

Fast Global Convergence of Online PCA

Zeyuan Allen-Zhu
zeyuan@csail.mit.edu
Princeton University

Yuanzhi Li
yuanzhil@cs.princeton.edu
Princeton University

July 26, 2016

Abstract

We study online principle component analysis (PCA), that is to find the top k eigenvectors of a $d \times d$ hidden matrix Σ with online data samples drawn from covariance matrix Σ . We provide *global* convergence for the low-rank generalization of Oja’s algorithm, which is popularly used in practice but lacks theoretical understanding.

Our convergence rate matches the lower bound in terms of the dependency on error, on eigengap and on dimension d ; in addition, our convergence rate can be made gap-free, that is proportional to the approximation error and independent of the eigengap.

In contrast, for general rank k , before our work (1) it was open to design any algorithm with efficient global convergence rate [9]; and (2) it was open to design any algorithm with (even local) gap-free convergence rate [8].

1 Introduction

Principle component analysis (PCA) is the problem of finding the subspace of largest variance in a dataset consisting of vectors, and is a fundamental tool used to analyze and visualize data in machine learning, computer vision, statistics, and operations research. In the big-data scenario, since it can be unrealistic to store the entire dataset, it is interesting and more challenging to study the online model (a.k.a. the stochastic model or the streaming model) of PCA.

Suppose the data vectors $x \in \mathbb{R}^d$ are drawn i.i.d. from an unknown distribution with covariance matrix $\Sigma = \mathbb{E}[xx^\top] \in \mathbb{R}^{d \times d}$, and the vectors are presented to the algorithm in an online manner. Suppose without loss of generality that the Euclidean norm $\|x\|_2 \leq 1$ for such random vectors, and we are interested in approximately computing the top k eigenvectors of Σ . We are interested in algorithms with memory storage $O(dk)$, the same as the memory needed to store any k vectors in d dimensions. We call this the *online k -PCA problem*.

For online k -PCA, the popular and natural extension of Oja’s algorithm originally designed for the $k = 1$ case works as follows. Beginning with a random Gaussian matrix $\mathbf{Q}_0 \in \mathbb{R}^{d \times k}$ (each entry i.i.d. $\sim \mathcal{N}(0, 1)$), it repeatedly applies

$$\text{rank-}k \text{ Oja's algorithm: } \quad \mathbf{Q}_t \leftarrow (\mathbf{I} + \eta_t x_t x_t^\top) \mathbf{Q}_{t-1}, \quad \mathbf{Q}_t = \text{QR}(\mathbf{Q}_t) \quad (1.1)$$

where $\eta_t > 0$ is some learning rate that may depend on t , vector x_t is the random data vector obtained in iteration t , and $\text{QR}(\mathbf{Q}_t)$ is an arbitrary QR decomposition that orthonormalize the column vectors of \mathbf{Q}_t (i.e., the Gram-Schmidt Orthogonalization).

Although Oja’s algorithm works reasonably well in practice, very limited theoretical results are known for its convergence in the $k > 1$ case. Even worse, little is known for *any* algorithm that solves online PCA in the $k > 1$. Specifically, there are three major challenges for this problem:

	Paper	Global Convergence	Is It “Efficient”?	Local Convergence
$k = 1$ gap- dependent	Shamir [15]	$\tilde{O}(\frac{d}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$	no	$\tilde{O}(\frac{1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
	Sa et al. [14]	$\tilde{O}(\frac{d}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$	no	$\tilde{O}(\frac{d}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
	Li et al. [10] ^a	$\tilde{O}(\frac{d\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$	no	$\tilde{O}(\frac{d\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
	Jain et al. [9]	$\tilde{O}(\frac{\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$	yes	$\tilde{O}(\frac{\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
	this paper: Theorem 1	$\tilde{O}(\frac{\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$	yes	$\tilde{O}(\frac{\lambda_1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
$k = 1$ gap-free	Shamir [15] (Remark 1.3)	$\tilde{O}(\frac{d}{\rho^2} \cdot \frac{1}{\varepsilon^2})$	no	$\tilde{O}(\frac{1}{\rho^2} \cdot \frac{1}{\varepsilon^2})$
	this paper: Theorem 2	$\tilde{O}(\frac{1}{\rho^2} \cdot \frac{1}{\varepsilon})$	yes	$\tilde{O}(\frac{1}{\rho^2} \cdot \frac{1}{\varepsilon})$
$k \geq 1$ gap- dependent	Hardt-Price [8] ^b	$\tilde{O}(\frac{d\lambda_k}{\text{gap}^3} \cdot \frac{1}{\varepsilon})$	no	$\tilde{O}(\frac{d\lambda_k}{\text{gap}^3} \cdot \frac{1}{\varepsilon})$
	Shamir [16]	unknown	n/a	$O(\frac{1}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
	this paper: Theorem 1	$\tilde{O}(\frac{\lambda_1 + \dots + \lambda_k}{\text{gap}^2} \cdot (\frac{1}{\varepsilon} + k))$	yes	$\tilde{O}(\frac{\lambda_1 + \dots + \lambda_k}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$
$k \geq 1$ gap-free	this paper: Theorem 2	$\tilde{O}(\frac{k}{\rho^2} \cdot \frac{1}{\varepsilon})$	yes	$\tilde{O}(\frac{k}{\rho^2} \cdot \frac{1}{\varepsilon})$

Table 1: Sampling complexity comparison. Since we assumed $\|x\| \leq 1$ for each sample vector, we have $\lambda_i \in [0, 1/i]$ and $\lambda_1 + \dots + \lambda_k \leq 1$. We define $\text{gap} = \lambda_{k+1} - \lambda_k \in [0, 1/k]$. We assume $\varepsilon \in (0, 1)$.

- Gap-dependent convergence: $\|\mathbf{Q}_T^\top \mathbf{Z}\|_F^2 \leq \varepsilon$ where \mathbf{Z} consists of the last $d - k$ eigenvectors.
- Gap-free convergence: $\|\mathbf{Q}_T^\top \mathbf{W}\|_F^2 \leq \varepsilon$ where \mathbf{W} consists of all eigenvectors with values no more than $\lambda_k - \rho$.
- We say a global convergence is “efficient” if it only (poly-)logarithmically depend on the dimension d .

^aLi et al. proved their result under a stronger 4-th moment assumption, and obtained running time $\tilde{O}(\frac{d\lambda_1^2}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$, a factor $\lambda_1 \in (0, 1)$ faster than what we show in this table. We believe their running time will be slowed down at least by a factor λ_1 if the 4-th moment assumption is removed.

^bTheir result gives a guarantee on the spectral norm $\|\mathbf{Q}_T^\top \mathbf{W}\|_2^2$ so we increased it by a factor k for a fair comparison.

1. Provide an *efficient* convergence rate that only logarithmically dependent on the dimension d .
2. Provide a *gap-free* convergence rate that is independent of the eigengap.
3. Provide a *global* convergence rate so the algorithm can start from a random initial point.

In the case of $k > 1$, to the best of our knowledge, there is no convergence result that is gap-free. In the gap-dependent regime, the convergence result of Shamir [16] is efficient but not global. The convergence result of Hardt and Price [8] is global but not efficient. We discuss them more formally below (and see Table 1):

- Shamir [16] provided implicitly a *local* but efficient convergence result for Oja’s algorithm,¹ which requires a very accurate starting matrix \mathbf{Q}_0 : his theorem relies on \mathbf{Q}_0 being correlated with the top k eigenvectors by a correlation value at least $k-1/2$. If using random initialization, this event happens with probability at most $2^{-\Omega(d)}$.

¹The original method of Shamir [16] is an offline method that uses variance reduction. We have translated his result into an online setting which requires a lot of extra work including the martingale techniques we used in this paper.

- Hardt and Price [8] analyzed a variant of Oja’s algorithm² and obtained a global convergence that is not efficient: it linearly depends on the dimension d . Their result also has a cubic dependency on the gap between the k -th and $(k+1)$ -th eigenvalue which is not optimal. They raised an open question regarding how to provide any convergence result that is gap-free.
- In practice, researchers observed that it is advantageous to choose the learning rate η_t to be high at the beginning, and then gradually decreasing (c.f. [19]). To the best of our knowledge, there is no theoretical support behind this learning rate scheme for general k .

In sum, it remains open before our work to obtain an efficient and global convergence rate, or any gap-free convergence rate.

Special Case of $k = 1$. The seminal work by Jain, Jin, Kakade, Netrapalli and Sidford [9] obtained a convergence result that is both efficient and global (but not gap-free) for online 1-PCA. Shamir [15] obtained the first gap-free result for online 1-PCA, but his result is not efficient. Both these results are based on Oja’s algorithm, and it remains open before our work to obtain a gap-free result that is also efficient even when $k = 1$.

1.1 Our Results

In this paper we analyze the rank- k variant of Oja’s algorithm (1.1) —or Oja’s algorithm for short. We present convergence results that are *global*, *efficient* and *gap-free*.

Gap-Dependent Online k -PCA. We prove the following theorem in this paper:

Theorem 1 (gap-dependent online k -PCA). *Letting $\text{gap} \stackrel{\text{def}}{=} \lambda_k - \lambda_{k+1} \in (0, \frac{1}{k}]$ and $\Lambda \stackrel{\text{def}}{=} \sum_{i=1}^k \lambda_i \in (0, 1]$, for every $\varepsilon, p \in (0, 1)$ define learning rates*

$$T_0 = \tilde{\Theta}\left(\frac{k\Lambda}{\text{gap}^2 p^2}\right), \quad T_1 = \tilde{\Theta}\left(\frac{\Lambda}{\text{gap}^2 p^2}\right), \quad \eta_t = \begin{cases} \tilde{\Theta}\left(\frac{1}{\text{gap} \cdot T_0}\right) & 1 \leq t \leq T_0; \\ \tilde{\Theta}\left(\frac{1}{\text{gap} \cdot T_1}\right) & T_0 < t \leq T_0 + T_1; \\ \tilde{\Theta}\left(\frac{1}{\text{gap} \cdot t}\right) & t > T_0 + T_1. \end{cases}^3$$

Let \mathbf{Z} be the column orthonormal matrix consisting of all eigenvectors of Σ with values no more than λ_{k+1} . Then, the output $\mathbf{Q}_T \in \mathbb{R}^{d \times k}$ of Oja’s algorithm satisfies:

$$\text{for every}^4 \quad T = T_0 + T_1 + \tilde{\Theta}\left(\frac{T_1}{\varepsilon}\right) \quad \text{it satisfies} \quad \Pr[\|\mathbf{Z}^\top \mathbf{Q}_T\|_F^2 \geq \varepsilon] \leq p.$$

Above, $\tilde{\Theta}$ hides poly-log factors in $\frac{1}{\Lambda}, \frac{1}{p}, \frac{1}{\text{gap}}$ and d .

In other words, after a warm up phase of length T_0 , we obtain a $\frac{\Lambda}{\text{gap}^2 T}$ convergence rate for the quantity $\|\mathbf{Z}^\top \mathbf{Q}_T\|_F^2$. We make several observations (see also Table 1):

- In the $k = 1$ case, Theorem 1 matches the best known result of Jain et al. [9].
- In the $k > 1$ case, Theorem 1 gives the first efficient global convergence rate.
- In the $k > 1$ case, even in terms of local convergence rate, Theorem 1 is faster than the best known result of Shamir [16] by a factor $\lambda_1 + \dots + \lambda_k \in (0, 1)$.

²They used multiple samples in each iteration and over congested the dimension from k to $2k$.

³The intermediate stage $[T_0, T_0 + T_1]$ is completely unnecessary; we add this phase only to simplify proofs.

⁴Theorem also applies to every $T \geq T_0 + T_1 + \tilde{\Omega}(T_1/\varepsilon)$ by making η_t poly-logarithmically dependent on T .

- The learning rates η_t are constants for $t \leq T_0$ and inversely proportional to $1/t$ for large t . To the best of our knowledge, this is the first theoretical justification of this popular learning rate choices researchers have used in practice for general k .

Remark 1.1. The quantity $\|\mathbf{Z}^\top \mathbf{Q}_T\|_F^2$ captures the correlation between the resulting matrix $\mathbf{Q}_T \in \mathbb{R}^{d \times k}$ and the smallest $d - k$ eigenvectors of $\mathbf{\Sigma}$. It is a natural generalization of the sin-square quantity widely used in the $k = 1$ case, because if $k = 1$ then $\|\mathbf{Z}^\top \mathbf{Q}_T\|_F^2 = \sin^2(q, v_1)$ where q is the only column of \mathbf{Q} and v_1 is the leading eigenvector of $\mathbf{\Sigma}$.

Remark 1.2. We are aware of an information-theoretical lower bound of $\Omega(\frac{k\lambda_k}{\text{gap}^2} \cdot \frac{1}{\varepsilon})$ for gap-dependent online k -PCA. [11] Therefore, the local convergence in Theorem 1 is optimal up to logarithmic factors (at least when $\lambda_1 = \dots = \lambda_k$).

Gap-Free Online k -PCA. When the eigengap is small which is usually true in practical applications, it is desirable to obtain gap-free convergence rates [13, 15]. We prove the following theorem, and thus fully answer the open question of Hardt and Price [8] regarding how to obtain gap-free convergence rate for online k -PCA.

Theorem 2 (gap-free online k -PCA). *For every $\rho, \varepsilon, p \in (0, 1)$, define learning rates*

$$T_0 = \tilde{\Theta}\left(\frac{k}{\rho^2 \cdot p^2}\right), \quad \eta_t = \begin{cases} \tilde{\Theta}\left(\frac{1}{\rho \cdot T_0}\right) & t \leq T_0; \\ \tilde{\Theta}\left(\frac{1}{\rho \cdot t}\right) & t > T_0. \end{cases}$$

Let \mathbf{W} be the column orthonormal matrix consisting of all eigenvectors of $\mathbf{\Sigma}$ with values no more than $\lambda_k - \rho$. Then, the output $\mathbf{Q}_T \in \mathbb{R}^{d \times k}$ of Oja's algorithm satisfies:

$$\text{for every}^5 \quad T = T_0 + \tilde{\Theta}\left(\frac{T_0}{\varepsilon}\right) \quad \text{it satisfies} \quad \Pr[\|\mathbf{W}^\top \mathbf{Q}_T\|_F^2 \geq \varepsilon] \leq p.$$

Above, $\tilde{\Theta}$ hides poly-log factors in $\frac{1}{p}$, $\frac{1}{\rho}$ and d .

Note that the above theorem is a *double approximation*. The number of iterations depend both on ρ and ε , where ε is an upper bound on the correlation between \mathbf{Q}_T and all eigenvectors in \mathbf{W} (which depends on ρ). This is the first known gap-free result for the $k > 1$ case.

One may also be interested in single-approximation guarantees, such as the rayleigh-quotient guarantee. Note that a single-approximation guarantee by definition loses information about the $\varepsilon - \rho$ tradeoff; furthermore, (good) single-approximation guarantees are not easy to obtain.⁶

We show in this paper the following theorem regarding rayleigh-quotient guarantee:

Theorem 3 (gap-free rayleigh-quotient guarantee). *In the same setting as Theorem 2, we have for every $T = \tilde{\Theta}(\frac{k}{p^2 \rho})$, letting q_i be the i -th column of the output matrix \mathbf{Q}_T , then*

$$\Pr\left[\forall i \in [k], \quad q_i^\top \mathbf{\Sigma} q_i \geq \lambda_i - \tilde{\Theta}(\rho)\right] \geq 1 - p.$$

Again, $\tilde{\Theta}$ hides poly-log factors in $\frac{1}{p}$, $\frac{1}{\rho}$ and d .

Remark 1.3. The only gap-free result known before our work is Shamir [15] — and it is only for $k = 1$ and not efficient due to its heavy initialization. Shamir's result is in terms of Rayleigh quotient but

⁵Theorem also applies to every $T \geq T_0 + \tilde{\Omega}(T_0/\varepsilon)$ by making η_t poly-logarithmically dependent on T .

⁶For instance, as pointed out by the authors of [9], a direct translation from a correlation-type convergence to a rayleigh-quotient type convergence loses a factor on the approximation error. They even raised it as an open question regarding how to design a direct proof without sacrificing this loss. Thus, our next theorem answers this open question (at least in the gap-free case).

not double-approximation. If the initialization phase is ignored, Shamir’s local convergence rate in terms of Rayleigh quotient in fact matches our *global* convergence rate in Theorem 3. However, if one translates his result into double approximation, his running time will lose a factor ε . This is why in Table 1 Shamir’s result [15] is in terms of $1/\varepsilon^2$ as opposed to $1/\varepsilon$.

Other Related Results. Mitliagkas et al. [12] obtained an online PCA result but in the restricted spiked covariance model. Balsubramani et al. [3] analyzed a modified variant of Oja’s algorithm and needed an extra $O(d^5)$ factor in the complexity.

The offline problem of PCA (or more generally of SVD) can be efficiently solved via iterative algorithms that are based on variance-reduction techniques on top of stochastic gradient methods [2, 16] (see also [5, 6] for the $k = 1$ case); these methods do multiple passes on the input data so are not relevant in our online setting. Offline PCA can also be solved via power method or block Krylov method [13], but since each iteration of these methods relies on one full pass on the dataset, they are not suitable for online setting either. Other offline problems and efficient algorithms relevant to PCA include canonical correlation analysis and generalized eigenvector decomposition [1, 7, 18].

We emphasize that the offline problem is *much easier* to solve and one can efficiently (although non-trivially) reduce a general k -PCA problem to k times of 1-PCA using the techniques of [2]. However, this is *not the case* in our online setting because one would have to lose a $\text{poly}(k)$ factor in the iteration complexity and sampling complexity.

2 Preliminaries

We denote by $1 \geq \lambda_1 \geq \dots \geq \lambda_d \geq 0$ the eigenvalues of the positive semidefinite (PSD) matrix Σ , and since we have assumed $\|x\| \leq 1$ for each online data sample, it must satisfies $\lambda_1 + \dots + \lambda_d = \text{Tr}(\Sigma) \leq 1$ and thus each $\lambda_i \leq 1/i$. We define $\text{gap} \stackrel{\text{def}}{=} \lambda_k - \lambda_{k+1} \in [0, \frac{1}{k}]$.

We denote by $\mathbf{V} \in \mathbb{R}^{d \times k}$ the matrix of the first k eigenvectors of Σ (in the non-increasing order eigenvalues) and $\mathbf{Z} \in \mathbb{R}^{d \times (d-k)}$ the last $d - k$ eigenvectors (also in the non-increasing order eigenvalues). For every parameter $\rho > 0$ in our gap-free setting, we also define $\mathbf{W} \in \mathbb{R}^{d \times r}$ to be column orthonormal matrix consisting of all eigenvectors of Σ with values no more than $\lambda_k - \rho$. It is clear that $r \leq d - k$.

We write $\Sigma_{\leq k} = \mathbf{V} \text{Diag}\{\lambda_1, \dots, \lambda_k\} \mathbf{V}^\top$ and $\Sigma_{> k} \stackrel{\text{def}}{=} \mathbf{Z} \text{Diag}\{\lambda_{k+1}, \dots, \lambda_d\} \mathbf{Z}^\top$ so $\Sigma = \Sigma_{\leq k} + \Sigma_{> k}$.

For a vector y , we sometimes denote by $y[i]$ or $y^{(i)}$ the i -th coordinate of y . We may use different notations in different lemmas in order to obtain the cleanest representations; when we do so, we shall clearly point out in the statement of the lemmas.

We denote by $\mathbf{P}_t \stackrel{\text{def}}{=} \prod_{s=1}^t (\mathbf{I} + \eta_s x_s x_s^\top)$ where x_s is the s -th data sample and η_s is the learning rate of iteration s . We denote by $\mathbf{Q} \in \mathbb{R}^{d \times k}$ (or \mathbf{Q}_0) the random initial matrix, and by $\mathbf{Q}_t \stackrel{\text{def}}{=} \text{QR}((\mathbf{I} + \eta_t x_t x_t^\top) \mathbf{Q}_{t-1}) = \text{QR}(\mathbf{P}_t \mathbf{Q}_0)$ for every $t \geq 1$.⁷ We use the notation \mathcal{F}_t to denote the sigma-algebra generated by x_t . We denote $\mathcal{F}_{\leq t}$ to be the sigma-algebra generated by x_1, \dots, x_t , i.e. $\mathcal{F}_{\leq t} = \vee_{s=1}^t \mathcal{F}_s$. In other words, whenever we condition on $\mathcal{F}_{\leq t}$ it means we have fixed x_1, \dots, x_t .

For a vector x we denote by $\|x\|$ or $\|x\|_2$ the Euclidean norm of x . We denote by $\|\mathbf{A}\|_{S_1}$ the Schatten-1 norm of matrix \mathbf{A} which is the summation of the (nonnegative) singular values of \mathbf{A} . It satisfies the following simple properties:

Proposition 2.1. *For not necessarily symmetric matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ we have*

$$(1): |\text{Tr}(\mathbf{A})| \leq \|\mathbf{A}\|_{S_1} \quad (2): |\text{Tr}(\mathbf{AB})| \leq \|\mathbf{AB}\|_{S_1} \leq \|\mathbf{A}\|_{S_1} \|\mathbf{B}\|_2 .$$

⁷The second equality is simple fact but anyways proved in Lemma 2.2 later.

$$(3): \mathbf{Tr}(\mathbf{AB}) \leq \|A\|_F \|B\|_F = (\mathbf{Tr}(\mathbf{A}^\top \mathbf{A}) \mathbf{Tr}(\mathbf{B}^\top \mathbf{B}))^{1/2}.$$

Proof. (1) is because $\mathbf{Tr}(\mathbf{A}) = \frac{1}{2} \mathbf{Tr}(\mathbf{A} + \mathbf{A}^\top) \leq \frac{1}{2} \|\mathbf{A} + \mathbf{A}^\top\|_{S_1} \leq \frac{1}{2} (\|\mathbf{A}\|_{S_1} + \|\mathbf{A}^\top\|_{S_1}) = \|\mathbf{A}\|_{S_1}$. (2) is because of (1) and the matrix Holder's inequality. (3) is owing to von Neumann's trace inequality (together with Cauchy's) which says $\mathbf{Tr}(\mathbf{AB}) \leq \sum_i \sigma_{A,i} \cdot \sigma_{B,i} \leq \|A\|_F \|B\|_F$. (Here, we have noted by $\sigma_{A,i}$ the i -th largest eigenvalue of A and similarly for B . \square)

2.1 A Matrix View of Oja's Algorithm

The following lemma tells us that we can push the QR orthogonalization step in Oja's algorithm to the end for analysis purpose only:

Lemma 2.2 (Oja's algorithm). *For every $s \in [d]$, every $\mathbf{X} \in \mathbb{R}^{d \times s}$, every $t \geq 1$, every $\mathbf{Q} \in \mathbb{R}^{d \times k}$, it satisfies $\|\mathbf{X}^\top \mathbf{Q}_t\|_F \leq \|\mathbf{X}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$.*

Proof of Lemma 2.2. Denote by $\tilde{\mathbf{Q}}_t = \mathbf{P}_t \mathbf{Q}$, we first observe that for every $t \geq 0$ $\mathbf{Q}_t = \tilde{\mathbf{Q}}_t \mathbf{R}_t$ for some (upper triangular) invertible matrix $\mathbf{R}_t \in \mathbb{R}^{k \times k}$. The claim is true for $t = 0$. Suppose it holds for t by induction, then

$$\mathbf{Q}_{t+1} = \text{QR}[(\mathbf{I} + \eta_{t+1} x_{t+1} x_{t+1}^\top) \mathbf{Q}_t] = (\mathbf{I} + \eta_{t+1} x_{t+1} x_{t+1}^\top) \mathbf{Q}_t \mathbf{S}_t$$

for some $\mathbf{S}_t \in \mathbb{R}^{k \times k}$ by the definition of QR (or Gram-Schmidt). This implies that

$$\mathbf{Q}_{t+1} = (\mathbf{I} + \eta_{t+1} x_{t+1} x_{t+1}^\top) \tilde{\mathbf{Q}}_t \mathbf{R}_t \mathbf{S}_t = \mathbf{P}_{t+1} \mathbf{Q} \mathbf{R}_t \mathbf{S}_t = \tilde{\mathbf{Q}}_{t+1} \mathbf{R}_t \mathbf{S}_t = \tilde{\mathbf{Q}}_{t+1} \mathbf{R}_{t+1}$$

if we define $\mathbf{R}_{t+1} = \mathbf{R}_t \mathbf{S}_t$. This completes the proof that $\mathbf{Q}_t = \tilde{\mathbf{Q}}_t \mathbf{R}_t$. As a result, since each \mathbf{Q}_t is column orthogonal for $t \geq 1$ (thus $\|\mathbf{V}^\top \mathbf{Q}_t\|_2 \leq 1$):

$$\|\mathbf{X}^\top \mathbf{Q}_t\|_F \leq \|\mathbf{X}^\top \mathbf{Q}_t (\mathbf{V}^\top \mathbf{Q}_t)^{-1}\|_F = \|\mathbf{X}^\top \tilde{\mathbf{Q}}_t \mathbf{R}_t (\mathbf{V}^\top \tilde{\mathbf{Q}}_t \mathbf{R}_t)^{-1}\|_F \leq \|\mathbf{X}^\top \tilde{\mathbf{Q}}_t (\mathbf{V}^\top \tilde{\mathbf{Q}}_t)^{-1}\|_F. \quad \square$$

Due to Lemma 2.2, we make an important observation that is in order to prove Theorem 1 and Theorem 2, it suffices to upper bound the quantity $\|\mathbf{X}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$ for $\mathbf{X} = \mathbf{W}$ or $\mathbf{X} = \mathbf{Z}$.

