

LONG-CONTEXT LINEAR SYSTEM IDENTIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper addresses the problem of long-context linear system identification, where the state x_t of a dynamical system at time t depends linearly on previous states x_s over a fixed context window of length p . We establish a sample complexity bound that matches the *i.i.d.* parametric rate up to logarithmic factors for a broad class of systems, extending previous works that considered only first-order dependencies. Our findings reveal a “learning-without-mixing” phenomenon, indicating that learning long-context linear autoregressive models is not hindered by slow mixing properties potentially associated with extended context windows. Additionally, we extend these results to (i) shared low-rank representations, where rank-regularized estimators improve rates with respect to dimensionality, and (ii) misspecified context lengths in strictly stable systems, where shorter contexts offer statistical advantages.

1 INTRODUCTION

System identification, which consists of estimating the parameters of a dynamical system from observations of its trajectories, is a fundamental problem in many fields such as econometrics, robotics, aeronautics, mechanical engineering, or reinforcement learning (Ljung, 1998; Gupta et al., 1976; Moerland et al., 2022). Recent theoretical advances focused on linear system identification, where observations are of the form:

$$x_t = A^* x_{t-1} + \xi_t, \quad (1)$$

for $t \geq 1$, with initialization $x_0 \in \mathbb{R}^d$, noise $\xi_t \in \mathbb{R}^d$ and design matrix $A^* \in \mathbb{R}^{d \times d}$. Linear system identification (Simpkins, 1999) has been thoroughly studied, with recent interest in sharp non-asymptotic rates (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019; Faradonbeh et al., 2018; Jedra & Proutiere, 2019). The existing analyses, however, focus solely on order-1 time dependency, in which the law of x_t only depends on the previous state x_{t-1} . For order- p time dependencies, the literature on non-asymptotic rates becomes surprisingly scarce, as existing techniques do not extend to $p > 1$.

We study this more general setting, where the state x_t depends on previous states x_s for s in a context window of length $p \in \mathbb{N}^*$, i.e.,

$$x_t = \sum_{k=1}^p A_k^* x_{t-k} + \xi_t, \quad (2)$$

for $t \geq p$, the initialization $x_0, \dots, x_{p-1} \in \mathbb{R}^d$, noise $\xi_t \in \mathbb{R}^d$ and design matrices $A_1^*, \dots, A_p^* \in \mathbb{R}^{d \times d}$. This classical p^{th} -order vector autoregression model (Box et al., 2015; Brockwell & Davis, 1991; Hamilton, 2020) is termed *long-context linear autoregressive model*. The term *linear* refers to the (noisy) linear relationship between iterates and *long-context* refers to the context length p . Recent advances in autoregressive models and architectures such as transformers (Vaswani et al., 2017; Dosovitskiy et al., 2020; El-Nouby et al., 2024) highlight the importance of long-context and its impact on learning. Developing a theoretical understanding of long-context linear autoregressive models is a necessary first step toward tackling these more complex architectures.

Motivated by empirical evidence that high-dimensional data may share some lower-dimensional representation (Bengio et al., 2013; Hospedales et al., 2022), several works additionally studied the problem of learning matrices A_k^* under the assumption that they are of low-rank (Alquier et al.,

2020; Basu et al., 2019), for order-1 autoregressive models. In the long-context setting, this problem is further motivated by the fact that if there exists a lower-dimensional representation of the autoregressive process, this translates into shared kernels for the matrices A_k^* .

Finally, a key challenge in long-context autoregressive models is *misspecification*: the system might have an *unknown* context window p as in Equation (2). p may be arbitrarily large and unknown by the statistician. She may then specify a context length p' that can be much smaller, thus yielding the following two fundamental questions: can useful structure still be learned under misspecified context lengths? And what advantages, if any, arise from model misspecification?

Our contributions in *long-context linear systems identification* are then threefold.

(i) We derive statistical rates on the recovery of matrices A_k^* in terms of Frobenius norm, which depends on the number of trajectories N and their length T , on the dimension d and the context length p . These rates reveal a “learning-without-mixing” phenomenon as they do not have a deflation in effective sample size due to the mixing time of the autoregressive process. This first contribution is an attempt to fill the gap in linear system identification for long context lengths.

(ii) We study statistical guarantees for learning the matrices A_k^* assuming that they are all of rank at most $r \ll d$. We prove that the statistical rate reduces, and that rank-regularized estimators adapt to the low-rank structure.

(iii) We study a scenario under which the model is *misspecified*. Fitting a linear model with context length $p' < p$ instead of p , we show that the first p' matrices are still learned. More importantly, the sample complexity of learning these matrices depends only on the misspecified context length, indicating that misspecification may benefit the model statistically, not just computationally.

Finally, we confirm these statistical rates through experiments that verify the scaling laws predicted by problem parameters. Due to space constraints, these experiments are provided in Appendix C.

