \sum_{t=0}^{t-1} \mathbf{M}_s) \mathbf{M}_t^{-1} \preceq \min\{\eta t \mathbf{I}_{r_u}, (2\boldsymbol{\Lambda}_1 (\boldsymbol{\Sigma}_{11} - \mathbf{M}_t))^{-1}\}$

By recalling the definitions  $\{\mathbf{R}_t\}_{t \in \mathbb{N}}$ ,  $\mathbf{\Lambda}_1$ ,  $\mathbf{\Sigma}_{11}$ , we define the event:

$$\mathcal{A}_{t+1} := \mathcal{A}_{t+1} \left( \mathbf{\Sigma}_{11}^{\frac{1}{2}} \mathbf{R}_t^{-\frac{1}{2}}, \mathbf{\Sigma}_{11}^{\frac{1}{2}} \mathbf{R}_t^{-\frac{1}{2}} \right) \cap \mathcal{A}_{t+1} \left( \mathbf{\Sigma}_{11} \mathbf{\Lambda}_1^{\frac{1}{2}} \mathbf{R}_t^{-\frac{1}{2}}, \mathbf{R}_t^{-\frac{1}{2}} \mathbf{\Lambda}_1^{\frac{1}{2}} \mathbf{\Sigma}_{11} \right) \cap \{|y_{t+1}| \leq L\}$$

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:

$$\begin{aligned} \mathbf{A}_{t+1} &:= \text{sym} \left( (\mathbf{I}_d - \mathbf{G}_t) \left( y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbb{1}\{\mathcal{A}_{t+1}\} - \mathbb{E}_t [y_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbb{1}\{\mathcal{A}_{t+1}\}] \right) \mathbf{G}_t \right) \\ \mathbf{B}_{t+1} &:= \mathbf{W}_t \left( \frac{\mathcal{P}_{t+1} \mathbb{1}\{\mathcal{B}_{t+1}\}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\mathcal{P}_{t+1} \mathbb{1}\{\mathcal{B}_{t+1}\}}{1 + c_{t+1}^2} \right] \right) \mathbf{W}_t^\top \\ \mathbf{C}_{t+1} &:= \left( \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbb{1}\{\mathcal{C}_{t+1}\} - \mathbb{E}_t \left[ \nabla_{\text{St}} \mathbf{L}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbb{1}\{\mathcal{C}_{t+1}\} \right] \right) \\ \mathbf{D}_{t+1} &:= \text{sym} \left( \left( \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbb{1}\{\mathcal{D}_{t+1}\}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \mathbb{1}\{\mathcal{D}_{t+1}\}}{1 + c_{t+1}^2} \right] \right) \mathbf{W}_t^\top \right) \\ \mathbf{F}_{t+1} &:= \left( \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbb{1}\{\mathcal{F}_{t+1}\}}{1 + c_{t+1}^2} - \mathbb{E}_t \left[ \frac{\nabla_{\text{St}} \mathbf{L}_{t+1} \mathcal{P}_{t+1} \nabla_{\text{St}} \mathbf{L}_{t+1}^\top \mathbb{1}\{\mathcal{F}_{t+1}\}}{1 + c_{t+1}^2} \right] \right). \end{aligned}$$

Let  $\mathbf{X} \in \left\{ \frac{\eta}{r_s} \mathbf{A}, \frac{\eta^2}{4r_s^2} \mathbf{B}, \frac{\eta^2}{4r_s^2} \mathbf{C}, \frac{\eta^3}{4r_s^3} \mathbf{D}, \frac{\eta^4}{16r_s^4} \mathbf{F} \right\}$  and

$$\mathbf{\Gamma}_1 := \mathbf{I}_{r_u}, \quad \mathbf{\Gamma}_2 := \mathbf{\Lambda}_1^{\frac{1}{2}} \mathbf{\Sigma}_{11}^{\frac{1}{2}}$$

For  $\ell \in \{1, 2\}$ , we define:

$$\text{Quad}_{k,t}^{(\ell)}(\mathbf{X}) := \sum_{j=1}^k \mathbb{E}_{j-1} \left[ \left( \mathbf{\Gamma}_\ell \mathbf{\Sigma}_{11}^{\frac{1}{2}} \mathbf{R}_t^{-\frac{1}{2}} \mathbf{X}_j \mathbf{R}_t^{-\frac{1}{2}} \mathbf{\Sigma}_{11}^{\frac{1}{2}} \mathbf{\Gamma}_\ell \right)^2 \right].$$

We have the following corollary.

**Corollary 37** *Let*

$$r_u = \log^{\frac{1}{\alpha+\beta}} d, \quad \eta \ll \frac{r_s}{d \text{tr}(\mathbf{\Sigma})}, \quad r_s \gg \frac{d}{\text{polylog } d}.$$

For  $\eta t \leq \frac{1}{2} \log d$ , we have:

$$\|\text{Quad}_{t,t}^{(\ell)}(\mathbf{X})\|_2 \leq \frac{C \log d}{\kappa_d^2} \frac{dr_u}{r_s} \begin{cases} \eta & \mathbf{X} = \frac{\eta}{r_s} \mathbf{A} \\ \eta^3 \text{tr}(\mathbf{\Sigma})^2, & \mathbf{X} = \frac{\eta^2}{4r_s^2} \mathbf{B} \\ \eta^3 \frac{d \log d \text{tr}(\mathbf{\Sigma})^2}{r_s}, & \mathbf{X} = \frac{\eta^2}{4r_s^2} \mathbf{C} \\ \eta^5 \text{tr}(\mathbf{\Sigma})^4, & \mathbf{X} = \frac{\eta^3}{16r_s^3} \mathbf{D} \\ \eta^7 \frac{d \log \text{tr}(\mathbf{\Sigma})^6}{r_s}, & \mathbf{X} = \frac{\eta^4}{16r_s^4} \mathbf{F}. \end{cases}$$

**Proof** By using the bounds in Proposition 35, the result can be proven. ■

By using matrix Bernstein inequality similar to [3, Propositions 21 and 22], we can derive the result.