3 Overview of Our Proofs and Techniques

Let us focus on the gap-dependent case first. Denoting in this section by $s_t \stackrel{\text{def}}{=} \|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$, owing to Lemma 2.2, we want to bound s_t in terms of x_t and $s_{t-1} = \|\mathbf{Z}^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_F$. A simple calculation using the Sherman-Morrison formula gives

$$\mathbb{E}[s_t^2] \leq (1 - \eta_t \text{gap}) \mathbb{E}[s_{t-1}^2] + \mathbb{E} \left[\left(\frac{\eta_t a_t}{1 - \eta_t a_t} \right)^2 \right] \quad \text{where} \quad a_t = \|x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_2 \quad (3.1)$$

At a first look, $\mathbb{E}[s_t^2]$ is decaying by multiplicative $(1 - \eta_t \text{gap})$ factor at every iteration; however, this bound could be *problematic* when $\eta_t a_t$ is close to 1 and thus we need to ensure $\eta_t \leq \frac{1}{a_t}$ with high probability for every step.

A naive bound on a_t gives $a_t \leq \|\mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_2 \leq s_t + 1$. However, since s_t can be as large as $\Omega(\sqrt{d})$ at $t = 0$ if random initialization is used, this would imply that η_t can be at most $1/\sqrt{d}$ and the resulting convergence rate would certainly be *not* efficient (i.e., at least proportional to d). This is why most known global results are not efficient (see Table 1). On the other hand, if one ignores initialization and starts from a point t_0 when $s_{t_0} \leq 1$ is already satisfied, then he or she can prove a *local* convergence rate that is efficient (c.f. [16]) but still slower than ours.

Our *first contribution* is the following crucial observation: for a random initial matrix \mathbf{Q} , $a_1 = \|x_1^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_2$ is actually quite small. We use a simple fact on the singular value distribution

of inverse-Wishart distribution to obtain that, with high probability, $a_1 = O(\sqrt{k})$. This implies, at least in the first iteration, we can set η_1 to be $\Omega(1/\sqrt{k})$ independent of the dimension d . However, in subsequent iterations, it is not clear whether a_t increases.

Our *second contribution* is to control a_t using the fact that a_t itself “forms another random process.” More precisely, denoting by $a_{t,s} = \|x_t^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2$ for $0 \leq s \leq t-1$, we wish to bound $a_{t,s}$ in terms of $a_{t,s-1}$ and show that it does not increase by much. (If we could achieve so, combining with the initialization $a_{t,0} \leq O(\sqrt{k})$ we would know that all $a_{t,s}$ are small for $s \leq t-1$.) Unfortunately, since x_t is not an eigenvector of Σ , the recursion one can obtain is (again using Sherman-Morrison)

$$\mathbb{E}[a_{t,s}^2] \leq (1 - \eta_s \lambda_k) \mathbb{E}[a_{t,s-1}^2] + \eta_s \lambda_k \mathbb{E}[b_{t,s-1}^2] + \mathbb{E}\left[\left(\frac{\eta_s a_s}{1 - \eta_s a_s}\right)^2\right] \quad (3.2)$$

where $b_{t,s} = \|x_t^\top \Sigma \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2$. Now two difficulties arise from this formula:

- $b_{t,s}$ can be very different from $a_{t,s}$ — in worse case, the ratio between them can be unbounded.
- the problematic term now becomes $a_s = a_{s,s-1}$ (rather than the original $a_t = a_{t,t-1}$ in (3.1)) which is not present in the chain $\{a_{t,s}\}_{s=1}^{t-1}$.

We solve both issues by considering a multi-dimensional random process $c_{t,s}$ with $c_{t,s}^{(i)} \stackrel{\text{def}}{=} \|x_t^\top \Sigma^i \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2$. Ignoring the last term, we can derive that

$$\forall t, \forall s \leq t-1, \quad \mathbb{E}[(c_{t,s}^{(i)})^2] \lesssim (1 - \eta_s \lambda_k) \mathbb{E}[(c_{t,s-1}^{(i)})^2] + \eta_s \lambda_k \mathbb{E}[(c_{t,s-1}^{(i+1)})^2] . \quad (3.3)$$

Our *third contribution* is a new random process concentration bound to control the change in this multi-dimensional chain (3.3). To achieve this, we also adapt the prove of standard Chernoff bound to multi dimensions (which is not the same as matrix concentration bound). After having this concentration result (see Section 6), all terms of $a_t = c_{t,t-1}^{(0)}$ can be simultaneously bounded by a constant, for every $t \in [T]$. This ensures that the problematic term in (3.1) is well-controlled.

The overall plan looks promising, however, there are holes in the above thought experiment.

- In order to apply any random process concentration bound (e.g., any martingale concentration), we need the process to not depend on the future. However, the random vector $c_{t,s}$ is not $\mathcal{F}_{\leq s}$ measurable but $\mathcal{F}_{\leq s} \vee \mathcal{F}_t$ measurable (i.e., it depends on x_t for a future $t > s$).
- Furthermore, the expectation bounds such as (3.1), (3.2), (3.3) only hold if $\mathbb{E}[x_t x_t] = \Sigma$; however, if we take away a failure event \mathcal{C} — \mathcal{C} may correspond to the event when a_t is large — the conditional expectation $\mathbb{E}[x_t x_t \mid \bar{\mathcal{C}}]$ becomes $\Sigma + \Delta$ where Δ is some error matrix. This can amplify the failure probability in next iteration.

Our *fourth contribution* is a “decoupling” framework to deal with the above issues (see Section D). At a high level, to deal with the first issue we fix x_t and study $\{c_{t,s}\}_{s=0,1,\dots,t-1}$ conditioning on x_t ; in this way the process decouples and each $c_{t,s}$ becomes $\mathcal{F}_{\leq s}$ measurable. We can do so because we can carefully ensure that the failure events only depend on x_s for $s \leq t-1$ but not on x_t . To deal with the second issue, we convert the random process into an unconditional random process (see (D.2)); this is a generalization of using stopping time on martingales. Using these tools, we manage to show that the failure probability only grows linearly with respect to T and henceforth bound the value of $c_{t,s}^{(i)}$ for all t, s and i .

Although each of our contributions is conceptually not a very big step, putting them together gives us a new way to analyze how certain property of a random initialization is preserved in all subsequent iterations, which we believe is useful in future research (especially when analyzing *any* high-rank online power-method type of algorithm).

Remark 3.1. The above ideas are insufficient for our gap-free results. In order to prove Theorem 2 and 3, in addition to s_t and $c_{t,s}$ discussed above, we also need to bound $s'_t \stackrel{\text{def}}{=} \|\mathbf{W}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$ where \mathbf{W} is a column orthonormal matrix consisting of all eigenvectors of Σ with values no more than $\lambda_k - \rho$, for some parameter ρ given to the algorithm. This is so because the interesting quantity in a gap-free case changes from s_t to s'_t according to Lemma 2.2. Similar to the gap-dependent case, to bound s'_t one has to bound $c_{t,s}$; however, the $c_{t,s}$ process also weakly depends on the original s_t . In sum, we have to bound s_t , s'_t , and $c_{t,s}$ all together.

Roadmap.

- Section 4 proves properties on the initial matrix \mathbf{Q} and corresponds to our first contribution.
- Section 5 gives expected guarantees on s_t and $a_{t,s}$ and corresponds to our second contribution.
- Section 6 provides concentration results which correspond to our third contribution.
- Appendix D gives the decoupling lemma which correspond to our fourth contribution.
- Section 7 gives main convergence lemmas to deal with iterations both before T_0 and after T_0 .
- Section 8 provides final remarks on how to translate Section 7 to our theorem statements.

Our proofs of nearly all technical lemmas and theorems are deferred to the appendix.

4 Random Initialization

Let $\mathbf{Q} \in \mathbb{R}^{d \times k}$ be a matrix with each entry i.i.d drawn from $\mathcal{N}(0, 1)$, the standard gaussian. Then,

Lemma 4.1. *For every $x \in \mathbb{R}^d$ that has Euclidean norm $\|x\|_2 \leq 1$, every PSD matrix \mathbf{A} , and every $\lambda \geq 1$, we have*

$$\Pr_{\mathbf{Q}} [x^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} \mathbf{A} \mathbf{Q}^\top \mathbf{Z} \mathbf{Z}^\top x \geq \text{Tr}(\mathbf{A}) + \lambda] \leq e^{-\frac{\lambda}{8\text{Tr}(\mathbf{A})}}.$$

Lemma 4.2. *Let \mathbf{Q} be our initial matrix, then for every $p \in (0, 1)$:*

$$\Pr_{\mathbf{Q}} \left[\text{Tr} \left[((\mathbf{V}^\top \mathbf{Q})^\top (\mathbf{V}^\top \mathbf{Q}))^{-1} \right] \geq \frac{\pi^2 ek}{3p} \right] \leq \frac{\sqrt{p}}{1-p}.$$

Combining them, one can obtain our main lemma for initialization:

Lemma 4.3 (initialization). *For every $p, q \in (0, 1)$, every $T \in \mathbb{N}^*$, every vector set $\{x_t\}_{t=1}^T$ with $\|x_t\|_2 \leq 1$, with probability at least $1 - p - 2q$ over the random choice of \mathbf{Q} , the following holds:*

$$\left\{ \begin{array}{l} \|\mathbf{Z}^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \frac{576dk}{p^2} \ln \frac{d}{p} \quad \text{and} \\ \Pr_{x_1, \dots, x_T} \left[\exists i \in [T], \exists t \in [T], \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \frac{18}{p} \left(2k \ln \frac{T}{q} \right)^{1/2} \right] \leq q \end{array} \right.$$

We remark here that the two statements of the above lemma correspond to s_0 and $c_{t,0}^{(i)}$ that we defined in Section 3.

5 Expected Results

In this section we provide formal statements of (3.1), (3.2), (3.3) which characterize to the behaviors of the random processes we are interested. Since the quantities $s_t, s'_t, c_{t,s}^{(i)}$ we discussed in Section 3 have the same form, below we provide a general lemma that talks about all of them at once.

Let $\mathbf{X} \in \mathbb{R}^{d \times r}$ be a generic matrix that shall later be chosen as either $\mathbf{X} = \mathbf{W}$ (corresponding to s'_t), $\mathbf{X} = \mathbf{Z}$ (corresponding to s_t), or $\mathbf{X} = [w]$ where $w \in \mathbb{R}^d$ is an arbitrary vector with norm at most 1 (corresponding to $c_{t,s}^{(i)}$). We introduce the following definitions that shall be used throughout this paper:

$$\begin{aligned} \mathbf{L}_t &= \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1} \in \mathbb{R}^{d \times k} & \mathbf{R}'_t &= \mathbf{X}^\top x_t x_t^\top \mathbf{L}_{t-1} \in \mathbb{R}^{r \times k} \\ \mathbf{S}_t &= \mathbf{X}^\top \mathbf{L}_t \in \mathbb{R}^{r \times k} & \mathbf{H}'_t &= \mathbf{V}^\top x_t x_t^\top \mathbf{L}_{t-1} \in \mathbb{R}^{k \times k} \end{aligned}$$

We present a generic lemma that holds for all of the three choices of \mathbf{X} :

Lemma 5.1. *For every $\mathbf{Q} \in \mathbb{R}^{d \times k}$ and every $t \in [T]$, suppose for $\phi_t \geq 0$, x_t satisfies:*

$$\|x_t^\top \mathbf{L}_{t-1}\|_2 = \|x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_2 \leq \phi_t \quad \text{and} \quad \eta_t \phi_t \leq \frac{1}{2}.$$

Then the following holds:

- (a) $\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)$
 $+ (12\eta_t^2 \|\mathbf{H}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2$
- (b) $|\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})|^2 \leq 243\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 12\eta_t^2 \|\mathbf{R}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 300\eta_t^4 \phi_t^2 \|\mathbf{R}'_t\|_2^2$
- (c) $|\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})| \leq 9\eta_t \phi_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 2\eta_t \phi_t \sqrt{\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})} + 10\eta_t^2 \phi_t^2$

Note that Lemma 5.1-(a) will be used to provide upper bounds on the quantities we care about (i.e., $s_t, s'_t, c_{t,s}^{(i)}$), while Lemma 5.1-(b) and Lemma 5.1-(c) provide variance and absolute difference bounds. We need the latter two bounds in order to provide concentration results.⁸

Taking expectation on top of Lemma 5.1-(a), one can verify that the following is true:

Corollary 5.2 (corollary of Lemma 5.1-(a)). *For every $t \in [T]$, suppose $\mathcal{C}_{\leq t}$ is an event that depends on random x_1, \dots, x_t and implies*

$$\|x_t^\top \mathbf{L}_{t-1}\|_2 = \|x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_2 \leq \phi_t \quad \text{where} \quad \eta_t \phi_t \leq \frac{1}{2}.$$

If $\mathbb{E}[x_t x_t^\top \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] = \Sigma + \Delta$, then we have:

(a) When $\mathbf{X} = \mathbf{Z}$,

$$\begin{aligned} \mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - 2\eta_t \text{gap} + 14\eta_t^2 \phi_t^2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\eta_t^2 \phi_t^2 \\ &\quad + 2\eta_t \|\Delta\|_2 \left([\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{3/2} + 2\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + [\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} \right) \end{aligned}$$

(b) When $\mathbf{X} = \mathbf{W}$,

$$\begin{aligned} \mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - 2\eta_t \rho + 14\eta_t^2 \phi_t^2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\eta_t^2 \phi_t^2 \\ &\quad + 2\eta_t \|\Delta\|_2 \left([\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) \left(1 + [\text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2} \right) \end{aligned}$$

(c) When $\mathbf{X} = [w] \in \mathbb{R}^{d \times 1}$ where w is a vector with Euclidean norm at most 1,

$$\begin{aligned} \mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - \eta_t \lambda_k + 14\eta_t^2 \phi_t^2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\eta_t^2 \phi_t^2 + \frac{\eta_t}{\lambda_k} \|w^\top \Sigma \mathbf{L}_{t-1}\|_2^2 \\ &\quad + 2\eta_t \|\Delta\|_2 \left([\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) \left(1 + [\text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2} \right) \end{aligned}$$

⁸Recall that even in the simplest martingale concentration, one needs upper bounds on the absolute difference between consecutive variables; furthermore, the concentration can be tightened if one also has an (expected) variance upper bound between variables.

6 Martingale Concentrations

We prove in the appendix the following two martingale concentration lemmas. Both of them are stated in their most general form for the purpose of this paper. The first lemma is for 1-d martingales and the second is for multi-d martingales.

At a high level, Lemma 6.1 will only be used to analyze the sequences s_t or s'_t (see Section 3) after warm start — that is, after $t \geq T_0$. Our Lemma 6.2 can be used to analyze $c_{t,s}$ as well as s_t and s'_t before warm start.

Lemma 6.1 (1-d martingale). *Let $\{z_t\}_{t=t_0}^\infty$ be a non-negative random process with starting time $t_0 \in \mathbb{N}^*$. Suppose there exists $\delta > 0$, $\kappa \geq 2$, and $\tau_t = \frac{1}{\delta t}$ such that*

$$\forall t \geq t_0: \left\{ \begin{array}{l} \mathbb{E}[z_{t+1} \mid \mathcal{F}_{\leq t}] \leq (1 - \delta\tau_t)z_t + \tau_t^2 \\ \mathbb{E}[(z_{t+1} - z_t)^2 \mid \mathcal{F}_{\leq t}] \leq \tau_t^2 z_t + \kappa^2 \tau_t^4 \\ |z_{t+1} - z_t| \leq \kappa \tau_t \sqrt{z_t} + \kappa^2 \tau_t^2 \end{array} \right\} \quad (6.1)$$

If there exists $\phi \geq 36$ satisfying $\frac{t_0}{\ln^2 t_0} \geq 7.5\kappa^2(\phi + 1)$ with $z_{t_0} \leq \frac{\phi \ln^2 t_0}{2\delta^2 t_0}$, we have:

$$\Pr \left[\exists t \geq t_0, z_t > \frac{(\phi+1)\ln^2 t}{\delta^2 t} \right] \leq \frac{\exp\{-\left(\frac{\phi}{36}-1\right)\ln t_0\}}{\frac{\phi}{36}-1}.$$

Lemma 6.2 (multi-dimensional martingale). *Let $\{z_t\}_{t=0}^T$ be a random process where each $z_t \in \mathbb{R}_{\geq 0}^D$ is $\mathcal{F}_{\leq t}$ -measurable. Suppose there exist nonnegative parameters $\{\beta_t, \delta_t, \tau_t\}_{t=0}^{T-1}$ satisfying $\kappa \geq 0$ and $\kappa\tau_t \leq 1/6$ such that, $\forall i \in [D], \forall t \in \{0, 1, \dots, T-1\}$,*

(denoting by $[z_t]_i$ is the i -th coordinate of z_t and $[z_t]_{D+1} = 0$)

$$\left. \begin{array}{l} \mathbb{E} [[z_{t+1}]_i \mid \mathcal{F}_{\leq t}] \leq (1 - \beta_t - \delta_t + \tau_t^2) [z_t]_i + \delta_t [z_t]_{i+1} + \tau_t^2, \\ \mathbb{E} [|[z_{t+1}]_i - [z_t]_i|^2 \mid \mathcal{F}_{\leq t}] \leq \tau_t^2 ([z_t]_i^2 + [z_t]_i) + \kappa^2 \tau_t^4, \text{ and} \\ |[z_{t+1}]_i - [z_t]_i| \leq \kappa \tau_t ([z_t]_i + \sqrt{[z_t]_i}) + \kappa^2 \tau_t^2. \end{array} \right\} \quad (6.2)$$

Then, we have: for every $\lambda > 0$, every $p \in [1, \min_{s \in [t]} \{\frac{1}{6\kappa\tau_{s-1}}\}]$:

$$\Pr [[z_t]_1 \geq \lambda] \leq \lambda^{-p} \left(\max_{j \in [t+1]} \{ [z_0]_j^p \} \exp \left\{ \sum_{s=0}^{t-1} 5p^2 \tau_s^2 - p\beta_s \right\} \right. \\ \left. + 1.4 \sum_{s=0}^{t-1} \exp \left\{ \sum_{u=s+1}^{t-1} 5p^2 \tau_u^2 - p\beta_u \right\} \right).$$

The above two lemmas are stated in the most general way in order to be used towards all of our three theorems each requiring different parameter choices of $\beta_t, \delta_t, \tau_t, \kappa$. For instance, to prove Theorem 2 it suffices to use $\kappa = O(1)$.

6.1 Martingale Corollaries

We provide below four instantiations of these lemmas, each of them can be verified by plugging in the specific parameters.

Corollary 6.3 (1-d martingale). *Consider the same setting as Lemma 6.1. Suppose $p \in (0, \frac{1}{e^2})$, $\delta \leq \frac{1}{\sqrt{8}}$, $\tau_t = \frac{1}{\delta t}$, $\kappa \in [2, \frac{1}{\sqrt{2\delta}}]$, $\frac{t_0}{\ln^2 t_0} \geq \frac{9\ln(1/p)}{\delta^2}$, and $z_{t_0} \leq 2$ we have:*

$$\Pr \left[\exists t \geq t_0, z_t > \frac{5(t_0/\ln^2 t_0)}{t/\ln^2 t} \right] \leq p.$$

Corollary 6.4 (multi-d martingale). *Consider the same setting as Lemma 6.2. Suppose $\kappa = 1$, then for every $t \in [T]$ and $q \in (0, 1)$,*

$$\text{if } \sum_{s=0}^{t-1} \tau_s^2 \leq \frac{1}{100} \ln^{-2} \frac{4t}{q} \quad \text{then} \quad \Pr \left[[z_t]_1 \geq 2 \max \left\{ 1, \max_{j \in [t+1]} \{[z_0]_j\} \right\} \right] \leq q .$$

Corollary 6.5 (multi-d martingale). *Consider the same setting as Lemma 6.2. For every $q \in (0, 1)$, letting $l \stackrel{\text{def}}{=} 12 \ln \frac{4t}{q}$, suppose for every $s \in \{0, 1, \dots, t-1\}$ it satisfies $\beta_s \geq l\tau_s^2$ and $\kappa\tau_sl \leq 1$. Then,*

$$\Pr \left[[z_t]_1 \geq 2 \max \left\{ 1, \max_{j \in [t+1]} \{[z_0]_j\} \right\} \right] \leq q .$$

Corollary 6.6 (multi-d martingale). *Consider the same setting as Lemma 6.2. Given $q \in (0, 1)$, suppose there exists parameter $\gamma \geq 1$ such that, denoting by $l \stackrel{\text{def}}{=} 10\gamma \ln \frac{3t}{q}$,*

$$\sum_{s=0}^{t-1} \beta_s - l\tau_s^2 \geq \ln \left(\max_{j \in [t+1]} \{[z_0]_j\} \right) \quad \text{and} \quad \forall s \in \{0, 1, \dots, t-1\} : \beta_s \geq l\tau_s^2 \bigwedge \kappa\tau_s \leq \frac{1}{12 \ln \frac{3t}{q}} .$$

Then, we have

$$\Pr \left[[z_t]_1 \geq 2/\gamma \right] \leq q .$$

7 Main Lemmas

In this section we present our main lemmas. These lemmas can be proved by combining (1) the expectation results in Section 5, (2) the martingale concentrations in Section 6, and (3) our decoupling lemma in Appendix D.

Before Warm Start. Our first lemma describes the behavior of quantities $s_t = \|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$ and $s'_t = \|\mathbf{W}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F$ (defined in Section 3) before warm start. At a high level, it shows if s_t starts from $s_0^2 \leq \Xi_{\mathbf{Z}}$, under mild conditions and with high probability, s_t^2 never increases to more than $2\Xi_{\mathbf{Z}}$. The other sequence $(s'_t)^2$ also never increases to more than $2\Xi_{\mathbf{Z}}$ because $s'_t \leq s_t$, but most importantly, $(s'_t)^2$ drops below 2 after $t \geq T_0$. This means we can choose T_0 as a warm start and proceed to derive a stronger convergence from T_0 (and this is the goal of our next lemma).