2 RELATED WORKS

In multivariate linear regression, one observes $\{(x_i, y_i)\}_{i=1}^N$ from the model $y_i = A^* x_i + \xi_i$, where matrix $A^* \in \mathbb{R}^{d \times d}$ and the sequences of noise ξ_i and inputs x_i are *i.i.d.*. The number of samples N needs to scale at least as d^2 for a good estimation of A^* with ordinary least squares estimator (Hsu et al., 2012; Wainwright, 2019) in Frobenius norm— $\|A^* - \hat{A}\|_F^2 \ll 1$. However, in many domains, data is sequential, violating the *i.i.d.* assumption. In such domains, classical non-*i.i.d.* formulations, such as vector autoregressive models or discrete-time linear dynamical systems (LDS), as seen in Equation (1), are often employed. Most works used to deal with the non-*i.i.d.*-ness of the data through mixing time arguments that fall short when the spectral radius of A^* reaches 1, leading to rates of the form $\|\hat{A} - A^*\|_F^2 = \mathcal{O}(d^2/(n(1 - \rho)))$ or $\|\hat{A} - A^*\|_{\text{op}}^2 = \mathcal{O}(d/(n(1 - \rho)))$ for some spectral quantity $1 - \rho$ related to the mixing time of the process. These rates apply to the OLS estimator (Faradonbeh et al., 2018) and online settings (Hardt et al., 2018; Even, 2023) alike.

Simchowitz et al. (2018); Sarkar & Rakhlin (2019) have developed excitation-based arguments to leverage mixing-time independent statistical bounds for the OLS estimator, while Hazan et al. (2017); Jain et al. (2021) respectively used spectral filtering and reverse experience replay in the online setting to obtain such bounds. The estimation of low-rank features has been studied by Basu et al. (2019); Alquier et al. (2020) via nuclear norm regularization. Finally, learning parameters of dynamical systems from N trajectories of length T has previously been considered by Tu et al. (2023) in a more general framework than Equation (1).

Layers of complexity can be added to the LDS described in Equation (1). Mania et al. (2022); Foster et al. (2020) considered non-linear dynamics, that write respectively as $x_{t+1} = A^* \phi(x_t, u_t) + \xi_t$ and $x_{t+1} = f^*(x_t) + \xi_t$, where in the former A^* is to be estimated and ϕ is a known non-linearity, while in the latter f^* is to be estimated. Kostic et al. (2022) recently provided a general framework using Koopman operators, to estimate the parameters of some general Markov chain. Giraud et al. (2015) considered time-varying systems, with arbitrary context lengths, while Bacchiocchi et al. (2024) studied autoregressive bandits. Ziemann & Tu (2022) provided a framework for learning non-parametric dynamical systems with “little mixing”: as their rates are not hindered by slow mixing after a burn-in time (that may itself depend on mixing properties). We refer the reader to

Tsiamis et al. (2023) for a survey on recent advances on non-asymptotic system identification of LDS such as in Equation (1). Surprisingly, there does not seem to be much known about *long-context LDS* in Equation (2), the counterparts of LDS in Equation (1) with a context window $p > 1$.

3 PROBLEM SETTING

For a matrix $M \in \mathbb{R}^{d_1 \times d_2}$ with singular values $\sigma_1, \dots, \sigma_{\min\{d_1, d_2\}}$, we denote its squared Frobenius norm as $\|M\|_F^2 = \sum_{(i,j)} M_{ij}^2 = \sum_{\ell} \sigma_{\ell}^2$, operator norm as $\|M\|_{\text{op}} = \max_{\ell} |\sigma_{\ell}|$, and nuclear norm as $\|M\|_* = \sum_{\ell} |\sigma_{\ell}|$. I_d and 0_d denotes the identity and the null $d \times d$ matrices, respectively. $\mathbf{A} = (A_1, \dots, A_p)$ denotes a rectangular matrix of size $d \times pd$ where each A_i is $d \times d$ block.

3.1 DATA GENERATION PROCESS

Let $d, p \in \mathbb{N}^*$ be the dimension of the state space and the context length, respectively. Consider the following linear autoregressive process:

$$\forall t > 0 : \quad x_t = \sum_{k=1}^p A_k^* x_{t-k} + \xi_t, \quad (3)$$

where $x_s = 0$ for any $s \leq 0$ and the noise ξ_t is independent of the x_s, ξ_s for $s < t$. This is a particular instance of the general linear autoregressive model in Equation (2) with initial conditions x_0, \dots, x_{p-1} set to 0 and the independent noise structure. We assume sub-Gaussian noise:

Assumption 3.1. For all t , the noise ξ_t is centered and isotropic:

$$\mathbb{E}[\xi_t] = 0, \quad \mathbb{E}[\xi_t \xi_t^\top] = \sigma^2 I_d,$$

and each coordinate of ξ_t is independent and $c^2 \sigma^2$ -sub-Gaussian (Wainwright, 2019, Chapter 2) for some $c \geq 1$:

$$\forall i \in [d] \quad \|(\xi_t)_i\|_{\psi_2} \leq c^2 \sigma^2, \quad \text{where} \quad \|x\|_{\psi_2} = \sup_{k \geq 1} k^{-1/2} \mathbb{E}[|x|^k]^{1/k}.$$

Let $\text{AR}(\mathbf{A}^*, \sigma^2)$ denote the law of the sequence defined in Equation (3) where \mathbf{A}^* denotes (A_1^*, \dots, A_p^*) for brevity. Given N independent sequences of length $T > p$:

$$\left\{ x_t^{(n)}, n \in [N], t \in [T] \right\}, \quad \text{where} \quad (x_t^{(n)})_{t \in [T]} \stackrel{i.i.d.}{\sim} \text{AR}(\mathbf{A}^*, \sigma^2),$$

the goal of long-context linear system identification is to estimate the matrices $A_k^*, k \in [p]$.