#### D.4. Alignment step

As detailed in [3, Section E], by using the parameterization  $\mathbf{M} = \mathbf{\Omega}\mathbf{\Omega}^\top$ , we can reduce the alignment step to a linear regression problem over the space of  $r_s \times r_s$  symmetric matrices, where the  $\ell_2$  regularization on  $\mathbf{\Omega}$  translates to nuclear norm regularization on  $\mathbf{M}$ . The problem then reduces to finding the best sparse approximation of the labels via  $\ell_1$  regularization. Because  $\mathbf{W}_t$  aligns with the top eigenvector of  $\mathbf{\Sigma}$ , we can show that with an appropriate regularization parameter,  $T_{\text{align}} = \text{polylog}(d)$  samples are sufficient to drive the risk  $\mathcal{R}(t)$  arbitrarily small (at a rate of  $\mathcal{O}(1/\log d)$ ). Computationally, we can solve this Lasso objective up to a  $\mathcal{O}(1/\log d)$  error tolerance with a cost of  $\mathcal{O}(r_s^2 \text{polylog } d)$ . This overhead is negligible compared to the first phase of Algorithm 1, which requires a matrix orthogonalization (e.g., QR decomposition) at every iteration.

#### Appendix E. Some moment bounds and concentration inequalities

**Lemma 38 (Hypercontractivity)** *Let  $P_k : \mathbf{R}^d \rightarrow \mathbf{R}$  be a polynomial of degree- $k$  and  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ . For  $q \geq 2$ , we have  $\mathbb{E}[P_k(\mathbf{x})^q]^{1/q} \leq (q-1)^{k/2} \mathbb{E}[P_k(\mathbf{x})^2]^{1/2}$ .*

**Lemma 39** *Let  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $\mathbf{S} \in \mathbf{R}^{d \times d}$  be a symmetric matrix. For  $u > 0$ ,*

$$\mathbb{P}\left[|\mathbf{x}^\top \mathbf{S} \mathbf{x} - \text{tr}(\mathbf{S})| \geq 2\|\mathbf{S}\|_F u + 2\|\mathbf{S}\|_2 u^2\right] \leq 2e^{-u^2}.$$

**Proof** We note that  $\mathbf{x}^\top \mathbf{S} \mathbf{x} - \text{tr}(\mathbf{S})$  has the same distribution with  $\sum_{i=1}^d \lambda_i(\mathbf{S})(Z_i^2 - 1)$ , where  $Z_i \sim_{iid} \mathcal{N}(0, 1)$ . By using the Laurent-Massart lemma [14], we have the result. ■

**Corollary 40** *Let  $y = \mathbf{x}^\top \mathbf{S} \mathbf{x} - \text{tr}(\mathbf{S})$  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  $\mathbb{E}[|y|^2] = 2\|\mathbf{S}\|_F^2$ , we have the result. ■

**Corollary 41** *For  $\mathbf{A} \in \mathbf{R}^{d \times r}$ ,  $p \geq 2$  and  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ , we have  $\mathbb{E}[\|\mathbf{A}^\top \mathbf{x}\|_2^{2p}]^{\frac{1}{p}} \leq \sqrt{3}(p-1) \text{tr}(\mathbf{A}^\top \mathbf{A})$ .*

**Proof** By Lemma 38, we have  $\mathbb{E}[\|\mathbf{A}^\top \mathbf{x}\|_2^{2p}]^{\frac{1}{p}} \leq (p-1)\mathbb{E}[\|\mathbf{A}^\top \mathbf{x}\|_2^4]^{\frac{1}{2}}$ . For  $\mathbf{S} = \mathbf{A}\mathbf{A}^\top$ , we have

$$\mathbb{E}[\|\mathbf{A}^\top \mathbf{x}\|_2^4] = \mathbb{E}[(\mathbf{x}^\top \mathbf{S} \mathbf{x})^2] = \text{tr}(\mathbb{E}[(\mathbf{x}^\top \mathbf{S} \mathbf{x})\mathbf{x}\mathbf{x}^\top]\mathbf{S}).$$

We have

$$\mathbb{E}[(\mathbf{x}^\top \mathbf{S} \mathbf{x})\mathbf{x}\mathbf{x}^\top] = \text{tr}(\mathbf{S})\mathbf{I}_d + 2\mathbf{S} \Rightarrow \mathbb{E}[\|\mathbf{A}^\top \mathbf{x}\|_2^4] = \text{tr}(\mathbf{S})^2 + 2\|\mathbf{S}\|_F^2 \stackrel{(a)}{\leq} 3\text{tr}(\mathbf{S})^2,$$

where (a) follows that  $\mathbf{S}$  is positive semi-definite. Since  $\text{tr}(\mathbf{S}) = \text{tr}(\mathbf{A}^\top \mathbf{A})$ , we have the statement. ■