We emphasize that although we are only interested in s_t and s'_t , our proof of the lemma also needs to bound the multi-dimensional $c_{t,s}$ sequence discussed in Section 3.

Lemma 7.1 (before warm start). *For every $\rho \in (0, 1)$, $q \in (0, \frac{1}{2}]$, $\Xi_{\mathbf{Z}} \geq 2$, $\Xi_x \geq 2$, and fixed matrix $\mathbf{Q} \in \mathbb{R}^{d \times k}$, suppose it satisfies*

- $\|\mathbf{Z}^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \Xi_{\mathbf{Z}}$, and
- $\Pr_{x_t} \left[\forall j \in [T], \|x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma/\lambda_{k+1})^{j-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_2 \leq \Xi_x \right] \geq 1 - q^2/2$ for every $t \in [T]$.

Suppose also the learning rates $\{\eta_s\}_{s \in [T]}$ satisfy

$$\forall s \in [T], q\Xi_{\mathbf{Z}}^{3/2} \leq \eta_s \leq \frac{\rho}{4000\Xi_x^2 \ln \frac{24T}{q^2}} \quad \text{and} \quad \sum_{t=1}^T \eta_t^2 \Xi_x^2 \leq \frac{1}{100 \ln^2 \frac{32T}{q^2}} .$$

$$\exists T_0 \in [T] \text{ such that } \sum_{t=1}^{T_0} \eta_t \geq \frac{\ln(3\Xi_{\mathbf{Z}})}{\rho}$$

Then, for every $t \in [T-1]$, with probability at least $1 - 2qT$ (over the randomness of x_1, \dots, x_t):

- $\|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq 2\Xi_{\mathbf{Z}}$, and
- if $t \geq T_0$ then $\|\mathbf{W}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq 2$.

Note that the following learning rates satisfy the above lemma:

Parameter 7.2. There exist constants $C_1, C_2 > 0$ such that for every $q > 0$ that is sufficiently small (meaning $q < 1/\text{poly}(T, \Xi_{\mathbf{Z}}, \Xi_x, 1/\rho)$), the following parameters satisfy Lemma 7.1:

$$T_0 \geq \frac{\Xi_x^2 \ln^2 \frac{T}{q} \ln^2(\Xi_{\mathbf{Z}})}{C_1 \rho^2} \quad \text{and} \quad \eta_t = C_2 \cdot \begin{cases} \frac{1}{\sqrt{T_0} \Xi_x \ln \frac{T}{q}} & t \leq T_0; \\ \frac{1}{t \cdot \rho} & t > T_0. \end{cases},$$

After Warm Start. Our second lemma asks for a stronger assumption on the learning rates and shows that after warm start (i.e., for $t \geq T_0$), the quantity $(s'_t)^2$ scales essentially inversely to $1/t$.

Lemma 7.3 (after warm start). *In the same setting as Lemma 7.1, if there exists $\delta \leq 1/\sqrt{8}$ s.t.*

$$\frac{T_0}{\ln^2 T_0} \geq \frac{9 \ln(8/q^2)}{\delta^2}, \quad \forall s \in \{T_0+1, \dots, T\}: \quad 2\eta_s \rho - 56\eta_s^2 \Xi_x^2 \geq \frac{1}{s-1} \quad \text{and} \quad \eta_s \leq \frac{1}{20(s-1)\delta \Xi_x},$$

then, with probability at least $1 - 2qT$ (over the randomness of x_1, \dots, x_T):

- $\|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq 2\Xi_{\mathbf{Z}}$ for every $t \in \{T_0, \dots, T\}$, and
- $\|\mathbf{W}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq \frac{5T_0/\ln^2(T_0)}{t/\ln^2 t}$ for every $t \in \{T_0, \dots, T\}$.

Parameter 7.4. There exist constants $C_1, C_2, C_3 > 0$ such that for every $q > 0$ that is sufficiently small (meaning $q < 1/\text{poly}(T, \Xi_{\mathbf{Z}}, \Xi_x, 1/\rho)$), the following parameters satisfy both Lemma 7.1 and Lemma 7.3:

$$\frac{T_0}{\ln^2(T_0)} = C_1 \cdot \frac{\Xi_x^2 \ln^2 \frac{T}{q} \ln^2 \Xi_{\mathbf{Z}}}{\rho^2}, \quad \eta_t = C_2 \cdot \begin{cases} \frac{\ln \Xi_{\mathbf{Z}}}{T_0 \cdot \rho} & t \leq T_0; \\ \frac{1}{t \cdot \rho} & t > T_0. \end{cases}, \quad \text{and} \quad \delta = C_3 \cdot \frac{\rho}{\Xi_x}.$$

8 Putting Everything Together

Using our learning rates choices Parameter 7.4 and main lemmas in Section 7, it is not hard to

- prove exactly Theorem 2 (see Appendix I.1), and
- prove a weaker version of Theorem 1 where $\Lambda = \lambda_1 + \dots + \lambda_k$ is replaced 1.

Improvement 1. To further improve Theorem 1 so that the factor Λ shows up in the convergence (e.g., shows up in T_0), we need tighter martingale concentrations on our random variables and below we discuss the main intuition.

Recall that all martingale concentrations for a random process $\{z_t\}_t$ require some upper bound between consecutive variables $|z_t - z_{t+1}|$. If this upper bound is a probability-one absolute one, that is, $|z_t - z_{t+1}| \leq M$, then an Azuma-type of concentration can be proved. However, Azuma concentration is not tight: if one knows a better bound on $\mathbb{E}[|z_{t+1} - z_t|^2 \mid z_t]$, he or she can replace M^2 with this expected bound and get a tighter concentration. See for instance the survey [4].

The same issue also shows up in online PCA. Our Lemma 5.1-(b) corresponds to a probability-one absolute bound on $|z_t - z_{t+1}|$; if one replaces it with a tighter (but very sophisticated) expected bound, the concentration result can be further improved and this improvement translates to faster running time on Oja's algorithm (through our same framework used in Section 7). We present such expected bounds in Appendix F, and prove similar versions of Lemma 7.1 and Lemma 7.3 in Appendix G. Combining them one can obtain the exact statement of Theorem 1, and the final proof is included in Appendix I.2.

Remark 8.1. This factor Λ improvement is only possible in the gap-dependent case and does not show up in gap-free running times to the best of our knowledge.

Improvement 2. In order to prove Theorem 3 which is the rayleigh-quotient guarantee in gap-free online PCA, we want to strengthen Lemma 7.1 so that it provides guarantee essentially of the form:

$$\text{for every } \gamma \geq 1: \quad \|\mathbf{W}_\gamma^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq 2/\gamma, \quad (8.1)$$

where \mathbf{W}_γ is the column orthonormal matrix consisting of all eigenvectors of $\mathbf{\Sigma}$ with eigenvalues $\leq \lambda_k - \gamma \cdot \rho$. For obvious reason Lemma 7.1 is a special case of (8.1) when restricting only to $\gamma = 1$. It is a simple exercise to show that (8.1) implies our desired rayleigh-quotient guarantee (via an Abel transformation and an integral computation, see Appendix I.3).

Therefore, it suffices to prove (8.1). If one were allowed to magically change learning rates and apply Lemma 7.1 multiple times, then (8.1) would be trivial to prove: just replace \mathbf{W} with \mathbf{W}_γ and replacing ρ with $\gamma \cdot \rho$ and repeatedly apply Lemma 7.1. Unfortunately, the difficulty arises because we want to prove (8.1) for all $\gamma \geq 1$ but with a *fixed* set of learning rates η_t .

We proved in this paper that, using the same learning rates in Parameter 7.4, together with a more general martingale concentration lemma (i.e., Corollary 6.6 with $\gamma \geq 1$), one can obtain (8.1). This proof follows from the same structure as that of Lemma 7.1 except for the change in how we apply Corollary 6.6. We include the details in Appendix H.

APPENDIX

A Random Initialization (Missing Proofs for Section 4)

Proof of Lemma 4.1. Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}_\mathbf{A}\mathbf{U}^\top$ be the eigendecomposition of \mathbf{A} , and we denote by $\mathbf{Q}_z = \mathbf{Z}^\top \mathbf{Q} \mathbf{U} \in \mathbb{R}^{(d-k) \times d}$. Since a random Gaussian matrix is rotation invariant, and since \mathbf{U} is unitary and \mathbf{Z} is column orthonormal, we know that each entry of \mathbf{Q}_z draw i.i.d. from $\mathcal{N}(0, 1)$.

Next, since we have $\|\mathbf{Z}^\top x\|_2 \leq 1$, it satisfies that $y = x^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} \mathbf{U}$ is a vector with each coordinate i independently drawn from distribution $\mathcal{N}(0, \sigma_i)$ for $\sigma_i \leq 1$. This implies

$$x^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} \mathbf{A} \mathbf{Q}^\top \mathbf{Z} \mathbf{Z}^\top x = y^\top \mathbf{\Sigma}_\mathbf{A} y = \sum_{i=1}^k [\mathbf{\Sigma}_\mathbf{A}]_{i,i} (y_i)^2.$$

Now, $\sum_{i \in [k]} [\mathbf{\Sigma}_\mathbf{A}]_{i,i} (y_i)^2$ is a subexponential distribution⁹ with parameter (σ^2, b) where $\sigma^2, b \leq 4 \sum_{i=1}^k [\mathbf{\Sigma}_\mathbf{A}]_{i,i}$. Using the subexponential concentration bound, we have for every $\lambda \geq 1$,

$$\Pr \left[\sum_{i=1}^k [\mathbf{\Sigma}_\mathbf{A}]_{i,i} (y_i)^2 \geq \sum_{i=1}^k [\mathbf{\Sigma}_\mathbf{A}]_{i,i} + \lambda \right] \leq \exp \left\{ -\frac{\lambda}{8 \sum_{i=1}^k [\mathbf{\Sigma}_\mathbf{A}]_{i,i}} \right\}.$$

After rearranging, we have

$$\Pr[x^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} \mathbf{A} \mathbf{Q}^\top \mathbf{Z} \mathbf{Z}^\top x \geq \text{Tr}(\mathbf{A}) + \lambda] \leq e^{-\frac{\lambda}{8 \text{Tr}(\mathbf{A})}}. \quad \square$$

The following lemma is on the singular value distribution of a random Gaussian matrix:

Lemma A.1 (Theorem 1.2 of [17]). *Let $\mathbf{Q} \in \mathbb{R}^{k \times k}$ be a random matrix with each entry i.i.d. drawn from $\mathcal{N}(0, 1)$, and $\sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_k$ be its singular values. We have for every $j \in [k]$ and $\alpha \geq 0$:*

$$\Pr \left[\sigma_j \leq \frac{\alpha j}{\sqrt{k}} \right] \leq \left((2e)^{1/2} \alpha \right)^{j^2}.$$

Proof of Lemma 4.2. Using Lemma A.1, we know that

$$\begin{aligned} \Pr \left[\text{Tr} \left[\left((\mathbf{V}^\top \mathbf{Q})^\top (\mathbf{V}^\top \mathbf{Q}) \right)^{-1} \right] \geq \frac{\pi^2 ek}{3p} \right] &\leq \Pr \left[\exists j \in [k], \sigma_j^{-2} (\mathbf{V}^\top \mathbf{Q}) \geq \frac{2ek}{j^2 p} \right] \\ &= \Pr \left[\exists j \in [k], \sigma_j (\mathbf{V}^\top \mathbf{Q}) \leq \frac{j \sqrt{p}}{\sqrt{2ek}} \right] \leq \sum_{j=1}^k p^{j^2/2} \leq \frac{\sqrt{p}}{1-p}. \quad \square \end{aligned}$$

Proof of Lemma 4.3. Applying Lemma 4.2 with the choice of probability $= \frac{p^2}{4}$, we know that

$$\Pr_{\mathbf{Q}} \left[\text{Tr}(\mathbf{A}) \geq \frac{36k}{p^2} \right] \leq p \quad \text{where} \quad \mathbf{A} \stackrel{\text{def}}{=} \left((\mathbf{V}^\top \mathbf{Q})^\top (\mathbf{V}^\top \mathbf{Q}) \right)^{-1}.$$

Conditioning on event $\mathcal{C} = \left\{ \text{Tr}(\mathbf{A}) \leq \frac{36k}{p^2} \right\}$, and setting $r = \frac{36k}{p^2}$, we have for every fixed x_1, \dots, x_T

⁹Recall that a random variable X is (σ^2, b) -subexponential if $\log \mathbb{E} \exp(\lambda(X - \mathbb{E} X)) \leq \lambda^2 \sigma^2 / 2$ for all $\lambda \in [0, 1/b]$. The squared standard Gaussian variable is $(4, 4)$ -subexponential.

and fixed $i \in [T]$, it satisfies

$$\begin{aligned}
& \Pr_{\mathbf{Q}} \left[\left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \left(18r \ln \frac{T}{q} \right)^{1/2} \mid \mathcal{C}, x_t \right] \\
& \stackrel{\textcircled{1}}{\leq} \Pr \left[\left\| y_t \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \left(18r \ln \frac{T}{q} \right)^{1/2} \mid \mathcal{C}, x_1, \dots, x_t \right] \\
& \stackrel{\textcircled{2}}{\leq} \Pr \left[y_t^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{Q} \mathbf{A} \mathbf{Q}^\top \mathbf{Z} \mathbf{Z}^\top y_t \geq 9r \ln \frac{T^2}{q^2} \mid \mathcal{C}, x_1, \dots, x_t \right] \stackrel{\textcircled{3}}{\leq} \frac{q^2}{T^2} .
\end{aligned}$$

Above, $\textcircled{1}$ uses the definition $y_t \stackrel{\text{def}}{=} x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1}$; $\textcircled{2}$ is from the definition of \mathbf{A} ; and $\textcircled{3}$ is owing to Lemma 4.1 together with the fact that $\|y_t\|_2 \leq \|x_t\|_2 \cdot \left\| \left(\frac{\mathbf{Z} \mathbf{Z}^\top \Sigma}{\lambda_{k+1}} \right)^{i-1} \right\|_2 \leq 1$ and the fact that $\mathbf{Z}^\top \mathbf{Q}$ is independent of $\mathbf{V}^\top \mathbf{Q}$. Next, define event

$$\mathcal{C}_2 = \left\{ \exists i \in [T], \exists t \in [T], \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \left(18r \ln \frac{T}{q} \right)^{1/2} \right\} .$$

The above derivation, after taking union bound, implies that for every fixed x_1, \dots, x_T , it satisfies $\Pr_{\mathbf{Q}}[\mathcal{C}_2 \mid \mathcal{C}, x_1, \dots, x_T] \leq q^2$. Therefore, denoting by $\mathbb{1}_{\mathcal{C}_2}$ the indicator function of event \mathcal{C}_2 ,

$$\begin{aligned}
\Pr_{\mathbf{Q}} \left[\Pr_{x_1, \dots, x_T} [\mathcal{C}_2 \mid \mathbf{Q}] \geq q \mid \mathcal{C} \right] & \leq \frac{1}{q} \mathbb{E}_{\mathbf{Q}} \left[\Pr_{x_1, \dots, x_T} [\mathcal{C}_2 \mid \mathbf{Q}] \mid \mathcal{C} \right] \\
& = \frac{1}{q} \mathbb{E}_{\mathbf{Q}} \left[\mathbb{E}_{x_1, \dots, x_T} [\mathbb{1}_{\mathcal{C}_2} \mid \mathbf{Q}] \mid \mathcal{C} \right] \\
& = \frac{1}{q} \mathbb{E}_{x_1, \dots, x_T} \left[\mathbb{E}_{\mathbf{Q}} [\mathbb{1}_{\mathcal{C}_2} \mid \mathcal{C}, x_1, \dots, x_T] \right] \\
& = \frac{1}{q} \mathbb{E}_{x_1, \dots, x_T} \left[\Pr_{\mathbf{Q}} [\mathcal{C}_2 \mid \mathcal{C}, x_1, \dots, x_T] \right] \leq q .
\end{aligned}$$

Above, the first inequality uses Markov's bound. In an analogous manner, we define event

$$\mathcal{C}_3 = \left\{ \exists j \in [d], j \geq k+1, \|v_j^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_2 \geq \left(18r \ln \frac{d}{p} \right)^{1/2} \right\}$$

where v_j is the j -th eigenvector of Σ corresponding to eigenvalue λ_j . A completely analogous proof as the lines above also shows $\Pr_{\mathbf{Q}}[\mathcal{C}_3 \mid \mathcal{C}] \leq q$. Finally, using union bound

$$\Pr_{\mathbf{Q}} \left[\mathcal{C}_3 \bigwedge_{x_1, \dots, x_T} \Pr_{x_1, \dots, x_T} [\mathcal{C}_2 \mid \mathbf{Q}] \geq q \right] \leq \Pr_{\mathbf{Q}}[\mathcal{C}_3 \mid \mathcal{C}] + \Pr_{\mathbf{Q}} \left[\Pr_{x_1, \dots, x_T} [\mathcal{C}_2 \mid \mathbf{Q}] \geq q \mid \mathcal{C} \right] + \Pr_{\mathbf{Q}}[\mathcal{C}] \leq p + 2q ,$$

we conclude that with probability at least $1 - p - 2q$ over the random choice of \mathbf{Q} , it satisfies

- $\Pr_{x_1, \dots, x_T} [\mathcal{C}_2 \mid \mathbf{Q}] < q$, and
- $\overline{\mathcal{C}_3}$ holds (which implies $\|\mathbf{Z}^\top \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 < 18rd \ln \frac{d}{p}$ as desired).

□

B Expected Results (Missing Proofs for Section 5)

Proof of Lemma 5.1. We first notice that

$$\begin{aligned}
\mathbf{X}^\top \mathbf{P}_t \mathbf{Q} &= \mathbf{X}^\top \mathbf{P}_{t-1} \mathbf{Q} + \eta_t \mathbf{X}^\top x_t x_t^\top \mathbf{P}_{t-1} \mathbf{Q} \quad \text{and} \\
\mathbf{V}^\top \mathbf{P}_t \mathbf{Q} &= \mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q} + \eta_t \mathbf{V}^\top x_t x_t^\top \mathbf{P}_{t-1} \mathbf{Q} ,
\end{aligned}$$

where the second equality further implies (using the Sherman-Morrison formula) that

$$\begin{aligned} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1} &= (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1} - \frac{\eta_t (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1} \mathbf{V}^\top x_t x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}}{1 + \eta_t x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1} \mathbf{V}^\top x_t} \\ &= (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1} - (\eta_t - \alpha_t \eta_t^2) (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1} \mathbf{H}'_t, \end{aligned}$$

and above we denote by $\alpha_t \stackrel{\text{def}}{=} \frac{\psi_t}{1 + \eta_t \psi_t}$ where $\psi_t \stackrel{\text{def}}{=} x_t^\top \mathbf{L}_{t-1} \mathbf{V}^\top x_t$. Therefore, we can write

$$\begin{aligned} \mathbf{S}_t &= \mathbf{X}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1} \\ &= \mathbf{S}_{t-1} - (\eta_t - \alpha_t \eta_t^2) \mathbf{S}_{t-1} \mathbf{H}'_t + \eta_t \mathbf{R}'_t - (\eta_t^2 - \alpha_t \eta_t^3) \mathbf{R}'_t \mathbf{H}'_t \\ &= \mathbf{S}_{t-1} - (\eta_t - \alpha_t \eta_t^2) \mathbf{S}_{t-1} \mathbf{H}'_t + (\eta_t - \psi_t \eta_t^2 + \alpha_t \psi_t \eta_t^3) \mathbf{R}'_t = \mathbf{S}_{t-1} - \eta_t \mathbf{S}_{t-1} \mathbf{H}_t + \eta_t \mathbf{R}_t. \end{aligned}$$

Above, in the last equality we have denoted by $\mathbf{H}_t = (1 - \alpha_t \eta_t) \mathbf{H}'_t$ and $\mathbf{R}_t = (1 - \psi_t \eta_t + \alpha_t \psi_t \eta_t^2) \mathbf{R}'_t$ to simplify the notations. We now proceed and compute

$$\begin{aligned} \text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) &= \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}_t) \\ &\quad + \eta_t^2 \text{Tr}(\mathbf{H}_t^\top \mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}_t) + \eta_t^2 \text{Tr}(\mathbf{R}_t^\top \mathbf{R}_t) - 2\eta_t^2 \text{Tr}(\mathbf{R}_t^\top \mathbf{S}_{t-1} \mathbf{H}_t) \\ &\stackrel{\textcircled{1}}{\leq} \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}_t) \\ &\quad + 2\eta_t^2 \text{Tr}(\mathbf{H}_t^\top \mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}_t) + 2\eta_t^2 \text{Tr}(\mathbf{R}_t^\top \mathbf{R}_t) \\ &\stackrel{\textcircled{2}}{\leq} \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}_t) \\ &\quad + 2\eta_t^2 (1 - \alpha_t \eta_t)^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 2\eta_t^2 (1 - \psi_t \eta_t + \alpha_t \psi_t \eta_t^2)^2 \|\mathbf{R}'_t\|_2^2 \\ &\stackrel{\textcircled{3}}{\leq} \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\ &\quad + 2\eta_t^2 |\alpha_t| \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) \right| + 2\eta_t (\eta_t |\psi_t| + \eta_t^2 |\alpha_t| |\psi_t|) \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \right| \\ &\quad + 2\eta_t^2 (1 + 2\phi_t \eta_t)^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 2\eta_t^2 (1 + \phi_t \eta_t + 2\phi_t^2 \eta_t^2)^2 \|\mathbf{R}'_t\|_2^2 \\ &\stackrel{\textcircled{4}}{\leq} \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\ &\quad + 4\eta_t^2 \|\mathbf{H}'_t\|_2 \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) \right| + 4\eta_t^2 \|\mathbf{H}'_t\|_2 \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \right| \\ &\quad + 8\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 \\ &\stackrel{\textcircled{5}}{\leq} \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\ &\quad + 4\eta_t^2 \|\mathbf{H}'_t\|_2 \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \right| + 12\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2. \quad (\text{B.1}) \end{aligned}$$

Above, ① is because $2\text{Tr}(\mathbf{A}^\top \mathbf{B}) \leq \text{Tr}(\mathbf{A}^\top \mathbf{A}) + \text{Tr}(\mathbf{B}^\top \mathbf{B})$ which is Young's inequality in the matrix case; ② and ③ are both because $\mathbf{H}_t = (1 - \alpha_t \eta_t) \mathbf{H}'_t$ and $\mathbf{R}_t = (1 - \psi_t \eta_t + \alpha_t \psi_t \eta_t^2) \mathbf{R}'_t$; ④ follow from the parameter properties $|\psi_t| \leq \|\mathbf{H}'_t\|_2 \leq \phi_t$, $|\alpha_t| \leq 2\|\mathbf{H}'_t\|_2 \leq 2\phi_t$, and $0 \leq \eta_t \phi_t \leq \frac{1}{2}$; ⑤ follows from $|\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t)| \leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \|\mathbf{H}'_t\|_2$ which uses Proposition 2.1.

Next, Proposition 2.1 tells us

$$|\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| \leq \|\mathbf{R}'_t\|_{S_1} \|\mathbf{S}_{t-1}\|_2 \leq \|\mathbf{R}'_t\|_2 \sqrt{\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})} \leq \frac{\|\mathbf{R}'_t\|_2}{2} \left(\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 1 \right), \quad (\text{B.2})$$

(the second inequality is because \mathbf{R}'_t is rank 1, and the spectral norm of a matrix is no greater than

its Frobenius norm.) we can further simplify the upper bound in (B.1) as

$$\begin{aligned}
\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) &\leq \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\
&\quad + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2 \left(\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 1 \right) + 12\eta_t^2 \|\mathbf{H}'_t\|_2^2 \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 \\
&= \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\
&\quad + (12\eta_t^2 \|\mathbf{H}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2 .
\end{aligned}$$

This finishes the proof of Lemma 5.1-(a).