Lastly, we assume a condition on the design matrices $A_k^*, k \in [p]$ that amounts to an operator norm bound. First, we define the following linear operators for any matrix $\mathbf{A} \in \mathbb{R}^{d \times pd}$:

Definition 3.2. Let $M_{\mathbf{A}} \in \mathbb{R}^{Td \times Td}$ be the block-matrix with block entries of size $d \times d$:

$$M_{\mathbf{A}}^{(i,j)} = A_{i-j}, \quad \text{for all } 1 \leq j < i \leq j+p \leq T, \quad \text{and} \quad M_{\mathbf{A}}^{(i,j)} = 0_d, \quad \text{otherwise}.$$

Definition 3.3. Let $L_{\star} \in \mathbb{R}^{Td \times Td}$ be the block-matrix with block entries of size $d \times d$:

$$\begin{aligned} L_{\star}^{(1,1)} &= I_d \quad \text{and} \quad L_{\star}^{(i,1)} = \sum_{k=1}^{\max\{i-1, p\}} A_k^* L_{\star}^{(i-k,1)} \quad 1 < i \leq T, \\ L_{\star}^{(i,j)} &= L_{\star}^{(i-j+1,1)} \quad \text{for all } 1 \leq i \leq j \leq p, \quad \text{and} \quad L_{\star}^{(i,j)} = 0_d \quad \text{otherwise}. \end{aligned}$$

$M_{\mathbf{A}}$ executes predictions from the given data with \mathbf{A} and L_{\star} generates the data from the noise. That is, letting $(M_{\mathbf{A}})_t, (L_{\star})_t : \mathbb{R}^d \times \mathbb{R}^{Td}$ be the t^{th} block-row of $M_{\mathbf{A}}$ and L_{\star} , respectively, we have:

$$(M_{\mathbf{A}})_t \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_T^{(n)} \end{pmatrix} = \sum_{k=1}^p A_k x_{t-k}^{(n)}, \quad (L_{\star})_t \begin{pmatrix} \xi_1^{(n)} \\ \vdots \\ \xi_T^{(n)} \end{pmatrix} = x_t^{(n)}, \quad \text{with } x_s = 0 \text{ for } s \leq 0.$$

Therefore, the operator norm of $M_{\mathbf{A}}$ is a measure of the worst-case growth of the predictions. Moreover, $M_{\mathbf{A}^*}$ is linked to the data-generating operator L_* :

$$L_* = I_{Td} + M_{\mathbf{A}^*} L_* \implies L_* = (I_{Td} - M_{\mathbf{A}^*})^{-1} = I_{Td} + \sum_{i=1}^{T-1} (M_{\mathbf{A}^*})^i.$$

We assume the following conditions on the design matrices:

Assumption 3.4. *There exists a known constant $D > 0$ such that $\|M_{\mathbf{A}^*}\|_{\text{op}} \leq D$.*

Assumption 3.4 is not restrictive as D is arbitrary and only needs to be an upper bound on $\|M_{\mathbf{A}^*}\|_{\text{op}}$. However, the *knowledge* of D is necessary, as it is used to confine the estimator in Section 3.2.

As the operator $M_{\mathbf{A}^*}$ is a derived object over the full trajectory, it is important to relate Assumption 3.4 to conditions on the design matrices A_k^* . In Proposition 3.5 below, we provide two different assumptions on the design matrices that ensure the boundedness of the operator norm of $M_{\mathbf{A}^*}$ with the same constant. Both conditions *imply* Assumption 3.4.

Proposition 3.5. *Assumption 3.4 holds if one of the following holds:*

$$(i) \quad \sum_{i=1}^p \|A_i^*\|_{\text{op}} \leq D, \quad (ii) \quad \|\mathbf{A}^*\|_{\text{op}} \leq \frac{D}{\sqrt{p}}.$$

There is no direct assumption on L_* ; yet, our results depend on well-behavedness of κ , the logarithm of the condition number of L_* , which is related to Γ_t that appears in Simchowitz et al. (2018); Sarkar & Rakhlin (2019). κ is related to the system stability, as explained in Section 5.

Definition 3.6. *Let κ be the logarithm of the condition number of L_* , i.e., $\kappa := \ln \left(\frac{\|L_*\|_{\text{op}}}{\sigma_{\min}(L_*)} \right)$.*

3.2 CONSTRAINED LEAST SQUARES

A natural estimator is the *Ordinary Least Square* (OLS), defined as any minimizer of the square loss:

$$\hat{\mathbf{A}}_{\text{OLS}} \in \operatorname{argmin}_{\mathbf{A}} \mathcal{L}(\mathbf{A}), \quad \text{where} \quad \mathcal{L}(\mathbf{A}) := \frac{1}{NT} \sum_{n=1}^N \sum_{t=p}^T \left\| x_t^{(n)} - \sum_{k=1}^p A_k x_{t-k}^{(n)} \right\|^2. \quad (4)$$

The OLS estimator has been considered in previous works (Simchowitz et al., 2018; Alquier et al., 2020; Faradonbeh et al., 2018; Sarkar & Rakhlin, 2019), albeit in the $p = 1$ case. Most of these works provide estimation rates on $\|\hat{\mathbf{A}} - \mathbf{A}^*\|_{\text{op}}$ or $\|\hat{\mathbf{A}} - \mathbf{A}^*\|_F$, for marginally stable systems, i.e., under the assumption that $\rho(A) \leq 1$ (Alquier et al., 2020; Simchowitz et al., 2018; Basu et al., 2019) and in the general case (Sarkar & Rakhlin, 2019).

Instead of directly considering the OLS estimator, we consider the empirical minimizer of the square loss under a restricted set of matrices \mathbf{A} that have a bounded operator norm:

$$\hat{\mathbf{A}} \in \operatorname{argmin}_{\mathbf{A}=(A_1, \dots, A_p)} \{ \mathcal{L}(\mathbf{A}) \mid \|M_{\mathbf{A}}\|_{\text{op}} \leq D \}. \quad (5)$$

Note that the set

$$\mathcal{A}(D) := \{ \mathbf{A} = (A_1, \dots, A_p) \mid \|M_{\mathbf{A}}\|_{\text{op}} \leq D \},$$

is bounded, closed and convex. Hence, the empirical minimizer of the square loss over $\mathcal{A}(D)$ can be computed with projected gradient descent (Duchi et al., 2008) or the Frank-Wolfe algorithm (Jaggi, 2013) as done for ℓ^1 constrained optimization. To avoid projecting onto the set $\mathcal{A}(D)$, following Proposition 3.5, it is possible to restrict $\mathcal{A}(D)$ further into

$$\mathcal{A}(D)' := \left\{ \mathbf{A} \mid \sum_{i=1}^p \|A_i\|_{\text{op}}^2 \leq D^2 \right\}, \quad \text{and} \quad \mathcal{A}(D)'' := \left\{ \mathbf{A} \mid \|\mathbf{A}\|_{\text{op}} \leq \frac{D}{\sqrt{p}} \right\}.$$

in order to ensure a condition directly on design matrices. Then, the empirical minimizer of the square loss over $\mathcal{A}(D)'$ or $\mathcal{A}(D)''$ can again be computed via projected gradient descent or the Frank-Wolfe algorithm, with simplified projection steps.

Lastly, we briefly remark that the diameter constraint in Equation (5) can be removed, i.e., $\mathcal{A}(D)$ replaced by $\mathcal{A}(\infty)$, under an additional assumption on NT . This is explained in detail in Section 5.

3.3 LOW-RANK ASSUMPTION

A common assumption in multi-task and meta-learning is that high-dimensional data often shares a representation in a smaller space (Bengio et al., 2013; Tripuraneni et al., 2021; Hospedales et al., 2022; Boursier et al., 2022; Collins et al., 2022; Yüksel et al., 2024). The following low-rank assumptions are crucial, as they significantly improve the statistical complexity of the problem.

Assumption 3.7. For all $k \in [p]$, $\text{rank}(A_k^*) \leq r$.

Assumption 3.8. There exists an orthonormal matrix $P^* \in \mathbb{R}^{r \times d}$ and matrices $B_1^*, \dots, B_p^* \in \mathbb{R}^{d \times r}$ such that $A_k^* = B_k^* P^*$ for all $k \in [p]$.

Note that Assumption 3.8 is an instance of Assumption 3.7. The factorization $A_k^* = Q^* C_k^*$ is another subcase of Assumption 3.7, but is not considered as it leads to iterates that directly lie in the subspace spanned by Q^* and hence Q^* can be learned by treating iterates $x_t^{(n)}$ as independent. In order to benefit from the low-rank structure, we consider the following regularized estimator:

$$\hat{\mathbf{A}} \in \operatorname{argmin}_{\mathbf{A} \in \mathcal{A}_r(D)} \mathcal{L}(\mathbf{A}), \quad \text{where } \mathcal{A}_r(D) := \{\mathbf{A} \in \mathcal{A}(D) \mid \forall k \in [p], \text{rank}(A_k) \leq r\}. \quad (6)$$

3.4 MISSPECIFICATION

The context length of the generative autoregressive process might be unbounded, too large for an efficient estimation, or apriori unknown. In any case, practitioners still have to set a context length $p' \in \mathbb{N}^*$ for the estimator, which might differ from the true p . In this scenario, we need an additional boundedness assumption that relates the first p' matrices of the ground truth.

Assumption 3.9. There exist a constant D' such that

$$\left\| (M_{\mathbf{A}^*} - M_{\mathbf{A}_{1:p'}^*}) L_* \right\|_{\text{op}} \leq D', \quad \text{where } \mathbf{A}_{1:p'}^* = (A_1^*, \dots, A_{p'}^*, 0_d, \dots, 0_d).$$

Instead of the estimator defined in Equation (6), we consider the following misspecified estimator:

$$\hat{\mathbf{A}} \in \operatorname{argmin}_{\mathbf{A} \in \mathcal{A}_{r,p'}(D)} \mathcal{L}(\mathbf{A}), \quad \text{where } \mathcal{A}_{r,p'}(D) := \{\mathbf{A} \in \mathcal{A}_r(D) \mid \forall p' < k \leq p, A_k = 0_d\}. \quad (7)$$

Assumption 3.9 is a strong assumption as it requires that L_* is well-behaved regardless of the sequence length T . Consequently, the misspecification results are more stringent than other results and apply to a smaller class of systems that still includes strictly stable systems as discussed in Section 5.