A completely symmetric analysis of the above derivation also gives

$$\begin{aligned}
\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) &\geq \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\
&\quad - (12\eta_t^2 \|\mathbf{H}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 - 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2 ,
\end{aligned}$$

and thus combining the upper and lower bounds we have

$$\begin{aligned}
|\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})| &\leq 2\eta_t |\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t)| + 2\eta_t |\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| \\
&\quad + (12\eta_t^2 \|\mathbf{H}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2
\end{aligned} \tag{B.3}$$

$$\begin{aligned}
&\stackrel{\textcircled{1}}{\leq} (2\eta_t \|\mathbf{H}'_t\|_2 + 12\eta_t^2 \|\mathbf{H}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 2\eta_t \|\mathbf{R}'_t\|_2 \sqrt{\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})} \\
&\quad + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 + 2\eta_t^2 \|\mathbf{R}'_t\|_2 \|\mathbf{H}'_t\|_2
\end{aligned} \tag{B.4}$$

$$\stackrel{\textcircled{2}}{\leq} 9\eta_t \|\mathbf{H}'_t\|_2 \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 2\eta_t \|\mathbf{R}'_t\|_2 \sqrt{\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})} + 10\eta_t^2 \phi_t \|\mathbf{R}'_t\|_2 . \tag{B.5}$$

Above, $\textcircled{1}$ again uses Proposition 2.1 and (B.2); $\textcircled{2}$ uses $\eta_t \phi_t \leq 1/2$ and $\|\mathbf{H}'_t\|_2, \|\mathbf{R}'_t\|_2 \leq \phi_t$.

Finally, if we take square on both sides of (B.5), we have (using again $\eta_t \|\mathbf{R}'_t\|_2 \leq \frac{1}{2}$):

$$|\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})|^2 \leq 243\eta_t^2 \|\mathbf{H}'_t\|_2^2 \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 12\eta_t^2 \|\mathbf{R}'_t\|_2^2 \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 300\eta_t^4 \phi_t^2 \|\mathbf{R}'_t\|_2^2$$

and this finishes the proof of Lemma 5.1-(b). If we continue to use $\|\mathbf{H}'_t\|_2, \|\mathbf{R}'_t\|_2 \leq \phi_t$ to upper bound the right hand side of (B.5), we finish the proof of Lemma 5.1-(c). \square

Proof of Corollary 5.2 from Lemma 5.1. According to the expectation we have $\mathbb{E}[\mathbf{H}'_t \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] = \mathbf{V}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{L}_{t-1}$ and $\mathbb{E}[\mathbf{R}'_t \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] = \mathbf{X}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{L}_{t-1}$. Now we consider the subcases separately:

(a) By Lemma 5.1-(a),

$$\begin{aligned}
\mathbb{E} \left[\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\stackrel{\textcircled{1}}{\leq} (1 + 14\eta_t^2 \phi_t^2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\eta_t^2 \phi_t^2 \\
&\quad - 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{L}_{t-1}) + 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{L}_{t-1}) \\
&\stackrel{\textcircled{2}}{\leq} (1 - 2\eta_t \text{gap} + 14\eta_t^2 \phi_t^2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\eta_t^2 \phi_t^2 \\
&\quad - 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \boldsymbol{\Delta} \mathbf{L}_{t-1}) + 2\eta_t \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \boldsymbol{\Delta} \mathbf{L}_{t-1})
\end{aligned} \tag{B.6}$$

Above, $\textcircled{1}$ uses $\|\mathbf{R}'_t\|_2, \|\mathbf{H}'_t\|_2 \leq \phi_t$, and $\textcircled{2}$ is because $\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \boldsymbol{\Sigma} \mathbf{L}_{t-1}) = \mathbf{Tr}(\mathbf{S}_{t-1}^\top \boldsymbol{\Sigma}_{>k} \mathbf{Z}^\top \mathbf{L}_{t-1}) = \mathbf{Tr}(\mathbf{S}_{t-1}^\top \boldsymbol{\Sigma}_{>k} \mathbf{S}_{t-1}) \leq \lambda_{k+1} \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})$, as well as $\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \boldsymbol{\Sigma} \mathbf{L}_{t-1}) = \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \boldsymbol{\Sigma}_{\leq k} \mathbf{V}^\top \mathbf{L}_{t-1}) = \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \boldsymbol{\Sigma}_{\leq k}) \geq \lambda_k \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})$.

Next, using the decomposition $\mathbf{I} = \mathbf{V}\mathbf{V}^\top + \mathbf{Z}\mathbf{Z}^\top$, $\|\mathbf{V}\|_2 \leq 1, \|\mathbf{Z}\|_2 \leq 1$, and Proposition 2.1

multiple times, we have

$$\begin{aligned}
\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{L}_{t-1}) &= \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1}) \\
&\leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{V}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{Z} \mathbf{Z} \mathbf{S}_{t-1}) \\
&\stackrel{\textcircled{1}}{\leq} \|\Delta\|_2 \left(\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + [\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{3/2} \right) \\
\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{L}_{t-1}) &= \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1}) \\
&\leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{V}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{Z} \mathbf{S}_{t-1}) \\
&\leq \|\Delta\|_2 \left(\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^{1/2} \right).
\end{aligned}$$

Above, $\textcircled{1}$ uses the fact that $\|\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top\|_{S_1} \leq \|\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top\|_{S_1} \|\mathbf{S}_{t-1}\|_2 \leq [\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{3/2}$. Plugging them into (B.6) finishes the proof of Corollary 5.2-(a).

- (b) In this case (B.6) also holds but one needs to replace gap with ρ because of the definitional difference between \mathbf{W} and \mathbf{Z} . We compute the following upper bounds similar to case (a):

$$\begin{aligned}
\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{L}_{t-1}) &= \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1}) \\
&\leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{V}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1}) \\
&\stackrel{\textcircled{1}}{\leq} \|\Delta\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \left(1 + [\text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2} \right) \\
\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{L}_{t-1}) &= \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1}) \\
&\leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{V}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top \Delta \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1}) \\
&\stackrel{\textcircled{2}}{\leq} \|\Delta\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^{1/2} \left(1 + \text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})^{1/2} \right) \quad (\text{B.7})
\end{aligned}$$

Above, $\textcircled{1}$ is because (using Proposition 2.1)

$$\begin{aligned}
\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1}) &\leq \text{Tr}((\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2)^{1/2} \cdot [\text{Tr}(\mathbf{V}^\top \Delta \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z} \mathbf{Z}^\top \Delta^\top \mathbf{V})]^{1/2} \\
&\leq \|\Delta\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \cdot [\text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2}
\end{aligned}$$

and $\textcircled{2}$ holds for a similar reason.

Putting these upper bounds into (B.6) finishes the proof of Corollary 5.2-(b).

- (c) When $\mathbf{X} = [w]$, a slightly different derivation of (B.6) gives

$$\begin{aligned}
\mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - 2\eta_t \lambda_k + 14\eta_t^2 \phi_t^2) \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 10\eta_t^2 \phi_t^2 \\
&\quad - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{L}_{t-1}) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top w^\top \Delta \mathbf{L}_{t-1}) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top w^\top \Sigma \mathbf{L}_{t-1}). \quad (\text{B.8})
\end{aligned}$$

Note that the third and fourth terms can be upper bounded similarly using (B.7). As for the fifth term, we have

$$\text{Tr}(\mathbf{S}_{t-1}^\top w^\top \Sigma \mathbf{L}_{t-1}) \leq \frac{\lambda_k}{2} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \frac{1}{2\lambda_k} \text{Tr}(w^\top \Sigma \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \Sigma w)$$

Putting these together, we have:

$$\begin{aligned}
\mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - \eta_t \lambda_k + 14\eta_t^2 \phi_t^2) \text{Tr}(\mathbf{S}_{t-1} \mathbf{S}_{t-1}^\top) + 10\eta_t^2 \phi_t^2 + \frac{\eta_t}{\lambda_k} \|w^\top \Sigma \mathbf{L}_{t-1}\|_2^2 \\
&\quad + 2\eta_t \|\Delta\|_2 \left([\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) \left(1 + [\text{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2} \right) \quad \square
\end{aligned}$$

C Martingale Concentrations (Missing Proofs for Section 6)

C.1 Proofs for One-Dimensional Martingale

Proof of Lemma 6.1. Define $y_t = \frac{\delta^2 t z_t}{\ln t} - \ln t$, then we have:

$$\begin{aligned} \mathbb{E}[y_{t+1} \mid \mathcal{F}_{\leq t}] &= \frac{\delta^2(t+1) \mathbb{E}[z_{t+1} \mid \mathcal{F}_{\leq t}]}{\ln(t+1)} - \ln(t+1) \\ &\leq \frac{\delta^2(t+1)(1 - \delta\tau_t)z_t}{\ln(t+1)} + \frac{\delta^2(t+1)\tau_t^2}{\ln(t+1)} - \ln(t+1) \\ &\leq \frac{\delta^2(t+1)(1 - \frac{1}{t})}{\ln(t+1)} z_t + \frac{t+1}{t^2 \ln(t+1)} - \ln(t+1) \stackrel{\textcircled{1}}{\leq} \frac{\delta^2 t z_t}{\ln t} - \ln t = y_t, \end{aligned}$$

where $\textcircled{1}$ is because for every $t \geq 4$ it satisfies $\frac{(t+1)(t-1)}{\ln(t+1)} \leq \frac{t^2}{\ln t}$ and $\frac{t+1}{t^2 \ln(t+1)} \leq \ln(1 + \frac{1}{t})$.

At the same time, we have

$$|y_{t+1} - y_t| \stackrel{\textcircled{2}}{\leq} \frac{\delta^2 t}{\ln t} |z_{t+1} - z_t| + \frac{\delta^2}{\ln t} z_{t+1} + \frac{1}{t}, \quad (\text{C.1})$$

where $\textcircled{2}$ is because for every $t \geq 3$ it satisfies $0 \leq \frac{t+1}{\ln(t+1)} - \frac{t}{\ln t} \leq \frac{1}{\ln t}$ and $\ln(t+1) - \ln(t) \leq 1/t$.

Taking square on both sides, we have

$$|y_{t+1} - y_t|^2 \leq 3 \left(\frac{\delta^2 t}{\ln t} \right)^2 |z_{t+1} - z_t|^2 + 3 \left(\frac{\delta^2}{\ln t} \right)^2 z_{t+1}^2 + \frac{3}{t^2}.$$

Taking expectation on both sides, we have

$$\begin{aligned} \mathbb{E}[|y_{t+1} - y_t|^2 \mid \mathcal{F}_{\leq t}] &\leq 3 \left(\frac{\delta^2 t}{\ln t} \right)^2 (\tau_t^2 z_t + \kappa^2 \tau_t^4) + 3 \frac{(y_t + \ln t)^2}{t^2} + \frac{3}{t^2} \\ &< \frac{3(y_t + \ln t)}{t \ln t} + \frac{3(y_t + \ln t)^2}{t^2} + \frac{3(1 + \kappa^2)}{t^2} \\ &\stackrel{\textcircled{3}}{\leq} \frac{3(\phi + 1)}{t} + \frac{3(\phi + 1)^2 \ln^2 t}{t^2} + \frac{15\kappa^2}{4t^2} \stackrel{\textcircled{4}}{\leq} \frac{4(\phi + 1)}{t}. \end{aligned}$$

Above, $\textcircled{3}$ uses $y_t \leq \phi \ln t$ and $\kappa \geq 2$; $\textcircled{4}$ uses $\frac{t}{\ln^2 t} \geq \frac{t_0}{\ln^2 t_0} \geq \max\{7.5\kappa^2, 6(\phi + 1)\}$ and $\ln t \geq 1$.

Therefore, if $y_t \leq \phi \ln t$ holds true for $t = t_0, \dots, T$ and $t_0 \geq 8$ (which implies $\frac{t}{\ln^2 t} \geq \frac{t_0}{\ln^2 t_0}$), then

$$\sum_{t=t_0}^T \mathbb{E}[|y_{t+1} - y_t|^2 \mid \mathcal{F}_{\leq t}] \leq \sum_{t=t_0}^T \frac{4(\phi + 1)}{t} \leq 4(\phi + 1) \int_{t=t_0-1}^T \frac{dt}{t} \leq 4(\phi + 1) \ln(T).$$

Now we can check about the absolute difference. We continue from (C.1) and derive that, if $y_t \leq \phi \ln t$, then

$$\begin{aligned} |y_{t+1} - y_t| &\leq \frac{\delta^2 t}{\ln t} |z_{t+1} - z_t| + \frac{\delta^2}{\ln t} z_{t+1} + \frac{1}{t} \leq \frac{\delta^2 t}{\ln t} (\kappa \tau_t \sqrt{z_t} + \kappa^2 \tau_t^2) + \frac{\delta^2}{\ln t} z_{t+1} + \frac{1}{t} \\ &\leq \kappa \left(\sqrt{\frac{y_t + \ln t}{t \ln t}} + \frac{\kappa}{t \ln t} + \frac{(y_t + \ln t)}{t} + \frac{1}{t} \right) \stackrel{\textcircled{5}}{\leq} \kappa \left(\sqrt{\frac{y_t + \ln t}{t \ln t}} + \frac{y_t + \ln t + \kappa}{t} \right) \\ &\stackrel{\textcircled{6}}{\leq} \kappa \left(\sqrt{\frac{(\phi + 1)}{t}} + \frac{(\phi + 1) \ln t + \kappa}{t} \right) \stackrel{\textcircled{7}}{\leq} 2\kappa \sqrt{\frac{(\phi + 1)}{t}} \end{aligned}$$

where $\textcircled{5}$ uses $\ln t \geq 2$ and $\kappa \geq 2$, $\textcircled{6}$ uses $y_t \leq \phi \ln t$, and $\textcircled{7}$ uses $\frac{t}{\ln^2 t} \geq \frac{t_0}{\ln^2 t_0} \geq 4 \max\{\phi + 1, \kappa\}$.

From the above inequality, we have that if $t_0 \geq 4\kappa^2(\phi + 1)$ and $y_t \leq \phi \ln t$ holds true for

$t = t_0, \dots, T-1$ then $|y_{t+1} - y_t| \leq 1$ for all $t = t_0, \dots, T-1$.

Finally, since we have assumed $\phi > 36$ and $z_{t_0} \leq \frac{\phi \ln^2 t_0}{2\delta^2 t_0}$ which implies $y_{t_0} \leq \frac{\phi \ln t_0}{2}$, we can apply martingale concentration inequality (c.f. [4, Theorem 18]):

$$\begin{aligned}
\Pr[\exists t \geq t_0, y_t > \phi \ln t] &\leq \sum_{T=t_0+1}^{\infty} \Pr[y_T > \phi \ln T; \forall t \in \{t_0, \dots, T-1\}, y_t \leq \phi \ln t] \\
&\leq \sum_{T=t_0+1}^{\infty} \Pr[y_T - y_{t_0} > \phi \ln T/2; \forall t \in \{t_0, \dots, T-1\}, y_t \leq \phi \ln t] \\
&\leq \sum_{T=t_0+1}^{\infty} \exp \left\{ \frac{-(\phi \ln T/2)^2}{2 \cdot 4(\phi+1) \ln(T-1) + \frac{2}{3}(\phi \ln T/2)} \right\} \\
&\leq \sum_{T=t_0+1}^{\infty} \exp \left\{ -\ln T \frac{\phi^2/4}{8(\phi+1) + \phi/3} \right\} \\
&\leq \int_{T=t_0}^{\infty} \exp \left\{ -\frac{\phi}{36} \ln T \right\} dT \leq \frac{\exp\{-(\frac{\phi}{36} - 1) \ln t_0\}}{\frac{\phi}{36} - 1}.
\end{aligned}$$

□

Proof of Corollary 6.3. Define $\phi \stackrel{\text{def}}{=} \frac{4\delta^2 t_0}{\ln^2 t_0} \geq 36 \ln \frac{1}{p} \geq 72$. It is easy to verify that $\frac{t_0}{\ln^2 t_0} \geq 7.5\kappa^2(\phi+1)$ (because $\kappa \leq 1/(\sqrt{2}\delta)$) and $z_{t_0} \leq \frac{\phi \ln^2 t_0}{2\delta^2 t_0} = 2$, so we can apply Lemma 6.1:

$$\Pr \left[\exists t \geq t_0, z_t > \frac{(\phi+1) \ln^2 t}{\delta^2 t} \right] \leq \frac{\exp \left\{ -(\frac{\phi}{36} - 1) \ln t_0 \right\}}{\frac{\phi}{36} - 1} \leq \exp \left\{ -\left(\frac{\phi}{36} - 1 \right) \ln t_0 \right\} \leq p,$$

where the last inequality uses $\ln t_0 \geq 2$ and $(\frac{\phi}{36} - 1) \ln t_0 \geq \frac{\phi}{36}$. Therefore, we conclude that

$$\Pr \left[\exists t \geq t_0, z_t > \frac{5(t_0/\ln^2 t_0)}{t/\ln^2 t} \right] \leq \Pr \left[\exists t \geq t_0, z_t > \frac{(\phi+1) \ln^2 t}{\delta^2 t} \right] \leq p.$$

□

C.2 Proofs for Multi-Dimensional Martingale

Proof of Corollary 6.4. We apply Lemma 6.2 with $\lambda = 2 \max \{1, \max_{j \in [t+1]} \{[z_0]_j\}\} \geq 2$. Using the fact that $\beta_t \geq 0$, we know that

$$\Pr[z_t]_1 \geq \lambda = \Pr \left[[z_t]_1 \geq 2 \left(\max_{j \in [t+1]} \{[z_0]_j\} + 1 \right) \right] \leq (1 + 1.4t) \exp \left\{ -\ln(2^p) + 5p^2 \sum_{s=0}^{t-1} \tau_s^2 \right\}$$

Denoting by $\alpha = \sum_{s=0}^{t-1} \tau_s^2$, we can take $p = \frac{1}{6\sqrt{\alpha}} \leq \min_{s \in [t]} \left\{ \frac{1}{6\kappa\tau_{s-1}} \right\}$ satisfying the assumption of Lemma 6.2. Therefore,

$$\Pr[z_t]_1 \geq \lambda \leq 4t \exp \left\{ -\frac{1}{9\sqrt{\alpha}} + \frac{5}{36} \right\} \leq q,$$

where the last inequality requires $\frac{1}{\sqrt{\alpha}} \geq 9 \left(\ln \frac{4t}{q} + \frac{5}{36} \right)$ which can be satisfied under our assumption $\alpha \leq \frac{1}{100} \ln^{-2} \frac{4t}{q}$. □

Proof of Corollary 6.5. We apply Lemma 6.2 with $\lambda = 2 \max \{1, \max_{j \in [t+1]} \{[z_0]_j\}\} \geq 2$, and $p = 2 \ln \frac{4t}{q} = \frac{l}{6}$. Note that the presumption $p \leq \min_{s \in [t]} \left\{ \frac{1}{6\kappa\tau_{s-1}} \right\}$ is satisfied because $\kappa\tau_s l \leq 1$.

The conclusion of Lemma 6.2 tells us that, since $p < \frac{l}{5}$ and $\beta_s \geq l\tau_s^2$ which together imply $\beta_t \geq 5p\tau_t^2$, we have

$$\mathbf{Pr} \left[[z_t]_1 \geq 2 \left(\max_{j \in [t+1]} \{[z_0]_j\} + 1 \right) \right] \leq (1 + 1.4t) \exp \{-p \ln 2\} \leq 4te^{-p \ln 2} < q . \quad \square$$

Proof of Corollary 6.6. We consider fixed $p = \frac{l}{5\gamma} = 2 \ln \frac{3t}{q}$. Let $y_t = \gamma \cdot z_t$, then y_t satisfies (6.2) with (using the fact that $\gamma \geq 1$)

$$\beta'_t = \beta_t, \quad \delta'_t = \delta_t, \quad (\tau'_t)^2 = \gamma\tau_t^2, \quad \kappa' = \kappa .$$

We denote by $b \stackrel{\text{def}}{=} \sum_{s=0}^{t-1} \beta_s = b$ and $a \stackrel{\text{def}}{=} \sum_{s=0}^{t-1} \tau_s^2$, and apply Lemma 6.2 on y_t with $\lambda = 2$. Using the fact that $\beta_s \geq l\tau_s^2 = 5zp\tau_s^2$ we know $p\beta'_t \geq 5p^2(\tau'_t)^2$. Therefore, for all $s \in \{0, 1, \dots, T-1\}$ we have

$$\mathbf{Pr} [[y_t]_1 \geq 2] \leq \exp \{-pb + 5p^2\gamma a + p \ln \Xi - p \ln 2\} + 1.4t \exp \{-p \ln 2\} , \quad (\text{C.2})$$

where we have denoted by $\Xi \stackrel{\text{def}}{=} \max_{j \in [t+1]} \{[z_0]_j\}$ for notational simplicity. Now, the choice $p = 2 \ln \frac{3t}{q}$ satisfies the presumption of Lemma 6.2 because we have assumed $\kappa\tau_s \leq \frac{1}{12 \ln \frac{3t}{q}}$. Therefore, we have

$$\begin{aligned} -pb + 5p^2\gamma a + p \ln \Xi - p \ln 2 &= p(-b + la + \ln \Xi - \ln 2) \leq \ln \frac{q}{2} \iff b - la \geq \ln \Xi \bigwedge p \geq 2 \ln \frac{3t}{q} \\ -p \ln 2 &\leq \ln \frac{q}{3t} \iff p \geq 2 \ln \frac{3t}{q} . \end{aligned}$$

Plugging them into (C.2) gives $\mathbf{Pr} \left[[z_t]_1 \geq \frac{2}{\gamma} \right] = \mathbf{Pr} [[y_t]_1 \geq 2] \leq \frac{q}{2} + \frac{q}{2} = q . \quad \square$

Proof of Lemma 6.2. Define vector s_t for every $t \in \{0, 1, \dots, T-1\}$ and $i \in [D]$, it satisfies $[s_t]_i \stackrel{\text{def}}{=} \frac{[z_{t+1}]_i}{[z_t]_i} - 1$. We have

$$\mathbb{E} [[s_t]_i \mid \mathcal{F}_{\leq t}] \leq -(\delta_t + \beta_t - \tau_t^2) + \delta_t \frac{[z_t]_{i+1}}{[z_t]_i} + \frac{\tau_t^2}{[z_t]_i} . \quad (\text{C.3})$$

In particular,

$$\text{if } [z_t]_i \geq 1, \text{ then } \mathbb{E} [[s_t]_i^2 \mid \mathcal{F}_{\leq t}] \leq \tau_t^2 + \frac{\tau_t^2}{[z_t]_i} + \frac{\kappa^2 \tau_t^4}{[z_t]_i^2} \leq (2 + (\tau_t \kappa)^2) \tau_t^2 \leq 3\tau_t^2 , \quad (\text{C.4})$$

$$|[s_t]_i| \leq \kappa\tau_t + \frac{\kappa\tau_t}{\sqrt{[z_t]_i}} + \frac{\kappa^2 \tau_t^2}{[z_t]_i} \leq \kappa\tau_t(2 + \kappa\tau_t) \leq 3\kappa\tau_t . \quad (\text{C.5})$$

We consider $[z_{t+1}]_i^p$ for some fixed value $p \geq 1$ and derive that (using (C.5))

$$\begin{aligned} \text{if } (\kappa\tau_t)p \leq \frac{1}{6} \text{ and } [z_t]_i \geq 1, \text{ then } [z_{t+1}]_i^p &= [z_t]_i^p (1 + [s_t]_i)^p = [z_t]_i^p \left(\sum_{q=0}^p \binom{p}{q} [s_t]_i^q \right) \\ &\leq [z_t]_i^p (1 + p[s_t]_i + p^2[s_t]_i^2) . \end{aligned}$$

After taking expectation, we have if $(\kappa\tau_t)p \leq \frac{1}{6}$ and $[z_t]_i \geq 1$, then

$$\begin{aligned}
\mathbb{E} [[z_{t+1}]_i^p \mid \mathcal{F}_{\leq t}] &\stackrel{\textcircled{1}}{\leq} [z_t]_i^p (1 + p \mathbb{E} [[s_t]_i \mid \mathcal{F}_{\leq t}] + 3p^2\tau_t^2) \\
&\stackrel{\textcircled{2}}{\leq} [z_t]_i^p \left(1 - p(\delta_t + \beta_t - \tau_t^2) + \delta_t p \frac{[z_t]_{i+1}}{[z_t]_i} + \frac{\tau_t^2 p}{[z_t]_i} + 3p^2\tau_t^2 \right) \\
&= [z_t]_i^p (1 - p(\delta_t + \beta_t - \tau_t^2) + 3p^2\tau_t^2) + \delta_t p [z_t]_i^{p-1} [z_t]_{i+1} + p\tau_t^2 [z_t]_i^{p-1} \\
&\stackrel{\textcircled{3}}{\leq} [z_t]_i^p (1 - p(\delta_t + \beta_t - \tau_t^2) + 3p^2\tau_t^2 + p\tau_t^2) + \delta_t p \left(\frac{p-1}{p} [z_t]_i^p + \frac{1}{p} [z_t]_{i+1}^p \right) \\
&= [z_t]_i^p (1 - \delta_t - p\beta_t + p\tau_t^2 + 3p^2\tau_t^2 + p\tau_t^2) + \delta_t [z_t]_{i+1}^p \\
&\stackrel{\textcircled{4}}{\leq} [z_t]_i^p (1 - \delta_t - p\beta_t + 5p^2\tau_t^2) + \delta_t [z_t]_{i+1}^p .
\end{aligned}$$

Above, ① uses (C.4); ② uses (C.3); ③ uses $[z_t]_i \geq 1$ and Young's inequality $ab \leq a^p/p + b^q/q$ for $1/p + 1/q = 1$; and ④ uses $p \geq 1$.