4 LONG-CONTEXT LINEAR SYSTEM IDENTIFICATION

In this section, we present statistical rates for the recovery of the design matrices in terms of Frobenius norm. Since the matrices \mathbf{A} lie in $\mathbb{R}^{d \times pd}$, the number of variables is pd^2 . In the *i.i.d.* setting, the rates of the form $\|\hat{\mathbf{A}} - \mathbf{A}^*\|_F^2 = \mathcal{O}(pd^2/(NT))$ are expected. The following theorem extends this rate for long-context linear dynamical system identification:

Theorem 4.1. Let Assumptions 3.1 and 3.4 hold. Then, for any $0 < \delta < e^{-1}$, there exists a constant $C(\delta) = \mathcal{O}(\ln \frac{1}{\delta})$ such that the estimator $\hat{\mathbf{A}}$ in Equation (5) verifies with probability $1 - \delta$:

$$\left\| \hat{\mathbf{A}} - \mathbf{A}^* \right\|_F^2 \leq C(\delta) D^2 \frac{\kappa^2 p d^2}{N(T-p)} \text{polylog}(p, d, N, T, \ln D). \quad (8)$$

The constant $C(\delta)$ depends mildly on the sub-Gaussianity constant c as described in Appendix B.5 and the rate is numerically verified in Figure 1. Theorem 4.1 exhibits several interesting features.

First, it shows that despite the temporal dependencies in the data, learning still occurs at a pace reminiscent of the *i.i.d.* setting, with a logarithmic term adjustment. This implies that the number of samples required to learn the system is approximately the same as in the *i.i.d.* setting, except for the logarithmic factor. Therefore, even though the data is sequential and only *i.i.d.* at the sequence level, the number of iterates $N(T-p)$ represents the *effective* data size.

Second, the rate in Equation (8) exhibits a linear dependency on the context length p instead of a quadratic dependency. This is only due to the number of parameters to be estimated, which is pd^2 .

instead of d^2 and not a deflation in T by a factor of p , which implies the context length does not affect the effective sample size. The additive factor in $T - p$ is due to the fact that first iterates do not depend on the full context length, and thus are not as informative as the later iterates. More detailed discussions of Theorem 4.1, in comparison with previous work, can be found in Section 5.

Low-rank setting. Next, we extend the results to the low-rank setting:

Theorem 4.2. *Let Assumptions 3.1, 3.4 and 3.7 hold. Then, for any $0 < \delta < e^{-1}$, there exists a constant $C(\delta) = \mathcal{O}(\ln \frac{1}{\delta})$ such that the estimator $\hat{\mathbf{A}}$ in Equation (6) verifies with probability $1 - \delta$:*

$$\|\hat{\mathbf{A}} - \mathbf{A}^*\|_F^2 \leq C(\delta) D^2 \frac{\kappa^2 p r d}{N(T - p)} \text{polylog}(p, d, r, N, T, \ln D). \quad (9)$$

The improved statistical rate depends on rd instead of d^2 . Note, however, that this estimator cannot be computed in polynomial time, since the underlying optimization problem involves a non-convex constraint on the rank of all A_k . Several heuristics exist to approximate this estimator. One approach is the Burer-Monteiro factorization (Burer & Monteiro, 2003; 2004), which involves parameterizing A_k as $A_k = B_k C_k$ with $B_k \in \mathbb{R}^{d \times r}$ and $C_k \in \mathbb{R}^{r \times d}$. This method relaxes the constraint to a convex set but results in a non-convex function. Another approach is *hard-thresholding* algorithms, which use projected (stochastic) gradient descent on the non-convex constraint set (Blumensath & Davies, 2009; Foucart & Subramanian, 2018).

Perhaps the most intuitive approach is to use nuclear norm regularization, which is a convex relaxation of the rank constraint:

$$\hat{\mathbf{A}} \in \text{argmin} \{ \mathcal{L}_\lambda(\mathbf{A}) \mid \mathbf{A} \in \mathcal{A}(D) \}, \text{ where } \mathcal{L}_\lambda(\mathbf{A}) = \mathcal{L}(\mathbf{A}) + \lambda \|\mathbf{A}\|_{*,\text{group}}, \quad (10)$$

and $\|\mathbf{A}\|_{*,\text{group}} = \sum_{k=1}^p \|A_k\|_*$ is the *group-nuclear norm*. We leave the analysis of the nuclear norm estimator for future work.

While the low-rank estimator cannot be computed easily, substituting the constraint $\forall k, \text{rank}(A_k) \leq r$ with $\text{rank}(\mathbf{A}) \leq r'$ enables a closed-form solution for the optimization problem (Bunea et al., 2011). However, the latter constraint effectively includes the former only when $r' \geq pr$, which would lead to suboptimal dependencies on the context length. These constraints are equivalent only if all A_k matrices project onto the same space: i.e., $A_k = QB_k$ for some $Q \in \mathbb{R}^{d \times r}$ and $B_k \in \mathbb{R}^{r \times r}$.

Misspecification. Lastly, we study linear long-context autoregressive prediction models under misspecified context lengths and show that partial learning still occurs for misspecified models:

Theorem 4.3. *Let Assumptions 3.1, 3.4, 3.7 and 3.9 hold. Then, for any $0 < \delta < e^{-1}$, there exists a constant $C(\delta) = \mathcal{O}(\ln \frac{1}{\delta})$ such that the estimator $\hat{\mathbf{A}}$ in Equation (7) verifies with probability $1 - \delta$:*

$$\|\hat{\mathbf{A}} - \mathbf{A}_{p'}^*\|_F^2 \leq C(\delta) D^2 (D' + 1)^2 \frac{\kappa^2 p' d r}{N(T - p)} \text{polylog}(p', d, r, N, T, \ln D). \quad (11)$$

For $r = d$, we recover Theorem 4.1 (full-rank setting) for misspecified context windows. The main improvement in that case of Theorem 4.3 over Theorem 4.1 is the dependency on p' instead of p . In practice, p can be much larger than p' and even on the order of T . In such a setting, learning all matrices A_k^* becomes impossible if N is not large enough and one does not take advantage of the length T of the sequences. One can instead misspecify the student with a context length of $p' \ll p$ such that $NT \gg p'd^2$, so that the first p' matrices are still learned.

Lastly, we briefly remark that Theorem 4.1 provides a rate for the case where $p < p'$. The latter case can be seen under a well-specified setting by rewriting the ground truth model as $\mathbf{A}^* = (A_1^*, \dots, A_p^*, 0_d, \dots, 0_d)$ where the last $p' - p$ indices are padded with null matrices. Learning in such a case is then answered by Theorem 4.1 with a worsened rate that depends on p' .

5 DISCUSSION

We now discuss the rates obtained in Section 4 and compare them with previous results obtained for linear dynamical systems. In particular, we comment the “learning-without-mixing” phenomenon, introduced by Simchowitz et al. (2018) for the first-order linear dynamical systems.

Adaptation of first-order techniques ($p = 1$) to the long-context setting. Here, we explain why techniques developed in the $p = 1$ setting, in particular those of (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019), do not work for the $p > 1$ setting and why, even if adapted, they would fail to achieve the desired sharp dependency on p .

Observe that the multi-step dynamics can be cast as a 1-step dynamic using *block companion matrices*. Let $X_t^{(n)} = (x_t^{(n),\top}, \dots, x_{t+p-1}^{(n),\top})^\top \in \mathbb{R}^{pd}$, $\Xi_t^{(n)} = (0, \dots, 0, (\xi_t^{(n)})^\top)^\top \in \mathbb{R}^{pd}$ and let $\mathcal{A}^* \in \mathbb{R}^{pd \times pd}$ be the companion matrix associated to \mathbf{A}^* :

$$\mathcal{A}^* = \begin{pmatrix} 0_d & I_d & \cdots & 0_d \\ \vdots & \ddots & \ddots & \vdots \\ 0_d & \cdots & 0_d & I_d \\ \mathbf{A}_p^* & \mathbf{A}_{p-1}^* & \cdots & \mathbf{A}_1^* \end{pmatrix}. \quad (12)$$

We have the relation $X_{t+1}^{(n)} = \mathcal{A}^* X_t^{(n)} + \Xi_t^{(n)}$, reducing the problem to the $p = 1$ case by increasing the dimension from d to pd . First, brute-force adapting previous results to this case (e.g. Basu et al., 2019; Simchowitz et al., 2018; Sarkar & Rakhlin, 2019) is not possible since these works assume that the noise covariance of the additive noise added at each step ($\Xi_t^{(n)}$ here) is the identity matrix, or at least is positive definite. In our case, the noise covariance is the $pd \times pd$ block-diagonal matrix, with $p - 1$ blocks equal to 0_d and the last one to I_d . The covariance matrix is thus non-invertible, preventing the use of previous works.

In addition, arguments based on system excitation (e.g. Basu et al., 2019; Simchowitz et al., 2018) are bound to incur an additional dependence on p , on top of the factors expected due to the dimensionality of the problem. In particular, as seen in the small-ball martingales argument by Simchowitz et al. (2018, Section 2.3), evaluating quantities like $\|(\mathcal{A} - \mathcal{A}^*)X_t^{(n)}\|^2$ for the (k, ν, q) -block martingale small-ball assumption requires $k \geq p$ as p represents the minimum number of steps for noise to propagate in every direction. Consequently, these analyses lead to a suboptimal p dependency.

Moreover, adapting the techniques developed in the $p = 1$ setting (Sarkar & Rakhlin, 2019) which relies on explicit factorization of the OLS estimator is challenging. In the $p > 1$ case, the higher-order dynamics complicate the factorization, and the data matrix takes a Toeplitz form, which is more difficult to handle.

Learning-without-mixing. We explain why our rates exhibit “learning-without-mixing”. We begin by defining “learning-with-mixing” and discussing the factors that influence the mixing time τ_{mix} . We then introduce the concept of “learning-without-mixing” as exemplified by Simchowitz et al. (2018) and show that our bounds exhibit similar properties.