On the other hand, if $(\kappa\tau_t)p \leq \frac{1}{6}$ but $[z_t]_i < 1$, we have the following simple bound (using $\kappa\tau_t \leq 1/6$):

$$[z_{t+1}]_i \leq (1 + \kappa\tau_t)[z_t]_i + \kappa\tau_t \sqrt{[z_t]_i} + \kappa^2\tau_t^2 \leq (1 + \kappa\tau_t) + (\kappa\tau_t) + \kappa^2\tau_t^2 < 1.4 .$$

Therefore, as long as $(\kappa\tau_t)p \leq \frac{1}{6}$ we always have

$$\mathbb{E} [[z_{t+1}]_i^p \mid \mathcal{F}_{\leq t}] \leq [z_t]_i^p (1 - \delta_t - p\beta_t + 5p^2\tau_t^2) + \delta_t [z_t]_{i+1}^p + 1.4 =: (1 - \alpha_t)[z_t]_i^p + \delta_t [z_t]_{i+1}^p + 1.4 ,$$

and in the last inequality we have denoted by $\alpha_t \stackrel{\text{def}}{=} \delta_t + p\beta_t - 5p^2\tau_t^2$. Telescoping this expectation, and choosing $i = 1$, we have whenever $p \in [1, \min_{s \in [t]} \{\frac{1}{6\kappa\tau_{s-1}}\}]$, it satisfies

$$\begin{aligned}
\mathbb{E} [[z_{t+1}]_1^p] &\leq \prod_{s=1}^t (1 - \alpha_s + \delta_s) \left(\max_{j \in [t+2]} \{[z_0]_j^p\} \right) + 1.4 \sum_{s=0}^t \left(\prod_{u=s+1}^t (1 - \alpha_u + \delta_u) \right) \\
&\leq \prod_{s=0}^t (1 - p\beta_s + 5p^2\tau_s^2) \left(\max_{j \in [t+2]} \{[z_0]_j^p\} \right) + 1.4 \sum_{s=0}^t \left(\prod_{u=s+1}^t (1 - p\beta_u + 5p^2\tau_u^2) \right) \\
&\leq \max_{j \in [t+2]} \{[z_0]_j^p\} \exp \left\{ -p \left(\sum_{s=0}^t \beta_s \right) + 5p^2 \sum_{s=0}^t \tau_s^2 \right\} + 1.4 \sum_{s=0}^t \exp \left\{ -p \left(\sum_{u=s+1}^t \beta_u \right) + 5p^2 \sum_{u=s+1}^t \tau_u^2 \right\} .
\end{aligned}$$

Finally, using Markov's inequality, we have for every $\lambda > 0$:

$$\begin{aligned}
\Pr [[z_{t+1}]_1 \geq \lambda] &\leq \lambda^{-p} \left(\max_{j \in [t+2]} \{[z_0]_j^p\} \exp \left\{ \sum_{s=0}^t 5p^2\tau_s^2 - p\beta_s \right\} \right. \\
&\quad \left. + 1.4 \sum_{s=0}^t \exp \left\{ \sum_{u=s+1}^t 5p^2\tau_u^2 - p\beta_u \right\} \right) . \quad \square
\end{aligned}$$

D Decoupling Lemmas

We prove the following general lemma. Let $x_1, \dots, x_T \in \Omega$ be random variables each i.i.d. drawn from some distribution \mathcal{D} . Let \mathcal{F}_t be the sigma-algebra generated by x_t , and denote by $\mathcal{F}_{\leq t} = \vee_{s=1}^t \mathcal{F}_s$.¹⁰

Lemma D.1 (decoupling lemma). *Consider a fixed value $q \in [0, 1)$. For every $t \in [T]$ and $s \in \{0, 1, \dots, t-1\}$, let $y_{t,s} \in \mathbb{R}^D$ be an $\mathcal{F}_t \vee \mathcal{F}_{\leq s}$ measurable random vector and let $\phi_{t,s} \in \mathbb{R}^D$ be a fixed*

¹⁰For the purpose of this paper, one can feel free view Ω as \mathbb{R}^d , each x_t as the t -th sample vector, and \mathcal{D} as the distribution with covariance matrix Σ .

vector. Let $D' \in [D]$. Define events (we denote by $^{(i)}$ the i -th coordinate)

$$\begin{aligned} \mathcal{C}'_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_{t-1}) \text{ satisfies } \Pr_{x_t} \left[\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{t-1} \right] \leq q \right\} \\ \mathcal{C}''_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_t) \text{ satisfies } \forall i \in [D'] : y_{t,t-1}^{(i)} \leq \phi_{t,t-1}^{(i)} \right\} \end{aligned}$$

and denote by $\mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$ and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$. Suppose the following three assumptions hold:

(A1) The random process $\{y_{t,s}\}_{t,s}$ satisfy that for every $i \in [D], t \in [T-1], s \in \{0, 1, \dots, t-2\}$

- (a) $\mathbb{E} [y_{t,s+1}^{(i)} \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s}] \leq f_s^{(i)}(y_{t,s}, q),$
- (b) $\mathbb{E} [|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}|^2 \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s}] \leq h_s^{(i)}(y_{t,s}, q),$ and
- (c) $|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}| \leq g_s^{(i)}(y_{t,s})$ whenever $\mathcal{C}_{\leq s}$ holds.

Above, for each $i \in [D]$ and $s \in \{0, 1, \dots, T-2\}$, we have $f_s, h_s : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}_{\geq 0}^D$, $g_s : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}^D$ are functions satisfying for every $x \in \mathbb{R}^d$,

- (d) $f_s^{(i)}(x, p), h_s^{(i)}(x, p)$ are monotone increasing in p , and
- (e) $|x^{(i)} - f_s^{(i)}(x, 0)|^2 \leq h_s^{(i)}(x, 0)$ and $|x^{(i)} - f_s^{(i)}(x, 0)| \leq g_s^{(i)}(x)$ whenever $f_s^{(i)}(x, 0) \leq x^{(i)}.$

(A2) Each $t \in [T]$ satisfies $\Pr_{x_t}[\bar{\mathcal{E}}_t] \leq q^2/2$ where event

$$\mathcal{E}_t \stackrel{\text{def}}{=} \{x_t \text{ satisfies } \forall i \in [D] : y_{t,0}^{(i)} \leq \phi_{t,0}^{(i)}\}.$$

(A3) For every $t \in [T]$, letting x_t be any vector satisfying \mathcal{E}_t , consider **any** random process $\{z_s\}_{s=0}^{t-1}$ where each $z_s \in \mathbb{R}_{\geq 0}^D$ is $\mathcal{F}_{\leq s}$ measurable with $z_0 = y_{t,0}$ as the starting vector. Suppose that whenever $\{z_s\}_{s=0}^{t-1}$ satisfies

$$\forall i \in [D], \forall s \in \{0, 1, \dots, t-2\} : \left\{ \begin{array}{ll} \mathbb{E} [z_{s+1}^{(i)} \mid \mathcal{F}_{\leq s}] & \leq f_s^{(i)}(z_s, q) \\ \mathbb{E} [|z_{s+1}^{(i)} - z_s^{(i)}|^2 \mid \mathcal{F}_{\leq s}] & \leq h_s^{(i)}(z_s, q) \\ |z_{s+1}^{(i)} - z_s^{(i)}| & \leq g_s^{(i)}(z_s) \end{array} \right\} \quad (\text{D.1})$$

then it holds $\Pr_{x_1, \dots, x_{t-1}}[\exists i \in [D'] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq q^2/2$.

Under the above two assumptions, we have for every $t \in [T]$, it satisfies $\Pr[\bar{\mathcal{C}}_t] \leq 2tq$.

Proof of Lemma D.1. We prove the lemma by induction. For the base case, by applying assumption (A2) we know that $\Pr_{x_1}[\exists i \in [D'] : y_{1,0}^{(i)} > \phi_{1,0}^{(i)}] \leq \Pr[\bar{\mathcal{E}}_1] \leq q^2/2 \leq q$ so event \mathcal{C}_1 holds with probability at least $1 - q$. In other words, $\Pr[\bar{\mathcal{C}}_{\leq 1}] = \Pr[\bar{\mathcal{C}}_1] \leq q < 2q$.

Suppose $\Pr[\bar{\mathcal{C}}_{\leq t-1}] \leq 2(t-1)q$ is true for some $t \geq 2$, we will prove $\Pr[\bar{\mathcal{C}}_{\leq t}] \leq 2tq$. Since it satisfies $\Pr[\bar{\mathcal{C}}_{\leq t}] \leq \Pr[\bar{\mathcal{C}}_{\leq t-1}] + \Pr[\bar{\mathcal{C}}_t]$, it suffices to prove that $\Pr[\bar{\mathcal{C}}_t] \leq 2q$.

Note also $\Pr[\bar{\mathcal{C}}_t] \leq \Pr[\bar{\mathcal{C}}'_t] + \Pr[\bar{\mathcal{C}}''_t \mid \mathcal{C}'_t]$ but the second quantity $\Pr[\bar{\mathcal{C}}''_t \mid \mathcal{C}'_t]$ is no more than q according to our definition of \mathcal{C}'_t and \mathcal{C}''_t . Therefore, in the rest of the proof, it suffices to show $\Pr[\bar{\mathcal{C}}'_t] \leq q$.

We use $y_{t,s}(x_t, x_{\leq s})$ to emphasize that $y_{t,s}$ is an $\mathcal{F}_t \times \mathcal{F}_{\leq s}$ measurable random vector. Let us now fix x_t to be a vector satisfying \mathcal{E}_t . Define $\{z_s\}_{s=0}^{t-1}$ to be a random process where each $z_s \in \mathbb{R}^D$ is $\mathcal{F}_{\leq s}$ measurable:

$$z_s^{(i)} = z_s^{(i)}(x_{\leq s}) \stackrel{\text{def}}{=} \begin{cases} y_{t,s}^{(i)}(x_t, x_{\leq s}) & \text{if } x_{\leq s} \text{ satisfies } \mathcal{C}_{\leq s}; \\ \min \left\{ f_{s-1}^{(i)}(z_{s-1}(x_{\leq s-1}), 0), z_{s-1}^{(i)}(x_{\leq s-1}) \right\} & \text{if } x_{\leq s} \text{ satisfies } \bar{\mathcal{C}}_{\leq s}. \end{cases} \quad (\text{D.2})$$

Then $z_s^{(i)}$ satisfies for every $i \in [D], s \leq \{0, 1, \dots, t-2\}$,

$$\begin{aligned}
\mathbb{E}[z_{s+1}^{(i)} \mid \mathcal{F}_{\leq s}] &= \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[z_{s+1}^{(i)} \mid \mathcal{C}_{\leq s+1}, \mathcal{F}_{\leq s}] + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[z_{s+1}^{(i)} \mid \overline{\mathcal{C}_{\leq s+1}}, \mathcal{F}_{\leq s}] \\
&\stackrel{\textcircled{1}}{\leq} \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[y_{t,s+1}^{(i)} \mid \mathcal{C}_{\leq s+1}, \mathcal{F}_{\leq s}] + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot f_s^{(i)}(z_s, 0) \\
&\stackrel{\textcircled{2}}{\leq} \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot f_s^{(i)}(y_{t,s}, q) + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot f_s^{(i)}(z_s, q) \\
&\stackrel{\textcircled{3}}{\leq} \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot f_s^{(i)}(y_{t,s}, q) + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot f_s^{(i)}(y_{t,s}, q) \tag{D.3} \\
&= f_s^{(i)}(z_s, q) \tag{D.4}
\end{aligned}$$

Above, $\textcircled{1}$ is because whenever $\mathcal{C}_{\leq s+1}$ holds it satisfies $z_{s+1}^{(i)} = y_{t,s+1}^{(i)}$, as well as whenever $\overline{\mathcal{C}_{\leq s+1}}$ holds it satisfies $z_{s+1}^{(i)} \leq f_s^{(i)}(z_s, 0)$; $\textcircled{2}$ uses assumptions (A1a) and (A1d) as well as the fact that we have fixed x_t ; $\textcircled{3}$ uses the fact that whenever $\Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] > 0$ it must hold that $\mathcal{C}_{\leq s}$ is satisfied, and therefore it satisfies $y_{t,s} = z_s$.

Similarly, we can also show for every $i \in [D], s \leq \{0, 1, \dots, t-2\}$,

$$\begin{aligned}
&\mathbb{E}[|z_{s+1}^{(i)} - z_s^{(i)}|^2 \mid \mathcal{F}_{\leq s}] \\
&= \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[|z_{s+1}^{(i)} - z_s^{(i)}|^2 \mid \mathcal{C}_{\leq s+1}, \mathcal{F}_{\leq s}] + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[|z_{s+1}^{(i)} - z_s^{(i)}|^2 \mid \overline{\mathcal{C}_{\leq s+1}}, \mathcal{F}_{\leq s}] \\
&\stackrel{\textcircled{1}}{\leq} \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot \mathbb{E}[|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}|^2 \mid \mathcal{C}_{\leq s+1}, \mathcal{F}_{\leq s}] + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot h_s^{(i)}(z_s, 0) \\
&\stackrel{\textcircled{2}}{\leq} \Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] \cdot h_s^{(i)}(y_{t,s}, q) + \Pr[\overline{\mathcal{C}_{\leq s+1}} \mid \mathcal{F}_{\leq s}] \cdot h_s^{(i)}(z_s, q) \\
&\stackrel{\textcircled{3}}{\leq} h_s^{(i)}(z_s, q) . \tag{D.5}
\end{aligned}$$

Above, $\textcircled{1}$ is because whenever $\mathcal{C}_{\leq s+1}$ holds it satisfies $z_{s+1}^{(i)} = y_{t,s+1}^{(i)}$ and $z_s^{(i)} = y_{t,s}^{(i)}$, together with whenever $\overline{\mathcal{C}_{\leq s+1}}$ holds it satisfies $|z_{s+1}^{(i)} - y_s^{(i)}|^2$ either equal zero or equal $|f_s^{(i)}(z_s, 0) - z_s^{(i)}|^2$, but in the latter case we must have $f_s^{(i)}(z_s, 0) < z_s^{(i)}$ (owing to (D.2)) and therefore it holds $|f_s^{(i)}(z_s, 0) - z_s^{(i)}|^2 \leq h_s^{(i)}(z_s, 0)$ using assumption (A1e). $\textcircled{2}$ uses assumptions (A1b) and (A1d) as well as the fact that we have fixed x_t . $\textcircled{3}$ uses the fact that whenever $\Pr[\mathcal{C}_{\leq s+1} \mid \mathcal{F}_{\leq s}] > 0$ then $\mathcal{C}_{\leq s}$ must hold, and therefore it satisfies $y_{t,s} = z_s$.

Finally, we also have

$$|z_{s+1}^{(i)} - z_s^{(i)}| \leq g_s^{(i)}(z_s^{(i)}) . \tag{D.6}$$

This is so because whenever $\mathcal{C}_{\leq s+1}$ holds it satisfies $|z_{s+1}^{(i)} - z_s^{(i)}| = |y_{t,s+1}^{(i)} - y_{t,s}^{(i)}|$ so we can apply assumption (A1c). Otherwise, $\overline{\mathcal{C}_{\leq s+1}}$ holds we either have $|z_{s+1}^{(i)} - z_s^{(i)}| = 0$ (so (D.6) trivially holds) or $|z_{s+1}^{(i)} - z_s^{(i)}| = |f_s^{(i)}(z_s, 0) - z_s^{(i)}|$, but in the latter case we must have $f_s^{(i)}(z_s, 0) < z_s^{(i)}$ (owing to (D.2)) so it must satisfy $|f_s^{(i)}(z_s, 0) - z_s^{(i)}| \leq g_s^{(i)}(z_s)$ using assumption (A1e).

We are now ready to apply assumption (A3), which together with (D.4), (D.5), (D.6), implies that (recalling we have fixed x_t to be any vector satisfying \mathcal{E}_t)

$$\Pr_{x_1, \dots, x_{t-1}} [\exists i \in [D'] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{E}_t] \leq q^2/2 .$$

This implies, after translating back to the random process $\{y_{t,s}\}$, we have

$$\begin{aligned} \Pr_{x_1, \dots, x_t} [\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)}] &\leq \Pr_{x_1, \dots, x_t} [\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{E}_t] + \Pr[\overline{\mathcal{E}_t}] \\ &\leq \Pr_{x_1, \dots, x_{t-1}} [\exists i \in [D'] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{E}_t] + q^2/2 \\ &\leq q^2/2 + q^2/2 = q^2. \end{aligned}$$

where the last inequality uses (A2). Finally, using Markov's inequality,

$$\begin{aligned} \Pr_{x_1, \dots, x_{t-1}} [\overline{\mathcal{C}'_t}] &= \Pr_{x_1, \dots, x_{t-1}} \left[\Pr_{x_t} [\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{\leq t-1}] > q \right] \\ &\leq \frac{1}{q} \cdot \mathbb{E}_{x_1, \dots, x_{t-1}} \left[\Pr_{x_t} [\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{\leq t-1}] \right] \\ &= \frac{1}{q} \cdot \Pr_{x_1, \dots, x_t} [\exists i \in [D'] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq q. \end{aligned}$$

Therefore, we finish proving $\Pr[\overline{\mathcal{C}'_t}] \leq q$ which implies $\Pr[\overline{\mathcal{C}_{\leq t}'}] \leq 2tq$ as desired. This finishes the proof of Lemma D.1. \square

E Main Lemmas (Missing Proofs for Section 7)

E.1 Before Warm Start

Proof of Lemma 7.1. For every $t \in [T]$ and $s \in \{0, 1, \dots, t-1\}$, consider random vectors $y_{t,s} \in \mathbb{R}^{T+2}$ defined as:

$$\begin{aligned} y_{t,s}^{(1)} &\stackrel{\text{def}}{=} \|\mathbf{Z}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2, \\ y_{t,s}^{(2)} &\stackrel{\text{def}}{=} \|\mathbf{W}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2, \\ y_{t,s}^{(3+j)} &\stackrel{\text{def}}{=} \begin{cases} \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^j \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1} \right\|_2^2, & \text{for } j \in \{0, 1, \dots, t-s-1\}; \\ (1 - \eta_s \lambda_k) \cdot y_{t,s-1}^{(3+j)}, & \text{for } j \in \{t-s, \dots, T-1\}. \end{cases} \end{aligned}$$

(In fact, we are only interested in $y_{t,s}^{(3+j)}$ for $j \leq t-s-1$, and can “almost” define $y_{t,s}^{(3+j)} = +\infty$ whenever $j \geq t-s$. However, we still decide to give such out-of-boundary variables meaningful values in order to make all of our vectors $y_{t,s}$ (and functions f, g, h defined later) to be of the same dimension $T+2$. This allows us to greatly simplify our notations.)

We consider upper bounds

$$\phi_{t,s}^{(1)} \stackrel{\text{def}}{=} 2\Xi_{\mathbf{Z}}, \quad \phi_{t,s}^{(2)} \stackrel{\text{def}}{=} \begin{cases} 2\Xi_{\mathbf{Z}} & s < T_0; \\ 2 & \text{otherwise.} \end{cases}, \quad \text{and } \phi_{t,s}^{(3+j)} \stackrel{\text{def}}{=} 2\Xi_x^2.$$

For each $t \in [T]$, define event \mathcal{C}'_t and \mathcal{C}''_t in the same way as decoupling Lemma D.1 (with $D' = 3$):

$$\begin{aligned} \mathcal{C}'_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_{t-1}) \text{ satisfies } \Pr_{x_t} [\exists i \in [3] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{t-1}] \leq q \right\} \\ \mathcal{C}''_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_t) \text{ satisfies } \forall i \in [3] : y_{t,t-1}^{(i)} \leq \phi_{t,t-1}^{(i)} \right\} \end{aligned}$$

and denote by $\mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$ and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$.

As a result, if $\mathcal{C}_{\leq s+1}$ holds, then we always have

$$\begin{aligned} \|x_{s+1}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2^2 &\leq (\|x_{s+1}^\top \mathbf{V} \mathbf{V}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2 + \|x_{s+1}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2)^2 \\ &\leq (1 + \phi_{s+1,s}^{(3)})^2 = (\sqrt{2}\Xi_x + 1)^2 \leq 4\Xi_x^2, \end{aligned}$$

where last inequality uses $\Xi_x \geq 2$. This allows us to later apply Corollary 5.2 with $\phi_t = 2\Xi_x$.

Verification of Assumption (A1) in Lemma D.1.

Suppose $\mathbb{E}[x_s x_s^\top \mid \mathcal{C}_{\leq s}, \mathcal{F}_{\leq s-1}] = \mathbf{\Sigma} + \mathbf{\Delta}$, and we want to bound $\|\mathbf{\Delta}\|_2$. Defining $q_1 \stackrel{\text{def}}{=} \Pr[\overline{\mathcal{C}}'_s \mid \mathcal{C}'_s, \mathcal{C}_{\leq s-1}, \mathcal{F}_{\leq s-1}]$, then we must have $q_1 \leq q$ according to the definition of \mathcal{C}'_s and $\overline{\mathcal{C}}'_s$. Using law of total expectation:

$$\mathbb{E}[x_s x_s^\top \mid \mathcal{C}'_s, \mathcal{C}_{\leq s-1}, \mathcal{F}_{\leq s-1}] = \mathbb{E}[x_s x_s^\top \mid \mathcal{C}_{\leq s}, \mathcal{F}_{\leq s-1}] \cdot (1 - q_1) + \mathbb{E}[x_s x_s^\top \mid \overline{\mathcal{C}}'_s, \mathcal{C}'_s, \mathcal{C}_{\leq s-1}, \mathcal{F}_{\leq s-1}] \cdot q_1,$$

and combining it with the fact that $0 \preceq x_s x_s^\top \preceq \mathbf{I}$ and $\mathbb{E}[x_s x_s^\top \mid \mathcal{C}'_s, \mathcal{C}_{\leq s-1}, \mathcal{F}_{\leq s-1}] = \mathbb{E}[x_s x_s^\top] = \mathbf{\Sigma}$, we have¹¹

$$\mathbf{\Sigma} \preceq (\mathbf{\Sigma} + \mathbf{\Delta})(1 - q_1) + q_1 \cdot \mathbf{I} \quad \text{and} \quad \mathbf{\Sigma} \succeq (\mathbf{\Sigma} + \mathbf{\Delta})(1 - q_1).$$

After rearranging, these two properties imply $\|\mathbf{\Delta}\|_2 \leq \frac{q_1}{1-q_1} \leq \frac{q}{1-q}$.