Let τ_{mix} be the mixing time of the Markov chain $(X_t^{(n)})_{t \geq 0}$. In the *i.i.d.* setting (for which $\tau_{\text{mix}} = 1$), the OLS estimator obtains the optimal rate $\|\hat{\mathbf{A}}_{\text{OLS}} - \mathbf{A}^*\|_F^2 = \mathcal{O}(pd^2/NT)$, since pd^2 is the dimension of the inputs. With non-*i.i.d.* but Markovian data, a naive strategy would be to emulate *i.i.d.*-ness and take only a sample every τ_{mix} steps of the trajectory to compute the OLS estimator, thus having data that are approximately *i.i.d.* while dividing the number of samples by τ_{mix} . This naive “learning-with-mixing” estimator would yield $\|\hat{\mathbf{A}}_{\text{naive}} - \mathbf{A}^*\|_F^2 = \tilde{\mathcal{O}}(\tau_{\text{mix}}pd^2/NT)$, where the mixing time appears as a cost of non-*i.i.d.*-ness.

In our case, two components contribute to the mixing time, τ_{mix} . The first component is related to the *stability* or the *excitability* of the system and scales as $1/(1 - \rho)$, where $\rho = \|M_{\mathbf{A}^*}\|_{\text{op}} < 1$. When $\rho \ll 1$, this component has no impact, while ρ tends to 1, the system is less stable and the Markov chain mixes more slowly. The second component is directly related to the *context length* p of the process. Regardless of the factor $1/(1 - \rho)$ above, the mixing time of our Markov chain is larger than p : since noise is added only in the last block in the recursion $X_{t+1}^{(n)} = \mathcal{A}^* X_t^{(n)} + \Xi_t^{(n)}$, starting from a given state, p iterations at least are needed to eventually forget this given state. The naive *learning-with-mixing* benchmark rate is thus $\|\hat{\mathbf{A}} - \mathbf{A}^*\|_F^2 \leq \max(1/(1 - \rho), p)pd^2/NT$.

In contrast, a rate of convergence that exhibits “learning-without-mixing” is a rate of the form $\|\hat{\mathbf{A}} - \mathbf{A}^*\|_F^2 \leq Cp d^2/NT$ where $C \ll \tau_{\text{mix}}$. Such a rate means that the matrix \mathbf{A}^* is learned without paying the cost of non-*i.i.d.*-ness. For instance, in the $p = 1$ case, the rate of Simchowitz et al. (2018) does not worsen as ρ tends to 1—in fact, $\rho \rightarrow 1$ actually improves their rates.

The bound presented in Theorem 4.1 takes the form $\tilde{\mathcal{O}}(D^2 \kappa^2 p d^2 / (N(T - p)))$ where $\tilde{\mathcal{O}}$ hides the logarithmic terms. Importantly, the dependencies on the underlying Markov chain are only through D and κ , which do not have a direct dependency on the mixing time. The dependency on D is merely an operator norm upper bound and does not diverge as the mixing time grows to ∞ . Similarly, κ is logarithmic in T for systems of interest, as we discuss below.

System stability and κ . We now explain the behavior of κ defined in Definition 3.6. First, by Lemma B.11, we have that $\sigma_{\min}(L_\star) \geq \frac{1}{D+1}$ and, thus, it is sufficient to upper bound

$$\zeta(T) := \sup_{i,j \in [T]} \|L_\star^{(i,j)}\|_{\text{op}} \geq \sup_{i,j \in [T]} \frac{\|L_\star^{(i,j)}\|_F}{\sqrt{d}} \geq \frac{\|L_\star\|_F}{\sqrt{dT}} \geq \frac{\|L_\star\|_{\text{op}}}{\sqrt{dT}}, \quad (13)$$

to control κ . Equation (13) implies that if the noise at step i contributes to step j , as measured by $L_\star^{(i,j)}$, at a polynomial rate in $(j - i)$, then κ grows at most logarithmically in T . For such a κ , the resulting dependency on T is of order $\ln T$ and mild. Instead, if it is exponential in $(j - i)$, then κ grows linearly in T and the dependency on T cancels out in the rate.

We use the quantity $\zeta(T)$ to define *strictly stable*, *marginally stable* and *explosive* systems:

Definition 5.1. An LDS as defined in Equation (3) is called

$$\begin{aligned} \text{strictly stable if: } & \zeta(T) = \mathcal{O}(\rho^T) \quad \text{for some } \rho < 1, \\ \text{marginally stable if: } & \zeta(T) = \mathcal{O}(T^k) \quad \text{for some } k \in \mathbb{N}, \\ \text{explosive if: } & \zeta(T) = \mathcal{O}(\rho^T) \quad \text{for some } \rho > 1. \end{aligned}$$

Definition 5.1 is similar to the notions of strictly stable, marginally stable and explosive systems considered in (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019) for $p = 1$. Let $\rho(A^\star) := \lambda_{\max}(A^\star)$ be the spectral radius of A^\star and $V\Lambda V^{-1}$ be the Jordan normal form of A^\star . Then,

$$\|L_\star^{(i,j)}\|_{\text{op}} = \|(A^\star)^{j-i}\|_{\text{op}} = \|V\Lambda^{j-i}V^{-1}\|_{\text{op}} \leq \|V\|_{\text{op}} \|\Lambda^{j-i}\|_{\text{op}} \|V^{-1}\|_{\text{op}}.$$

Note that $\|V\|_{\text{op}}$ and $\|V^{-1}\|_{\text{op}}$ are constants. For upper bounding $\|\Lambda^{j-i}\|_{\text{op}}$, consider the Jordan blocks $\{\Lambda_k\}$ of Λ , associated with the eigenvalues λ_k of A^\star . Then, $\|\Lambda^{j-i}\|_{\text{op}} \leq \sup_k \|\Lambda_k^{j-i}\|_{\text{op}}$ and

$$\begin{aligned} \|\Lambda_k^{j-i}\|_{\text{op}} &= \|(\lambda_k I_n + N_n)^{j-i}\|_{\text{op}} = \left\| \sum_{m=0}^{\max\{j-i, n-1\}} \lambda_k^m \binom{j-i}{m} N_n^m \right\|_{\text{op}} \\ &\leq \sum_{m=0}^{\max\{j-i, n-1\}} \rho(A^\star)^{j-i} \binom{j-i}{m}, \end{aligned}$$

where n is the block size for the Jordan block Λ_k . Note here that n does not scale with T .