Now, we can apply Corollary 5.2 and obtain for every $t \in [T]$, $s \in \{0, 1, \dots, t-2\}$, and every $j \in \{0, 1, \dots, T-1\}$, it satisfies¹²

$$\begin{aligned} \mathbb{E}[y_{t,s+1}^{(1)} \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s+1}] &\leq (1 + 56\eta_{s+1}^2 \Xi_x^2) y_{t,s}^{(1)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}, \\ \mathbb{E}[y_{t,s+1}^{(2)} \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s+1}] &\leq (1 - 2\eta_{s+1}\rho + 56\eta_{s+1}^2 \Xi_x^2) y_{t,s}^{(2)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}, \text{ and} \\ \mathbb{E}[y_{t,s+1}^{(3+j)} \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s+1}] &\leq (1 - \eta_{s+1}\lambda_k + 56\eta_{s+1}^2 \Xi_x^2) y_{t,s}^{(3+j)} + \eta_{s+1}\lambda_k y_{t,s}^{(3+j+1)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}. \end{aligned}$$

Moreover, for every $i \in [T+2]$, using Lemma 5.1-(c) with $\phi_t = 2\Xi_x$ we have whenever $\mathcal{C}_{\leq s+1}$ holds it satisfies

$$|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}| \leq 18\eta_{s+1}\Xi_x \cdot y_{t,s}^{(i)} + 4\eta_{s+1}\Xi_x \cdot \sqrt{y_{t,s}^{(i)}} + 40\eta_{s+1}^2 \Xi_x^2 \leq 20\eta_{s+1}\Xi_x \cdot y_{t,s}^{(i)} + 42\eta_{s+1}^2 \Xi_x^2.$$

Putting the above bounds together, one can verify that the random process $\{y_{t,s}\}_{t \in [T], s \leq t-1}$ satisfy assumption (A1) of Lemma D.1 with¹³

$$\begin{aligned} f_s^{(1)}(y, q) &= (1 + 56\eta_{s+1}^2 \Xi_x^2) y^{(1)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}, \\ f_s^{(2)}(y, q) &= (1 - 2\eta_{s+1}\rho + 56\eta_{s+1}^2 \Xi_x^2) y^{(2)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}, \\ f_s^{(3+j)}(y, q) &= (1 - \eta_{s+1}\lambda_k + 56\eta_{s+1}^2 \Xi_x^2) y^{(3+j)} + \eta_{s+1}\lambda_k y^{(3+j+1)} + 40\eta_{s+1}^2 \Xi_x^2 + 20\eta_{s+1} \frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}, \\ g_s^{(i)}(y) &= 20\eta_{s+1}\Xi_x \cdot y^{(i)} + 42\eta_{s+1}^2 \Xi_x^2, \text{ and} \\ h_s^{(i)}(y, q) &= (g_s^{(i)}(y))^2 \end{aligned}$$

¹¹Here, we use notation $\mathbf{A} \preceq \mathbf{B}$ to indicate spectral dominance: that is, $\mathbf{B} - \mathbf{A}$ is positive semidefinite.

¹²To verify these upper bounds, one needs to use $\|\mathbf{Z}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2 \leq 2\Xi_{\mathbf{Z}}$ which is included in event $\mathcal{C}_{\leq s+1}$. Also, whenever $j \geq t-s$ so $y_{t,s}^{(3+j)}$ is out of boundary, we also have $y_{t,s}^{(3+j)} \leq (1 - \eta_s \lambda_k) \cdot y_{t,s-1}^{(3+j)}$ and it satisfies all the upper bounds.

¹³The only part of (A1) that is non-trivial to verify is (A1e) for $g_s^{(i)}$. Whenever $f_s^{(i)}(x, 0) \leq x^{(i)}$, it satisfies

$$|f_s^{(i)}(x, 0) - x^{(i)}| \leq \begin{cases} 0, & \text{if } i = 1; \\ 2\eta_{s+1}\rho \cdot x^{(2)}, & \text{if } i = 2; \\ \eta_{s+1}\lambda_k \cdot x^{(i)}, & \text{if } i \geq 3. \end{cases} \leq 2\eta_{s+1} \cdot x^{(i)} \leq g_s^{(i)}(x),$$

where the second inequality uses $\rho, \lambda_k \leq 1$ and the last inequality uses $\Xi_x \geq 2$.

Verification of Assumption (A2) of Lemma D.1.

For coordinates $i = 1$ and $i = 2$, our assumption $\|\mathbf{Z}^\top \mathbf{Q}(\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \Xi_{\mathbf{Z}}$ implies $y_{t,0}^{(i)} \leq \Xi_{\mathbf{Z}} < \phi_{t,0}^{(i)}$. For coordinates $i \geq 3$, we have assumption $\Pr_{x_t} \left[\forall j \in [T], \|x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma/\lambda_{k+1})^{j-1} \mathbf{Q}(\mathbf{V}^\top \mathbf{Q})^{-1}\|_2 \leq \Xi_x \right] \geq 1 - q^2/2$. Together, event \mathcal{E}_t (recall $\mathcal{E}_t \stackrel{\text{def}}{=} \{x_t \text{ satisfies } \forall i \in [D]: y_{t,0}^{(i)} \leq \phi_{t,0}^{(i)}\}$) holds for all $t \in [T]$ with probability at least $1 - q^2/2$. In sum, assumption (A2) is satisfied in Lemma D.1.

Verification of Assumption (A3) of Lemma D.1.

For every $t \in [T]$, at a high level assumption (A3) is satisfied once we plug in the following three sets of parameter choices to Corollary 6.4 and Corollary 6.6: for every $s \in [T-1]$, define

$$\begin{aligned} \beta_{s,1} &= 0, & \delta_{s,1} &= 0, & \tau_{s,1} &= 20\eta_{s+1}\Xi_x \\ \beta_{s,2} &= 2\eta_{s+1}\rho, & \delta_{s,2} &= 0, & \tau_{s,2} &= 20\eta_{s+1}\Xi_x \\ \beta_{s,3} &= 0, & \delta_{s,3} &= \eta_{s+1}\lambda_k, & \tau_{s,3} &= 20\eta_{s+1}\Xi_x \end{aligned}$$

More specifically, for every $t \in [T]$, let $\{z_s\}_{s=0}^{t-1}$ be the *arbitrary* random vector satisfying (D.1) of Lemma D.1. Define $q_2 = q^2/8$.

- For coordinate $i = 1$ of $\{z_s\}_{s=0}^{t-1}$,
 - apply Corollary 6.4 with $\{\beta_{s,1}, \delta_{s,1}, \tau_{s,1}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, and $\kappa = 1$;
- For coordinate $i = 2$ of $\{z_s\}_{s=0}^{t-1}$,
 - if $t < T_0$, apply Corollary 6.4 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, and $\kappa = 1$;
 - if $t \geq T_0$, apply Corollary 6.6 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, $\gamma = 1$, and $\kappa = 1$;
- For coordinates $i = 3, 4, \dots, T+2$ of $\{z_s\}_{s=0}^{t-1}$,
 - apply Corollary 6.4 with $\{\beta_{s,3}, \delta_{s,3}, \tau_{s,3}\}_{s=0}^{t-2}$, $q = q_2$, $D = T$, and $\kappa = 1$.

One needs to verify that the assumptions of Corollary 6.4 and 6.6 are satisfied as follows. First of all, one can carefully check that our parameters β, δ, τ satisfy (6.2) with $\kappa = 1$ and this needs our assumption $q \leq \frac{\eta_{s+1}}{\Xi_{\mathbf{Z}}^{3/2}}$. Next, we can apply Corollary 6.4 because we have assumed $\sum_{s=0}^{T-1} \tau_{s,1}^2 \leq \frac{1}{100} \ln^{-2} \frac{4T}{q_2}$. To verify the presumption of Corollary 6.6 with $\gamma = 1$, we notice that

- our assumption $\eta_s \leq \frac{\rho}{4000 \cdot \Xi_x^2 \ln \frac{3T}{q_2}}$ implies $\beta_{s,2} \geq 10 \ln \frac{3T}{q_2} \cdot \tau_{s,2}^2$ and $\kappa \tau_s \leq \frac{1}{12 \ln \frac{3T}{q_2}}$ for every s ,
- our assumption $\sum_{s=0}^{T_0-1} \beta_{s,2} \geq 1 + \ln \Xi_{\mathbf{Z}}$ implies $\sum_{s=0}^{t-1} \beta_s - 10 \ln \frac{3t}{q_2} \tau_s^2 \geq \ln \Xi_{\mathbf{Z}} + 1 - 1 = \ln \Xi_{\mathbf{Z}}$ whenever $t > T_0$,

Therefore, the conclusion of Corollary 6.4 and Corollary 6.6 imply that

$$\Pr[\exists i \in [3] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq 3q_2 < q^2/2$$

so assumption (A3) of Lemma D.1 holds.

Application of Lemma D.1. Applying Lemma D.1, we have $\Pr[\overline{\mathcal{C}_T}] \leq 2qT$ which implies our desired bounds and this finishes the proof of Lemma 7.1. \square

E.2 After Warm Start

Proof of Lemma 7.3. For every $t \in [T]$ and $s \in \{0, 1, \dots, t-1\}$, consider *the same* random vectors $y_{t,s} \in \mathbb{R}^{T+2}$ defined in the proof of Lemma 7.1:

$$\begin{aligned} y_{t,s}^{(1)} &\stackrel{\text{def}}{=} \|\mathbf{Z}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2, \\ y_{t,s}^{(2)} &\stackrel{\text{def}}{=} \|\mathbf{W}^\top \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2, \\ y_{t,s}^{(3+j)} &\stackrel{\text{def}}{=} \begin{cases} \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^j \mathbf{P}_s \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1} \right\|_2^2, & \text{for } j \in \{0, 1, \dots, t-s-1\}; \\ (1 - \eta_s \lambda_k) \cdot y_{t,s-1}^{(3+j)}, & \text{for } j \in \{t-s, \dots, T-1\}. \end{cases} \end{aligned}$$

This time, we consider slightly different upper bounds

$$\phi_{t,s}^{(1)} \stackrel{\text{def}}{=} 2\Xi_{\mathbf{Z}}, \quad \phi_{t,s}^{(2)} \stackrel{\text{def}}{=} \begin{cases} 2\Xi_{\mathbf{Z}} & \text{if } s < T_0; \\ 2 & \text{if } s = T_0; \\ \frac{5T_0/\ln^2(T_0)}{s/\ln^2 s} & \text{if } s > T_0. \end{cases}, \quad \text{and } \phi_{t,s}^{(3+j)} \stackrel{\text{def}}{=} 2\Xi_x^2.$$

We stress that the only difference between the above upper bounds and the ones we used in the proof of Lemma 7.1 is the choice of $\phi_{t,s}^{(2)}$ for $s > T_0$. Instead of setting it to be constant 2 for all such s , we make it decrease almost linearly with respect to index s .

Again, define event

$$\begin{aligned} \mathcal{C}'_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_{t-1}) \text{ satisfies } \Pr_{x_t} [\exists i \in [3] : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{t-1}] \leq q \right\} \\ \mathcal{C}''_t &\stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_t) \text{ satisfies } \forall i \in [3] : y_{t,t-1}^{(i)} \leq \phi_{t,t-1}^{(i)} \right\} \end{aligned}$$

and denote by $\mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$ and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$.

We next want to apply the decoupling Lemma D.1.

Verification of Assumption (A1) in Lemma D.1.

The same functions $f_s^{(i)}$, $g_s^{(i)}$, and $h_s^{(i)}$ used in the proof of Lemma 7.1 still apply here. However, we want to make a minor change on $g_s^{(i)}$ whenever $s \geq T_0$.

Applying Lemma 5.1-(c) with $\phi_t = 2\Xi_x$, we have whenever $\mathcal{C}_{\leq s+1}$ holds for some $s \geq T_0$ (which implies $y_{t,s}^{(2)} \leq 5$),

$$|y_{t,s+1}^{(2)} - y_{t,s}^{(2)}| \leq 18\eta_{s+1}\Xi_x y_{t,s}^{(2)} + 4\eta_{s+1}\Xi_x \sqrt{y_{t,s}^{(2)}} + 40\eta_{s+1}^2 \Xi_x^2 \leq 45\eta_{s+1}\Xi_x \sqrt{y_{t,s}^{(2)}} + 40\eta_{s+1}^2 \Xi_x^2.$$

Therefore, we can choose

$$g_s^{(2)}(y) = 45\eta_{s+1}\Xi_x \sqrt{y^{(2)}} + 40\eta_{s+1}^2 \Xi_x^2$$

for all $s \geq T_0$ and this still satisfies assumption (A1) of Lemma D.1.¹⁴

Verification of Assumption (A2) of Lemma D.1.

This is the same as the proof of Lemma 7.1.

Verification of Assumption (A3) in Lemma D.1.

Again, for every $t \in [T]$, let $\{z_s\}_{s=0}^{t-1}$ be the *arbitrary* random vector satisfying (D.1) of Lemma D.1. Choosing $q_2 = q^2/8$ again, the same proof of Lemma 7.1 shows that

$$\Pr[\exists i \in \{1, 3\} : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq 2q_2.$$

¹⁴Similar to Footnote 13, we also need to verify (A1e) for $g_s^{(i)}$. Whenever $f_s^{(2)}(x, 0) \leq x^{(2)}$, it satisfies $|f_s^{(2)}(x, 0) - x^{(2)}| \leq 2\eta_{s+1} \cdot x^{(2)} \leq g_s^{(2)}(x)$, where the first inequality uses $\rho \leq 1$ and the second uses $\Xi_x \geq 2$.

Therefore, it suffices to prove that $\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)}] \leq 2q_2$.

We only need to focus on the case $t \geq T_0 + 2$, because otherwise if $t \leq T_0 + 1$ then $g_s^{(2)}$ is not changed for all $s \in \{0, \dots, t-2\}$ so the same proof of Lemma 7.1 also shows $\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)}] \leq q_2$.

When $t \geq T_0 + 2$, we can first apply the same proof of Lemma 7.1 (for $t = T_0 + 1$) to show that $\Pr[z_{T_0}^{(2)} > \phi_{T_0+1,T_0}^{(2)} = 2] \leq q_2$. Next, conditioning on $z_{T_0}^{(2)} \leq 2$ which happens with probability at least $1 - q_2$, we want to apply Corollary 6.3 with $\kappa = 2$ and $\tau_s = \frac{1}{\delta s}$.

More specifically, for every $t \in \{T_0 + 2, \dots, T\}$, we have shown that the random sequence $\{z_s^{(2)}\}_{s=T_0}^{t-1}$ satisfies (D.1) with

$$\begin{aligned} f_s^{(2)}(y, q) &= (1 - 2\eta_{s+1}\rho + 56\eta_{s+1}^2\Xi_x^2)y^{(2)} + 40\eta_{s+1}^2\Xi_x^2 + 20\eta_{s+1}\frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q} \\ g_s^{(2)}(y) &= 45\eta_{s+1}\Xi_x\sqrt{y^{(2)}} + 40\eta_{s+1}^2\Xi_x^2 \\ h_s^{(2)}(y, q) &= (g_s^{(2)}(y))^2 \end{aligned}$$

Therefore, $\{z_s^{(2)}\}_{s=T_0}^{t-1}$ also satisfies (6.1) with $\kappa = 2$ and $\tau_s = \frac{1}{\delta s}$ because the following holds from our assumptions:

$$\begin{aligned} q\Xi_{\mathbf{Z}}^{3/2} &\leq \eta_{s+1} & \delta\tau_s &= \frac{1}{s} \leq 2\eta_{s+1}\rho - 56\eta_{s+1}^2\Xi_x^2 \\ \tau_s^2 &= \frac{1}{\delta^2 s^2} \geq 60\eta_{s+1}^2\Xi_x^2 \geq 40\eta_{s+1}^2\Xi_x^2 + 20\eta_{s+1}\frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q} & \kappa\tau_s &= \frac{2}{\delta s} \geq 40\eta_{s+1}\Xi_x \end{aligned}$$

Now, we are ready to apply Corollary 6.3 with $q = q_2$, $t_0 = T_0$, and $\kappa = 2$. Because $q_2 \leq e^{-2}$, $z_{T_0}^{(2)} \leq 2$, $\delta \leq 1/\sqrt{8}$ and $\frac{T_0}{\ln^2 T_0} \geq \frac{9\ln(1/q_2)}{\delta^2}$, the conclusion of Corollary 6.3 tells us

$$\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)} \mid z_{T_0}^{(2)} \leq 2] \leq q_2.$$

By union bound, we have $\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)}] \leq q_2 + q_2 = 2q_2$ as desired.

Finally, we conclude (for every $t \geq T_0 + 2$) that

$$\Pr[\exists i \in [3] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq 4q_2 < q^2/2$$

so assumption (A3) of Lemma D.1 holds.

Application of Lemma D.1. Applying Lemma D.1, we have $\Pr[\overline{\mathcal{C}}_T] \leq 2qT$ which implies our desired bounds and this finishes the proof of Lemma 7.3. \square

F Improvement: Expectation Lemmas

Lemma F.1. For every $t \in [T]$, For every $t \in [T]$, let $\mathcal{C}_{\leq t}$ be any event that depends on random x_1, \dots, x_t and implies

$$\|x_t^\top \mathbf{L}_{t-1}\|_2 = \|x_t^\top \mathbf{P}_{t-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{t-1} \mathbf{Q})^{-1}\|_2 \leq \phi_t \quad \text{where} \quad \eta_t \phi_t \leq \frac{1}{2},$$

and $\mathbb{E}[x_t x_t^\top \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] = \mathbf{\Sigma} + \mathbf{\Delta}$. Denote by $\Gamma \stackrel{\text{def}}{=} \min\{\sum_{i=1}^k \lambda_i + \|\mathbf{\Delta}\|_2, 1\}$. We have:

(a) If $\mathbf{X} = [w] \in \mathbb{R}^{d \times 1}$ where w is a vector with Euclidean norm at most 1,

$$\begin{aligned} \mathbb{E} \left[\mathbf{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - \eta_t \lambda_k + 14\Gamma \eta_t^2 \phi_t^2) \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\Gamma \eta_t^2 \phi_t^2 + \frac{\eta_t}{\lambda_k} \|w^\top \mathbf{\Sigma} \mathbf{L}_{t-1}\|_2^2 \\ &\quad + 2\eta_t \|\mathbf{\Delta}\|_2 \left([\mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} + \mathbf{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) \left(1 + [\mathbf{Tr}(\mathbf{Z}^\top \mathbf{L}_{t-1} \mathbf{L}_{t-1}^\top \mathbf{Z})]^{1/2} \right) \end{aligned}$$

(b) If $\mathbf{X} = [w] \in \mathbb{R}^{d \times 1}$ where w is a vector with Euclidean norm at most 1,

$$\begin{aligned} \mathbb{E} \left[\left| \text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right|^2 \mid \mathcal{F}_{t-1}, \mathcal{C}_{\leq t} \right] \\ \leq 243\Gamma\eta_t^2\phi_t^2\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 12\Gamma\eta_t^2\phi_t^2\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 300\Gamma\eta_t^4\phi_t^4 \end{aligned}$$

(c) If $\mathbf{X} = \mathbf{W} = \mathbf{Z}$,

$$\begin{aligned} \mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{t-1}, \mathcal{C}_{\leq t} \right] \\ \leq (1 - 2\eta_t \text{gap} + 12\Gamma\eta_t^2\phi_t^2 + \eta_t^2(6\phi_t + 8)\lambda_{k+1}) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\Gamma\eta_t^2(2\phi_t + 8) \\ + 2\eta_t \|\Delta\|_2 \left(\eta_t(4 + \phi_t)k + [\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{3/2} + (5 + 4\eta_t)\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + [\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})]^{1/2} \right). \end{aligned}$$

(d) If $\mathbf{X} = \mathbf{W} = \mathbf{Z}$,

$$\begin{aligned} \mathbb{E}[\|\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})\|_2^2 \mid \mathcal{F}_{t-1}, \mathcal{C}_{\leq t}] \\ \leq 192\Gamma\eta_t^2\phi_t^2\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 4\eta_t^2(6\phi_t + 10)^2\lambda_{k+1}\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \\ + 192\Gamma\eta_t^4\phi_t^2 + \|\Delta\|_2 \cdot 4\eta_t^2(6\phi_t + 10)^2(k + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})) . \end{aligned}$$

Proof. The proof of the first two cases rely on the follow tighter upper bound when $\mathbf{X} = [w]$:

$$\mathbb{E}[\|\mathbf{H}'_t\|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}], \mathbb{E}[\|\mathbf{R}'_t\|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] \leq \min \left\{ \left(\sum_{i=1}^k \lambda_i + \|\Delta\|_2 \right), 1 \right\} \phi_t^2 = \Gamma\phi_t^2 \quad (\text{F.1})$$

as opposed to ϕ_t^2 that we have used in the past.

The proof of the last two cases rely on the following tighter upper bounds when $\mathbf{X} = \mathbf{Z} = \mathbf{W}$:

$$\begin{aligned} \mathbb{E}[\|\mathbf{H}'_t\|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] &\leq \phi_t^2 \cdot \mathbb{E}[\|x_t^\top \mathbf{V}\|_F^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] \leq \Gamma\phi_t^2 \\ \mathbb{E}[\|\mathbf{R}'_t\|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] &\leq \mathbb{E}[\|x_t^\top \mathbf{L}_{t-1}\|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] = \text{Tr}(\mathbf{L}_{t-1}^\top \Sigma \mathbf{L}_{t-1}) + \text{Tr}(\mathbf{L}_{t-1}^\top \Delta \mathbf{L}_{t-1}) \\ &\leq \text{Tr}(\mathbf{L}_{t-1}^\top (\mathbf{V} \Sigma_{\leq k} \mathbf{V}^\top + \mathbf{Z} \Sigma_{> k} \mathbf{Z}^\top) \mathbf{L}_{t-1}) + \|\Delta\|_2 \text{Tr}(\mathbf{L}_{t-1}^\top (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1}) \\ &\leq \Lambda + \lambda_{k+1} \|\mathbf{Z}^\top \mathbf{L}_{t-1}\|_F^2 + \|\Delta\|_2 \cdot (k + \text{Tr}(\mathbf{L}_{t-1}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1})) \\ &= \Lambda + \lambda_{k+1} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \|\Delta\|_2 \cdot (k + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})) . \end{aligned} \quad (\text{F.2})$$

(a) This follows from almost the same proof of Corollary 5.2-(c), except that one can replace the use of (B.8) with the following (owing to (F.1))

$$\begin{aligned} \mathbb{E} \left[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t} \right] &\leq (1 - 2\eta_t \lambda_k + 14\Gamma\eta_t^2\phi_t^2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 10\Gamma\eta_t^2\phi_t^2 \\ &\quad - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{V}^\top \Delta \mathbf{L}_{t-1}) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top w^\top \Delta \mathbf{L}_{t-1}) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top w^\top \Sigma \mathbf{L}_{t-1}) . \end{aligned}$$

(b) This follows directly from Lemma 5.1-(b) and (F.1).

(c) We first note that (B.1) implies

$$\begin{aligned} \text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) &\leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) - 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \\ &\quad + 4\eta_t^2 \phi_t \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \right| + 12\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 . \end{aligned} \quad (\text{F.3})$$

This time, we upper bound

$$\begin{aligned} |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| &= |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{L}_{t-1})| = |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top (\mathbf{V} \mathbf{V}^\top + \mathbf{Z} \mathbf{Z}^\top) \mathbf{L}_{t-1})| \\ &= |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{V})| \\ &\leq \frac{3}{2} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) + \frac{1}{2} \|x_t^\top \mathbf{V}\|_2^2 . \end{aligned} \quad (\text{F.4})$$

Denoting by $\Lambda = \sum_{i=1}^k \lambda_i \leq \Gamma$, we can take expectation and get:

$$\begin{aligned} \mathbb{E} [|\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| \mid \mathcal{F}_{t-1}, \mathcal{C}_{\leq t}] &\leq \frac{3}{2} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{Z} \mathbf{S}_{t-1}) + \frac{1}{2} \text{Tr}(\mathbf{V}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{V}) \\ &\leq \frac{3}{2} \lambda_{k+1} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \frac{1}{2} \Lambda + \|\boldsymbol{\Delta}\|_2 \cdot \left(\frac{3}{2} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \frac{k}{2} \right) \end{aligned} \quad (\text{F.5})$$

At this point, plugging (F.5), (F.2) into (F.3) and using the assumption $\eta_t \phi_t \leq 1/2$, we have

$$\begin{aligned} \mathbb{E}[\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] &\leq (1 + 12\Gamma \eta_t^2 \phi_t^2 + \eta_t^2 (6\phi_t + 8) \lambda_{k+1}) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \eta_t^2 (2\phi_t + 8) \Lambda \\ &\quad + \mathbb{E}[-2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + 2\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] \\ &\quad + 2\eta_t \|\boldsymbol{\Delta}\|_2 \cdot \left(\eta_t (4 + \phi_t) k + \eta_t (4 + 3\phi_t) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) . \end{aligned}$$

Now, using the proof of Corollary 5.2-(a) which gives an upper bound on the expected value of $-\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t) + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)$, we can obtain the desired bound.

(d) This time we use the following upper bound which comes from (F.4)

$$|\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| \leq \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) + \|\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t\|_2 .$$

Plugging this into (B.1), we obtain

$$\begin{aligned} |\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})| &\leq 2\eta_t |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} \mathbf{H}'_t)| + 2\eta_t |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| + 4\eta_t^2 \|\mathbf{H}'_t\|_2 \left| \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t) \right| \\ &\quad + 12\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 . \\ &\stackrel{\textcircled{1}}{\leq} 8\eta_t \|\mathbf{H}'_t\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 4\eta_t |\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{R}'_t)| + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 \\ &\leq 8\eta_t \|\mathbf{H}'_t\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + 4\eta_t \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) \\ &\quad + 4\eta_t \|\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t\|_2 + 8\eta_t^2 \|\mathbf{R}'_t\|_2^2 \\ &\stackrel{\textcircled{2}}{\leq} 8\eta_t \|\mathbf{H}'_t\|_2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \eta_t (6\phi_t + 10) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1})^{1/2} \\ &\quad + 8\eta_t^2 \phi_t \|\mathbf{R}'_t\|_2 \end{aligned}$$

Above, ① uses the fact that $\eta_t \|\mathbf{H}'_t\|_2 \leq \eta_t \phi_t \leq 1/2$, and ② uses

$$\text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) = \|x_t^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{L}_{t-1}\|_2^2 \leq 2\|x_t^\top \mathbf{L}_{t-1}\|_2^2 + 2\|x_t^\top \mathbf{V} \mathbf{V}^\top \mathbf{L}_{t-1}\|_2^2 \leq 2(\phi_t^2 + 1) \leq 2(\phi_t + 1)^2$$

Taking square on both sides, we have

$$\begin{aligned} |\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})|_2^2 &\leq 192\eta_t^2 \|\mathbf{H}'_t\|_2^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 3\eta_t^2 (6\phi_t + 10)^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top x_t x_t^\top \mathbf{Z} \mathbf{S}_{t-1}) \\ &\quad + 192\eta_t^4 \phi_t^2 \|\mathbf{R}'_t\|_2^2 \end{aligned}$$

Finally, taking expectation and using (F.2), we have (noticing that $\eta_t \phi_t \leq 1/2$)

$$\begin{aligned} &\mathbb{E}[|\text{Tr}(\mathbf{S}_t^\top \mathbf{S}_t) - \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})|_2^2 \mid \mathcal{F}_{\leq t-1}, \mathcal{C}_{\leq t}] \\ &\leq 192\Gamma \eta_t^2 \phi_t^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 3\eta_t^2 (6\phi_t + 10)^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{Z}^\top (\boldsymbol{\Sigma} + \boldsymbol{\Delta}) \mathbf{Z} \mathbf{S}_{t-1}) \\ &\quad + 192\eta_t^4 \phi_t^2 \left(\Lambda + \lambda_{k+1} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \|\boldsymbol{\Delta}\|_2 \cdot (k + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})) \right) . \\ &\leq 192\Gamma \eta_t^2 \phi_t^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 3\eta_t^2 (6\phi_t + 10)^2 \left((\lambda_{k+1} + \|\boldsymbol{\Delta}\|_2) \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \right) \\ &\quad + 192\eta_t^4 \phi_t^2 \left(\Lambda + \lambda_{k+1} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) + \|\boldsymbol{\Delta}\|_2 \cdot (k + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})) \right) . \\ &\leq 192\Gamma \eta_t^2 \phi_t^2 \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})^2 + 4\eta_t^2 (6\phi_t + 10)^2 \lambda_{k+1} \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1}) \\ &\quad + 192\Gamma \eta_t^4 \phi_t^2 + \|\boldsymbol{\Delta}\|_2 \cdot 4\eta_t^2 (6\phi_t + 10)^2 (k + \text{Tr}(\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1})) . \end{aligned}$$

□

G Main Lemma Improvement 1: Gap-Dependent Case

In this section, we improve our main lemmas to obtain an extra $\Lambda \stackrel{\text{def}}{=} \sum_{i=1}^k \lambda_i \in (0, 1)$ factor in the gap-dependent case (i.e., when $\rho = \text{gap}$ and $\mathbf{Z} = \mathbf{W}$). We strengthen both Lemma 7.1 and Lemma 7.3.

Since the proofs of these new lemmas are analogous to the ones we had before, we spend most of this section only emphasizing the differences. At a high level, whenever we apply martingale corollaries in the old proofs (with constant κ), we now want to apply them with $\kappa \approx 1/\sqrt{\Lambda}$. This makes the notations much heavier as compared to the original proofs. We recommend readers to first take a close look at our proofs of Lemma 7.1 and Lemma 7.3 before verifying the proofs in this section.

G.1 Before Warm Start

Lemma G.1 (before warm start). *Suppose $\mathbf{W} = \mathbf{Z}$ and gap is the k -th eigengap. For every $q \in (0, \frac{1}{2}]$, $\Xi_{\mathbf{Z}} \geq 2$, $\Xi_x \geq 2$, and fixed matrix $\mathbf{Q} \in \mathbb{R}^{d \times k}$, suppose it satisfies*

- $\|\mathbf{Z}^\top \mathbf{Q}(\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \Xi_{\mathbf{Z}}$, and
- $\Pr_{x_t} \left[\forall j \in [T], \|x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma/\lambda_{k+1})^{j-1} \mathbf{Q}(\mathbf{V}^\top \mathbf{Q})^{-1}\|_2 \leq \Xi_x \right] \geq 1 - q^2/2$ for every $t \in [T]$.

Suppose also the learning rates $\{\eta_s\}_{s \in [T]}$ satisfy

$$\forall s \in [T], \frac{2q(\Xi_{\mathbf{Z}}^{3/2} + k)}{\Xi_x \cdot \Lambda} \leq \eta_s \leq O\left(\frac{\text{gap}}{\Lambda \cdot \Xi_x^2 \ln \frac{T}{q}}\right) \quad \text{and} \quad \exists T_0 \in [T] \text{ such that } \sum_{t=1}^{T_0} \eta_t \geq \Omega\left(\frac{\ln(\Xi_{\mathbf{Z}})}{\text{gap}}\right) \quad (\text{G.1})$$

Then, for every $t \in [T - 1]$, we have with probability at least $1 - 2qT$ (over the randomness of x_1, \dots, x_t):

- if $t \geq T_0$ then $\|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q}(\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq 2$.

Proof of Lemma G.1. The proof is a non-trivial adaption of the proof of Lemma 7.1.

We again consider random vectors $y_{t,s} \in \mathbb{R}^{T+2}$ defined as (we ignore coordinate $i = 1$ throughout the proof because $\mathbf{W} = \mathbf{Z}$ in this section):

$$y_{t,s}^{(2)} \stackrel{\text{def}}{=} \|\mathbf{Z}^\top \mathbf{P}_s \mathbf{Q}(\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_F^2, \\ y_{t,s}^{(3+j)} \stackrel{\text{def}}{=} \begin{cases} \|x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma/\lambda_{k+1})^j \mathbf{P}_s \mathbf{Q}(\mathbf{V}^\top \mathbf{P}_s \mathbf{Q})^{-1}\|_2^2, & \text{for } j \in \{0, 1, \dots, t-s-1\}; \\ (1 - \eta_s \lambda_k) \cdot y_{t,s-1}^{(3+j)}, & \text{for } j \in \{t-s, \dots, T-1\}. \end{cases}$$

We again consider upper bounds

$$\phi_{t,s}^{(1)} \stackrel{\text{def}}{=} 2\Xi_{\mathbf{Z}}, \quad \phi_{t,s}^{(2)} \stackrel{\text{def}}{=} \begin{cases} 2\Xi_{\mathbf{Z}} & s < T_0; \\ 2 & \text{otherwise.} \end{cases}, \quad \text{and } \phi_{t,s}^{(3+j)} \stackrel{\text{def}}{=} 2\Xi_x^2.$$

For each $t \in [T]$, define event \mathcal{C}'_t and \mathcal{C}''_t in the same way as before:

$$\mathcal{C}'_t \stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_{t-1}) \text{ satisfies } \Pr_{x_t} \left[\exists i \in \{2, 3\} : y_{t,t-1}^{(i)} > \phi_{t,t-1}^{(i)} \mid \mathcal{F}_{t-1} \right] \leq q \right\} \\ \mathcal{C}''_t \stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_t) \text{ satisfies } \forall i \in \{2, 3\} : y_{t,t-1}^{(i)} \leq \phi_{t,t-1}^{(i)} \right\}$$

and denote by $\mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$ and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$.

Verification of Assumption (A1) in Lemma D.1.

Suppose $\mathbb{E}[x_s x_s^\top \mid \mathcal{C}_{\leq s}, \mathcal{F}_{\leq s-1}] = \Sigma + \Delta$, then we have $\|\Delta\|_2 \leq \frac{q_1}{1-q_1} \leq \frac{q}{1-q}$ using the same proof as before.

This time, we use Lemma F.1 to obtain the following tighter bounds for $i = 2, \dots, T+2$:¹⁵

$$\mathbb{E}[y_{t,s+1}^{(i)} \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s}] \leq f_s^{(i)}(y_{t,s}, q) \quad \text{and} \quad \mathbb{E}[|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}|^2 \mid \mathcal{F}_t, \mathcal{F}_{\leq s}, \mathcal{C}_{\leq s}] \leq h_s^{(i)}(y_{t,s}, q)$$

where for every $j \in \{0, \dots, T-2\}$,

$$\begin{aligned} f_s^{(2)}(y, q) &\stackrel{\text{def}}{=} (1 - 2\eta_{s+1}\text{gap} + O(\Lambda\eta_{s+1}^2\Xi_x^2))y^{(2)} + O(\Lambda\eta_{s+1}^2\Xi_x + \text{Err}) , \\ h_s^{(2)}(y, q) &\stackrel{\text{def}}{=} O(\Lambda\eta_{s+1}^2\Xi_x^2(y^{(2)})^2 + \Lambda\eta_{s+1}^2\Xi_x^2y^{(2)} + \Lambda\eta_{s+1}^4\Xi_x^2 + \text{Err}) , \\ f_s^{(3+j)}(y, q) &\stackrel{\text{def}}{=} (1 - \eta_{s+1}\lambda_k + O(\Lambda\eta_{s+1}^2\Xi_x^2))y^{(3+j)} + \eta_{s+1}\lambda_{k+1}y^{(3+j+1)} + O(\Lambda\eta_{s+1}^2\Xi_x^2 + \text{Err}) , \\ h_s^{(3+j)}(y, q) &\stackrel{\text{def}}{=} O(\Lambda\eta_{s+1}^2\Xi_x^2(y^{(3+j)})^2 + \Lambda\eta_{s+1}^2\Xi_x^2y^{(3+j)} + \Lambda\eta_{s+1}^4\Xi_x^4) . \end{aligned}$$

Above, we denote by $\text{Err} \stackrel{\text{def}}{=} \eta_{s+1}\Xi_x \frac{q(\Xi_x^{3/2} + k)}{1-q}$ the error term similar to the proof of Lemma 7.1.

Obviously if $\frac{2q(\Xi_x^{3/2} + k)}{\Lambda} \leq \eta_s$ is satisfied then the Err term can be absorbed into the big- O notation.

Moreover, for every $i \in \{2, \dots, T+2\}$, consider the same g_s as defined before

$$g_s^{(i)}(y) = 20\eta_{s+1}\Xi_x \cdot y^{(i)} + 42\eta_{s+1}^2\Xi_x^2 \tag{G.2}$$

and it satisfies whenever $\mathcal{C}_{\leq s+1}$ holds then $|y_{t,s+1}^{(i)} - y_{t,s}^{(i)}| \leq g_s^{(i)}(y_{t,s})$.

Putting the above bounds together, we finish verifying assumption (A1) of Lemma D.1 with.

Verification of Assumption (A2) of Lemma D.1.

This step is exactly the same as the proof of Lemma 7.1 so ignored here.

Verification of Assumption (A3) of Lemma D.1.

For every $t \in [T]$, at a high level assumption (A3) is satisfied once we plug in the following three sets of parameter choices to Corollary 6.5 and Corollary 6.6: define $\kappa \stackrel{\text{def}}{=} 1/\sqrt{\Lambda} > 1$ and for every $s \in [T-1]$,

$$\begin{aligned} \beta_{s,2} &= 2\eta_{s+1}\text{gap}, & \delta_{s,2} &= 0, & \tau_{s,2} &= O(\eta_{s+1}\Xi_x \cdot \sqrt{\Lambda}) \\ \beta_{s,3} &= \eta_{s+1}\text{gap}, & \delta_{s,3} &= \eta_{s+1}\lambda_k & \tau_{s,3} &= O(\eta_{s+1}\Xi_x \cdot \sqrt{\Lambda}) \end{aligned}$$

More specifically, for every $t \in [T]$, let $\{z_s\}_{s=0}^{t-1}$ be the *arbitrary* random vector satisfying (D.1) of Lemma D.1. Define $q_2 = q^2/8$.

- For coordinate $i = 2$ of $\{z_s\}_{s=0}^{t-1}$,
 - if $t < T_0$, apply Corollary 6.5 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, and κ ;
 - if $t \geq T_0$, apply Corollary 6.6 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, $\gamma = 1$, and κ ;
- For coordinates $i = 3, 4, \dots, T+2$ of $\{z_s\}_{s=0}^{t-1}$,
 - apply Corollary 6.5 with $\{\beta_{s,3}, \delta_{s,3}, \tau_{s,3}\}_{s=0}^{t-2}$, $q = q_2$, $D = T$, and κ .

¹⁵In order to obtain such bounds, one needs to use the fact that when $w = x_t \mathbf{Z} \mathbf{Z}^\top$, the quantity $\frac{\eta_t}{\lambda_k} \|w^\top \Sigma \mathbf{L}_{t-1}\|_2^2$ that appeared in Lemma F.1-(a) can be upper bounded by $\frac{\eta_t \lambda_{k+1}^2}{\lambda_k} \|w^\top \mathbf{L}_{t-1}\|_2^2 \leq \eta_t \lambda_{k+1} \|w^\top \mathbf{L}_{t-1}\|_2^2$.

Note that we can apply Corollary 6.5 because our assumption $\eta_s \leq O(\frac{\text{gap}}{\Lambda \Xi_x^2 \ln \frac{T}{q}})$ implies $\beta_s \geq 12 \ln \frac{4T}{q_2} \tau_s^2$ and $\kappa \tau_s \cdot 12 \ln \frac{4T}{q_2} \leq 1$ for both $(\beta_s, \tau_s) = (\beta_{s,2}, \tau_{s,2})$ and $(\beta_{s,3}, \tau_{s,3})$. We can apply Corollary 6.6 with $\gamma = 1$ because our assumption $\eta_s \leq O(\frac{\text{gap}}{\Lambda \Xi_x^2 \ln \frac{T}{q}})$ implies $\beta_{s,2} \geq 10 \ln \frac{3T}{q_2} \cdot \tau_{s,2}^2$ for every s , and our assumption $\sum_{s=0}^{T_0-1} \beta_{s,2} \geq 1 + \ln \Xi_{\mathbf{Z}}$ implies $\sum_{s=0}^{t-1} \beta_s - 10 \ln \frac{3t}{q_2} \tau_s^2 \geq \ln \Xi_{\mathbf{Z}} + 1 - 1 = \ln \Xi_{\mathbf{Z}}$ whenever $t > T_0$.

Therefore, the conclusion of Corollary 6.5 and Corollary 6.6 imply that

$$\Pr[\exists i \in [3] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq 3q_2 < q^2/2$$

so assumption (A3) of Lemma D.1 holds.

Application of Lemma D.1. Applying Lemma D.1, we have $\Pr[\overline{\mathcal{C}_T}] \leq 2qT$ which implies our desired bounds and this finishes the proof of Lemma G.1. \square

G.2 After Warm Start

We have the following lemma and corollary

Lemma G.2 (after warm start). *In the same setting as Lemma G.1, suppose in addition there exists $\delta \leq 1/\sqrt{8}$ such that*

$$\frac{T_0}{\ln^2 T_0} \geq \frac{9 \ln(8/q^2)}{\delta^2}, \quad \forall s \in \{T_0+1, \dots, T\}: \quad 2\eta_s \text{gap} - \eta_s^2 \Xi_x^2 \geq \frac{\Omega(1)}{s-1} \quad \text{and} \quad \eta_s \leq \frac{O(1)}{\sqrt{\Lambda}(s-1)\delta \Xi_x}.$$

Then, with probability at least $1 - 2qT$ (over the randomness of x_1, \dots, x_T):

$$\bullet \quad \|\mathbf{Z}^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq \frac{5T_0/\ln^2(T_0)}{t/\ln^2 t} \text{ for every } t \in \{T_0, \dots, T\}.$$

Proof of Lemma G.2. For every $t \in [T]$ and $s \in \{0, 1, \dots, t-1\}$, consider the same random vectors $y_{t,s} \in \mathbb{R}^{T+2}$ defined in the proof of Lemma 7.1. Also, consider the same upper bounds defined in the proof of Lemma 7.3:

$$\phi_{t,s}^{(2)} \stackrel{\text{def}}{=} \begin{cases} 2\Xi_{\mathbf{Z}} & \text{if } s < T_0; \\ 2 & \text{if } s = T_0; \\ \frac{5T_0/\ln^2(T_0)}{s/\ln^2 s} & \text{if } s > T_0. \end{cases}, \quad \text{and } \phi_{t,s}^{(3+j)} \stackrel{\text{def}}{=} 2\Xi_x^2.$$

Also consider the same events $\mathcal{C}'_t, \mathcal{C}''_t, \mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$ and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$ defined as before.

We next want to apply the decoupling Lemma D.1.

Verification of Assumption (A1) in Lemma D.1.

The same functions $f_s^{(i)}, g_s^{(i)}$, and $h_s^{(i)}$ used in the proof of Lemma G.1 still apply here. We make minor changes in the spirit as the proof of Lemma 7.3: whenever $s \geq T_0$, define

$$g_s^{(2)}(y) \stackrel{\text{def}}{=} 45\eta_{s+1}\Xi_x \sqrt{y^{(2)}} + 40\eta_{s+1}^2 \Xi_x^2 \quad \text{and} \quad h_s^{(2)}(y, q) \stackrel{\text{def}}{=} O\left(\Lambda \eta_{s+1}^2 \Xi_x^2 y^{(2)} + \Lambda \eta_{s+1}^4 \Xi_x^2 + \text{Err}\right).$$

Note that we can make this change for $g_s^{(2)}$ owing to exactly the same reason as the proof of Lemma 7.3. We can do so for $h_s^{(2)}$ because whenever $\mathcal{C}_{\leq s+1}$ holds for some $s \geq T_0$ (which implies $y_{t,s}^{(2)} \leq 5$), we have $(y^{(2)})^2 = O(y^{(2)})$ so the formulation of $h_s^{(2)}$ can be simplified as above.

These choices of $f_s^{(i)}, g_s^{(i)}$, and $h_s^{(i)}$ satisfy assumption (A1) of Lemma D.1.

Verification of Assumption (A2) of Lemma D.1.

Same as before.

Verification of Assumption (A3) in Lemma D.1.

Same as the proof of Lemma 7.3, for every $t \in [T]$, let $\{z_s\}_{s=0}^{t-1}$ be the *arbitrary* random vector satisfying (D.1) of Lemma D.1. Choosing $q_2 = q^2/8$ again, the same argument before indicates that it suffices to focus on $t \geq T_0 + 2$ and prove

$$\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)} \mid z_{T_0}^{(2)} \leq 2] \leq q_2. \quad (\text{G.3})$$

We next want to apply Corollary 6.3. Recall that for every $t \in \{T_0 + 2, \dots, T\}$, the random sequence $\{z_s^{(2)}\}_{s=T_0}^{t-1}$ satisfies (D.1) with

$$\begin{aligned} f_s^{(2)}(y, q) &\stackrel{\text{def}}{=} (1 - 2\eta_{s+1}\text{gap} + O(\Lambda\eta_{s+1}^2\Xi_x^2))y^{(2)} + O(\Lambda\eta_{s+1}^2\Xi_x + \text{Err}), \\ h_s^{(2)}(y, q) &\stackrel{\text{def}}{=} O(\Lambda\eta_{s+1}^2\Xi_x^2y^{(2)} + \Lambda\eta_{s+1}^4\Xi_x^2 + \text{Err}), \\ g_s^{(2)}(y) &\stackrel{\text{def}}{=} 45\eta_{s+1}\Xi_x\sqrt{y^{(2)}} + 40\eta_{s+1}^2\Xi_x^2 \end{aligned}$$

Therefore, $\{z_s^{(2)}\}_{s=T_0}^{t-1}$ satisfies (6.1) with $\kappa = 2/\sqrt{\Lambda}$ and $\tau_s = \frac{1}{\delta s}$ because the following holds from our assumptions:

$$\begin{aligned} q\Xi_{\mathbf{Z}}^{3/2} &\leq \eta_{s+1} & \delta\tau_s &= \frac{1}{s} \leq 2\eta_{s+1}\text{gap} - \Omega(\eta_{s+1}^2\Xi_x^2) \\ \tau_s^2 &= \frac{1}{\delta^2 s^2} \geq \Omega(\Lambda\eta_{s+1}^2\Xi_x^2) & \kappa^2\tau_s^4 &= \frac{1}{\Lambda\delta^4 s^4} \geq \Omega(\Lambda\eta_{s+1}^4\Xi_x^2) & \kappa\tau_s &= \frac{2}{\sqrt{\Lambda}\delta s} \geq \Omega(\eta_{s+1}\Xi_x) \end{aligned}$$

Finally, we are ready to apply Corollary 6.3 with $q = q_2$, $t_0 = T_0$, and $\kappa = 2/\sqrt{\Lambda}$. Because $q_2 \leq e^{-2}$, $z_{T_0}^{(2)} \leq 2$, $\delta \leq 1/\sqrt{8}$ and $\frac{T_0}{\ln^2 T_0} \geq \frac{9\ln(1/q_2)}{\delta^2}$, the conclusion of Corollary 6.3 tells us $\Pr[z_{t-1}^{(2)} > \phi_{t,t-1}^{(2)} \mid z_{T_0}^{(2)} \leq 2] \leq q_2$, which is exactly (G.3) so this finishes the verification of assumption (A3).

Application of Lemma D.1. Applying Lemma D.1, we have $\Pr[\overline{\mathcal{C}}_T] \leq 2qT$ which implies our desired bounds and this finishes the proof of Lemma G.2. \square

Parameter G.3. There exist constants $C_1, C_2, C_3 > 0$ such that for every $q > 0$ that is sufficiently small (meaning $q < 1/\text{poly}(T, \Xi_{\mathbf{Z}}, \Xi_x, 1/\text{gap})$), the following parameters both satisfy Lemma G.1 and Lemma G.2:

$$\frac{T_0}{\ln^2(T_0)} = C_1 \cdot \frac{\Lambda\Xi_x^2 \ln^2 \frac{T}{q} \ln^2 \Xi_{\mathbf{Z}}}{\text{gap}^2}, \quad \eta_t = C_2 \cdot \begin{cases} \frac{\ln \Xi_{\mathbf{Z}}}{T_0 \cdot \text{gap}} & t \leq T_0; \\ \frac{1}{t \cdot \text{gap}} & t > T_0. \end{cases}, \quad \text{and} \quad \delta = C_3 \cdot \frac{\text{gap}}{\sqrt{\Lambda}\Xi_x}.$$

H Main Lemma Improvement 2: Gap-Free Case

In this section we also sketch the proof to obtain Rayleigh quotient result. We will prove the following lemma which is a strengthened version of Lemma 7.1.

Lemma H.1 (before warm start). *In the same setting as Lemma 7.1, suppose we redefine $\mathbf{W} = \mathbf{W}_\gamma$ to be the column orthonormal matrix consisting of eigenvectors of Σ with values $\leq \lambda_k - \gamma \cdot \rho$.*

Then, for every $\gamma \in [1, 1/\rho]$, with probability at least $1 - 2qT$:

$$\forall t \in \{T_0, \dots, T\}, \|\mathbf{W}_\gamma^\top \mathbf{P}_t \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_t \mathbf{Q})^{-1}\|_F^2 \leq \frac{2}{\gamma}.$$

Proof of Lemma H.1. The proof is a non-trivial adaption of the proof of Lemma 7.1.

We redefine $\mathbf{W} = \mathbf{W}_\gamma$ and consider random vectors $y_{t,s} \in \mathbb{R}^{T+2}$ defined in the same way as the proof of Lemma 7.1. This time, we consider upper bounds

$$\phi_{t,s}^{(1)} \stackrel{\text{def}}{=} 2\Xi_{\mathbf{Z}}, \quad \phi_{t,s}^{(2)} \stackrel{\text{def}}{=} \begin{cases} 2\Xi_{\mathbf{Z}} & s < T_0; \\ 2/\gamma & s \geq T_0. \end{cases}, \quad \text{and } \phi_{t,s}^{(3+j)} \stackrel{\text{def}}{=} 2\Xi_x^2$$

so the only difference we make here is on coordinate $i = 2$ for $s \geq T_0$. For each $t \in [T]$, we also consider events $\mathcal{C}'_t, \mathcal{C}''_t, \mathcal{C}_t \stackrel{\text{def}}{=} \mathcal{C}'_t \wedge \mathcal{C}''_t$, and $\mathcal{C}_{\leq t} \stackrel{\text{def}}{=} \bigwedge_{s=1}^t \mathcal{C}_s$ defined in the same way as before.

Verification of Assumption (A1) in Lemma D.1.

We consider the same functions f_s, g_s, h_s as defined in the proof of Lemma 7.1, except that we replace ρ with $\gamma \cdot \rho$ because this time we have redefined $\mathbf{W} = \mathbf{W}_\gamma$ so that it consists of eigenvectors with values $\leq \lambda_k - \gamma \cdot \rho$. In other words, we redefine

$$f_s^{(2)}(y, q) = (1 - 2\eta_{s+1}\gamma\rho + 56\eta_{s+1}^2\Xi_x^2)y^{(2)} + 40\eta_{s+1}^2\Xi_x^2 + 20\eta_{s+1}\frac{q\Xi_{\mathbf{Z}}^{3/2}}{1-q}.$$

In the same way we can verify that these functions satisfy assumption (A1) of Lemma D.1.

Verification of Assumption (A2) of Lemma D.1.

This step is exactly the same as the proof of Lemma 7.1 so ignored here.

Verification of Assumption (A3) of Lemma D.1.

We consider the same parameters $\{\beta_s, \delta_s, \tau_s\}_s$ as Lemma 7.1 except that at coordinate $i = 2$ we replace ρ with $\gamma \cdot \rho$:

$$\beta_{s,2} = 2\eta_{s+1}\gamma\rho, \quad \delta_{s,2} = 0, \quad \tau_{s,2} = 20\eta_{s+1}\Xi_x.$$

Now, for every $t \in [T]$, let $\{z_s\}_{s=0}^{t-1}$ be the *arbitrary* random vector satisfying (D.1) of Lemma D.1. Letting $q_2 = q^2/8$, we can handle coordinates $i = 1$ and $i \geq 3$ in the same way as before. As for coordinate $i = 2$ of $\{z_s\}_{s=0}^{t-1}$,

- if $t < T_0$, apply Corollary 6.4 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, and $\kappa = 1$;
- if $t \geq T_0$, apply Corollary 6.6 with $\{\beta_{s,2}, \delta_{s,2}, \tau_{s,2}\}_{s=0}^{t-2}$, $q = q_2$, $D = 1$, $\gamma = \gamma$, and $\kappa = 1$;

Note that the $t < T_0$ case is exactly the same as before. When $t \geq T_0$, we again apply Corollary 6.6 but this time with value $\gamma \geq 1$ rather than $\gamma = 1$. Since this is the only difference here, we only need to verify the the presumptions of Corollary 6.6:

- our assumption $\eta_s \leq \frac{\rho}{4000 \cdot \Xi_x^2 \ln \frac{3T}{q_2}}$ implies $\beta_{s,2} \geq 20\gamma \ln \frac{3T}{q_2} \cdot \tau_{s,2}^2$ and $\kappa\tau_s \leq \frac{1}{12 \ln \frac{3T}{q_2}}$ for every s ,
- our assumption $\sum_{s=0}^{T_0-1} \beta_{s,2} \geq 1 + \ln \Xi_{\mathbf{Z}}$ implies $\sum_{s=0}^{t-1} \beta_{s,2} - 10\gamma \ln \frac{3t}{q_2} \tau_{s,2}^2 \geq \frac{1}{2} \sum_{s=0}^{t-1} \beta_{s,2} \geq \ln \Xi_{\mathbf{Z}}$ whenever $t > T_0$.

Therefore, in the same way as the old proof in Lemma 7.1, we can conclude using Corollary 6.4 and Corollary 6.6 that

$$\Pr[\exists i \in [3] : z_{t-1}^{(i)} > \phi_{t,t-1}^{(i)}] \leq 3q_2 < q^2/2.$$

This verifies assumption (A3) of Lemma D.1.

Application of Lemma D.1. Applying Lemma D.1, we have $\Pr[\overline{\mathcal{C}_T}] \leq 2qT$ which implies our desired bounds and this finishes the proof of Lemma H.1. \square

I Missing Proofs for Final Theorems

We prove Theorem 2 first, and then Theorem 1 and Theorem 3.

I.1 Proof of Theorem 2

Proof of Theorem 2. First for a sufficiently large constant C , we can apply Lemma 4.3 with $p' = \frac{p}{6}$ and $q = \min \left\{ \frac{1}{CT^2d^2}, \frac{p}{4T} \right\}$ and obtain: with probability at least $1 - p' - q^2 \geq 1 - p/2$ over the random choice of \mathbf{Q} , the following holds:

$$\begin{cases} \|(\mathbf{Z}^\top \mathbf{Q})(\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \frac{20736dk}{p^2} \ln \frac{6d}{p}, \text{ and} \\ \Pr_{x_1, \dots, x_T} \left[\exists i \in [T], \exists t \in [T], \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \frac{216 \sqrt{k \ln \frac{2T}{p}}}{p} \right] \leq \frac{q^2}{2}. \end{cases}$$

Denote by \mathcal{C}_1 the union of the above two events, and we have $\Pr_{\mathbf{Q}}[\mathcal{C}_1] \geq 1 - p/2$.

Now, for every fixed \mathbf{Q} , whenever \mathcal{C}_1 holds, we can let

$$\Xi_{\mathbf{Z}} = \frac{20736dk}{p^2} \ln \frac{6d}{p}, \quad \Xi_x = \frac{216 \sqrt{2k \ln \frac{2T}{p}}}{p},$$

so the initial conditions in Lemma 7.1 (and thus Lemma 7.3) is satisfied. Also, according to Parameter 7.4, our parameter choices satisfy the assumptions in Lemma 7.3. Finally, the conclusion of Lemma 7.3 immediately implies for every $T \geq T_0$

$$\Pr_{x_1, \dots, x_T} \left[\|\mathbf{W}^\top \mathbf{P}_T \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_T \mathbf{Q})^{-1}\|_F^2 = \tilde{O} \left(\frac{T_0}{T} \right) \mid \mathcal{C}_1 \right] \geq 1 - 2qT \geq 1 - \frac{p}{2}.$$

Union bounding this with event $\overline{\mathcal{C}_1}$, we have

$$\Pr_{\mathbf{Q}, x_1, \dots, x_T} \left[\|\mathbf{W}^\top \mathbf{P}_T \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_T \mathbf{Q})^{-1}\|_F^2 = \tilde{O} \left(\frac{T_0}{T} \right) \right] \geq 1 - p.$$

Combining this with Lemma 2.2 completes the proof. \square

I.2 Proof of Theorem 1

Proof of Theorem 1. First for a sufficiently large constant C , we can apply Lemma 4.3 on $p' = \frac{p}{6}$ and $q = \min \left\{ \frac{1}{CT^2d^2}, \frac{p}{8T} \right\}$ and obtain: with probability at least $1 - p' - q^2 \geq 1 - p/2$ over the random choice of \mathbf{Q} , the following holds:

$$\begin{cases} \|(\mathbf{Z}^\top \mathbf{Q})(\mathbf{V}^\top \mathbf{Q})^{-1}\|_F^2 \leq \frac{20736dk}{p^2} \ln \frac{6d}{p}, \text{ and} \\ \Pr_{x_1, \dots, x_T} \left[\exists i \in [T], \exists t \in [T], \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\Sigma / \lambda_{k+1})^{i-1} \mathbf{Q} (\mathbf{V}^\top \mathbf{Q})^{-1} \right\|_2 \geq \frac{216 \sqrt{k \ln \frac{2T}{p}}}{p} \right] \leq \frac{q^2}{2}. \end{cases}$$

Denote by \mathcal{C}_1 the union of the above two events, and we have $\Pr_{\mathbf{Q}}[\mathcal{C}_1] \geq 1 - p/2$.

Now, whenever \mathcal{C}_1 holds, we can set

$$\Xi_{\mathbf{Z}} = \frac{20736dk}{p^2} \ln \frac{6d}{p}, \quad \Xi_x = \frac{216 \sqrt{2k \ln \frac{2T}{p}}}{p}.$$

so the initial conditions in Lemma G.1 are satisfied. Also, according to Parameter G.3, our parameter choices satisfy the assumptions in Lemma G.1. Therefore, the conclusion of Lemma G.1 implies

$$\Pr_{x_1, \dots, x_{T_0}} \left[\|\mathbf{Z}^\top \mathbf{P}_{T_0} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{T_0} \mathbf{Q})^{-1}\|_F^2 \geq 2 \mid \mathcal{C}_1 \right] \leq 2qT.$$

We denote by \mathcal{C}_2 the event that $\|\mathbf{Z}^\top \mathbf{P}_{T_0} \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_{T_0} \mathbf{Q})^{-1}\|_F^2 \geq 2$. Note that \mathcal{C}_2 only depends on the randomness of \mathbf{Q} and x_1, \dots, x_{T_0} .

Whenever \mathcal{C}_2 holds, denoting by $\mathbf{Q}' = \mathbf{P}_{T_0} \mathbf{Q}$, we have:¹⁶

$$\begin{cases} \|\mathbf{Z}^\top \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1}\|_F^2 \leq 2, \text{ and} \\ \forall i \in [T], t \in \{T_0 + 1, \dots, T\} : \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^{i-1} \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1} \right\|_2 \leq 3 \end{cases}$$

We next want to apply Lemma 7.3 again but on x_{T_0}, \dots, x_T : we shift all the indices by $-T_0$, meaning that x_t now becomes x_{t-T_0} . This time we apply Lemma 7.3 with $\mathbf{Q} = \mathbf{Q}'$, $\Xi_{\mathbf{Z}} = 2$, and $\Xi_x = 3$. We use again the parameter choices of Parameter G.3 but this time we denote by T_1 this new T_0 and it satisfies:

$$\frac{T_1}{\ln^2(T_1)} = \Theta\left(\frac{\Lambda \Xi_x^2 \ln^2 \frac{T}{q} \ln^2 \Xi_{\mathbf{Z}}}{\text{gap}^2}\right) = \Theta\left(\frac{\Lambda \ln^2 \frac{T}{q}}{\text{gap}^2}\right).$$

The conclusion of Lemma 7.3 tells us that, denoting by $\mathbf{P}_{T_0:t} = \prod_{s=T_0+1}^t (\mathbf{I} + \eta_s x_s x_s^\top)$, we have for every $t \geq T_0 + T_1$,

$$\Pr_{x_{T_0+1}, \dots, x_t} \left[\|\mathbf{Z}^\top \mathbf{P}_{T_0:t} \mathbf{Q}' (\mathbf{V}^\top \mathbf{P}_{T_0:t} \mathbf{Q}')^{-1}\|_F^2 \geq \frac{5T_1 \ln T_1}{(t - T_0) \ln(t - T_0)} \mid \mathcal{C}_2 \right] \leq 2qT.$$

In other words, if $T \geq T_0 + T_1$, then

$$\begin{aligned} & \Pr_{\mathbf{Q}, x_1, \dots, x_T} \left[\|\mathbf{Z}^\top \mathbf{P}_T \mathbf{Q} (\mathbf{V}^\top \mathbf{P}_T \mathbf{Q})^{-1}\|_F^2 = \tilde{\Omega}\left(\frac{T_1}{T - T_0}\right) \right] \\ & \leq \Pr_{x_{T_0+1}, \dots, x_T} \left[\|\mathbf{W}^\top \mathbf{P}_{T_0:T} \mathbf{Q}' (\mathbf{V}^\top \mathbf{P}_T \mathbf{Q}')^{-1}\|_F^2 = \tilde{\Omega}\left(\frac{T_1}{T - T_0}\right) \mid \mathcal{C}_2 \right] + \Pr_{x_1, \dots, x_{T_0}} [\overline{\mathcal{C}_2} \mid \mathcal{C}_1] + \Pr_{\mathbf{Q}} [\overline{\mathcal{C}_1}] \\ & \leq 2qT + 2qT + p/2 \leq p. \end{aligned}$$

Combining this with Lemma 2.2 completes the proof. \square

I.3 Proof of Theorem 3

Proof of Theorem 3. Recall that we are using the same learning rates Parameter 7.4 as in Theorem 2. Therefore, the same proof of Theorem 2 ensures that the initialization assumptions in Lemma H.1 are satisfied so we can apply Lemma H.1.

We want to prove next the output matrix $\mathbf{Q}_T = [q_1, \dots, q_k] \in \mathbb{R}^{d \times k}$ satisfies

$$\text{with probability at least } 1 - (2kdT)q, \quad \forall i \in [k]: \quad q_i^\top \boldsymbol{\Sigma} q_i \geq \lambda_i - 3\rho \ln \frac{1}{\rho}.$$

For every $i \in [k]$, let $\mathbf{Q}_T^i \in \mathbb{R}^{d \times i}$ denote the first i -columns of \mathbf{Q}_T . By the property of Oja's algorithm, the same \mathbf{Q}_T^i would have been the output if we started from an $\mathbb{R}^{d \times i}$ random matrix \mathbf{Q}_0 for online i -PCA. In other words, we can write $\mathbf{Q}_T^i = [q_1, \dots, q_i]$.

Letting \mathbf{W}_γ^i be the column orthonormal matrix consisting of all eigenvectors of $\boldsymbol{\Sigma}$ with eigenvalue $\leq \lambda_i - \gamma \cdot \rho$, we applying Lemma H.1 (with $k = i$) and obtain:

$$\text{w.p. at least } 1 - 2qT: \quad \|(\mathbf{W}_\gamma^i)^\top \mathbf{Q}_T^i\|_F^2 \leq \|(\mathbf{W}_\gamma^i)^\top \mathbf{P}_T \mathbf{Q}^i (\mathbf{V}^\top \mathbf{P}_T \mathbf{Q}^i)^{-1}\|_F^2 \leq 2/\gamma.$$

¹⁶Note that the second line implies by the first line:

$$\begin{aligned} & \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^{i-1} \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1} \right\|_2 \\ & \leq \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^{i-1} \mathbf{Z} \mathbf{Z}^\top \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1} \right\|_2 + \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^{i-1} \mathbf{V} \mathbf{V}^\top \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1} \right\|_2 \\ & \leq \left\| x_t^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\Sigma} / \lambda_{k+1})^{i-1} \mathbf{Z} \right\|_2 \cdot \left\| \mathbf{Z}^\top \mathbf{Q}' (\mathbf{V}^\top \mathbf{Q}')^{-1} \right\|_2 + 1 \leq \sqrt{2} + 1 < 3. \end{aligned}$$

(Above, the first inequality uses Lemma 2.2.) This in particular implies $\|(\mathbf{W}_\gamma^i)^\top q_i\|_2^2 \leq \frac{2}{\gamma}$.

Let us define for each $i \in [k]$,

$$\Gamma_i \stackrel{\text{def}}{=} \left\{ \frac{\lambda_i - \lambda_j}{\rho} \mid \lambda_i - \lambda_j \geq \rho \right\} \subseteq \mathbb{R}_{\geq 1} \quad \text{and} \quad \gamma_{i,j} \stackrel{\text{def}}{=} \frac{\lambda_i - \lambda_j}{\rho} \in [1, \frac{1}{\rho}] .$$

By union bound,

$$\text{w.p. at least } 1 - 2qkdT, \quad \forall i \in [k], \forall \gamma \in \Gamma_i: \quad \|(\mathbf{W}_\gamma^i)^\top q_i\|_2^2 \leq 2/\gamma . \quad (\text{I.1})$$

We are now ready to bound Rayleigh quotient. For each $i \in [k]$, let i_0 be the index of the first (i.e., the largest) eigenvector with eigenvalue $\leq \lambda_i - \rho$ and define $b_{i,j} \stackrel{\text{def}}{=} \sum_{s=j}^d \langle q_i, v_s \rangle^2$ where v_j is the j -th largest eigenvector of Σ . It satisfies $b_{i,1} = 1$. By Abel's formula,

$$q_i^\top \Sigma q_i = \sum_{j=1}^d \lambda_j \langle q_i, v_j \rangle^2 \geq (\lambda_i - \rho) - \sum_{j=i_0+1}^d b_{i,j} (\lambda_{j-1} - \lambda_j) .$$

Note that for every $j \geq i_0 + 1$, we have $b_{i,j} \leq \|\mathbf{W}_{\gamma_{i,j}}^i q_i\|_2^2 \leq \frac{2}{\gamma_{i,j}}$ according to (I.1). Therefore,

$$\sum_{j=i_0+1}^d b_{i,j} (\lambda_{j-1} - \lambda_j) \leq \sum_{j=i_0+1}^d \frac{2}{\gamma_{i,j}} \rho (\gamma_{i,j} - \gamma_{i,j-1}) \leq 2\rho \int_1^{\frac{1}{\rho}} \frac{1}{\gamma} dz \leq 2\rho \ln \frac{1}{\rho} ,$$

which implies $q_i^\top \Sigma q_i \geq \lambda_i - 3\rho \ln \frac{1}{\rho}$. □

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Doubly Accelerated Methods for Faster CCA and Generalized Eigendecomposition. *ArXiv e-prints*, abs/1607.06017, July 2016.
- [2] Zeyuan Allen-Zhu and Yuanzhi Li. Even Faster SVD Decomposition Yet Without Agonizing Pain. *ArXiv e-prints*, abs/1607.03463, May 2016.
- [3] Akshay Balsubramani, Sanjoy Dasgupta, and Yoav Freund. The fast convergence of incremental pca. In *NIPS*, pages 3174–3182, 2013.
- [4] Fan Chung and Linyuan Lu. Concentration inequalities and martingale inequalities: a survey. *Internet Mathematics*, 3(1):79–127, 2006.
- [5] Dan Garber and Elad Hazan. Fast and simple PCA via convex optimization. *ArXiv e-prints*, September 2015.
- [6] Daniel Garber, Elad Hazan, Chi Jin, Sham M. Kakade, Cameron Musco, Praneeth Netrapalli, and Aaron Sidford. Robust shift-and-invert preconditioning: Faster and more sample efficient algorithms for eigenvector computation. In *ICML*, 2016.
- [7] Rong Ge, Chi Jin, Sham M. Kakade, Praneeth Netrapalli, and Aaron Sidford. Efficient Algorithms for Large-scale Generalized Eigenvector Computation and Canonical Correlation Analysis. *ArXiv e-prints*, abs/1604.03930, April 2016.
- [8] Moritz Hardt and Eric Price. The noisy power method: A meta algorithm with applications. In *NIPS*, pages 2861–2869, 2014.

- [9] Prateek Jain, Chi Jin, Sham M. Kakade, Praneeth Netrapalli, and Aaron Sidford. Streaming PCA: Matching Matrix Bernstein and Near-Optimal Finite Sample Guarantees for Oja’s Algorithm. In *COLT*, 2016.
- [10] Chris J. Li, Mengdi Wang, Han Liu, and Tong Zhang. Near-Optimal Stochastic Approximation for Online Principal Component Estimation. *ArXiv e-prints*, abs/1603.05305, March 2016.
- [11] Jieming Mao. private communication, 2016.
- [12] Ioannis Mitliagkas, Constantine Caramanis, and Prateek Jain. Memory limited, streaming pca. In *NIPS*, pages 2886–2894, 2013.
- [13] Cameron Musco and Christopher Musco. Randomized block krylov methods for stronger and faster approximate singular value decomposition. In *NIPS*, pages 1396–1404, 2015.
- [14] Christopher De Sa, Christopher Re, and Kunle Olukotun. Global convergence of stochastic gradient descent for some non-convex matrix problems. In *ICML*, pages 2332–2341, 2015.
- [15] Ohad Shamir. Convergence of stochastic gradient descent for pca. In *ICML*, 2016.
- [16] Ohad Shamir. Fast stochastic algorithms for svd and pca: Convergence properties and convexity. In *ICML*, 2016.
- [17] Stanislaw J Szarek. Condition numbers of random matrices. *Journal of Complexity*, 7(2):131–149, 1991.
- [18] Weiran Wang, Jiale Wang, Dan Garber, and Nathan Srebro. Efficient Globally Convergent Stochastic Optimization for Canonical Correlation Analysis. *ArXiv e-prints*, abs/1604.01870, April 2016.
- [19] Bo Xie, Yingyu Liang, and Le Song. Scale up nonlinear component analysis with doubly stochastic gradients. In *NIPS*, pages 2341–2349, 2015.