In particular, for *strictly stable* systems of Simchowitz et al. (2018); Sarkar & Rakhlin (2019) with $\rho < 1$, $\zeta(T) = \mathcal{O}(\rho^T)$. For *marginally stable* systems of Sarkar & Rakhlin (2019) with $\rho < 1 + \frac{\gamma}{T}$ with some constant $\gamma > 0$, $\rho(A^\star)^{j-i} \leq e^\gamma$ and $\zeta(T) = \mathcal{O}(T^k)$ for some fixed k that depends on the largest Jordan block of A^\star . For *explosive* systems of Sarkar & Rakhlin (2019) with $\rho > 1$, $\zeta(T) = \mathcal{O}(\rho^T)$. Thus, Definition 5.1 provides a general categorization of the systems based on the growth of $\zeta(T)$ in $p > 1$ case. Furthermore, our analysis yields sharp rates for *strictly stable* and *marginally stable* systems previously considered only in the $p = 1$ setting.

Search space diameter D . Our analysis is based on the assumption that the diameter D of the search space is bounded and, hence, not directly applicable to the OLS estimator in Equation (4). However, Corollary B.7 in Appendix B.1 extends the results of Theorems 4.1 to 4.3 to minimizers without a constraint on the diameter of the search space. This extension does not change the rates but requires the additional assumption that $NT = \tilde{\Omega}(p^2 dr)$.¹ In the case of Theorem 4.1, this corresponds to a result for the OLS estimator, but necessitating a number of samples quadratic in context length. Below, we comment on why the diameter restrictions is required when $NT \ll p^2 dr$.

¹We use the convention that $r = d, p' = p$ for Theorem 4.1.

As mentioned earlier in comparison with (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019), the simple OLS factorization in the $p = 1$ case does not generalize to the $p > 1$ and the data matrix has a Toeplitz structure that is more difficult to control. In order to deal with these issues, as explained in the sketch of proof in Section 6, we rely on techniques from empirical process theory. These techniques are applied to quantify the probability of the event in Equation (14), which hold for any empirical risk minimizer of the square loss. This leads us to the study of the concentration of the martingales defined in Equation (16) around their predictable variation, which is a key step in our analysis. A uniform concentration is possible only if there is a uniform lower bound on the variations of the martingales, which can be achieved using a set of well-behaved matrices $\|M_{\mathbf{A}} - M_{\mathbf{A}^*}\|_F / \|M_{\mathbf{A}} - M_{\mathbf{A}^*}\|_{\text{op}}$. In order to translate these conditions on the design matrices without additional dimensional dependencies, we introduce the operator norm constraint.

Lastly, it is possible to extend our analysis to unconstrained OLS by establishing a general coarse upper bound on the operator norm $\|M_{\hat{\mathbf{A}}}\|_{\text{op}} \leq K$. This allows us to consider uniform lower bounds to matrices \mathbf{A} with $\|M_{\mathbf{A}}\|_{\text{op}} \leq K$, which lead to a rate for the OLS estimator in a similar manner.

Upper bound on D' . The misspecification result in Theorem 4.3 requires the additional assumption given in Assumption 3.9. In Remark B.4, we show that a good upper bound on D' is possible when $D < 1$, i.e., the system is strictly stable, by using the bound $\|L_{\star}\|_{\text{op}} \leq 1/(1 - D)$. However, misspecification results are not, a priori, applicable to marginally stable systems, which limits the practical applicability of our results. We leave the investigation of misspecification results for marginally stable systems for future work.

6 SKETCH OF PROOF

We provide a sketch of proof for Theorem 4.1. The proofs of Theorems 4.2 and 4.3 are similar and can be found in Appendix B. In the following, $\Delta_{\mathbf{A}}$ is a shorthand for $M_{\mathbf{A}} - M_{\mathbf{A}^*}$ and $E \in \mathbb{R}^{Td \times N}$ is the matrix that collects the noise concatenated over time, as explained in Definition B.1.

The empirical risk minimizer $\hat{\mathbf{A}}$ satisfies the following optimality condition:

$$\mathcal{L}(\hat{\mathbf{A}}) \leq \mathcal{L}(\mathbf{A}^*), \quad \text{or written differently,} \quad \|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|^2 \leq 2 \text{Tr} (E^{\top} L_{\star}^{\top} \Delta_{\hat{\mathbf{A}}}^{\top} E), \quad (14)$$

due to the *well-specified* setting, i.e., $\mathbf{A}^* \in \mathcal{A}(D)$. The condition in Equation (14) is of interest as

$$\forall \mathbf{A} \in \mathcal{A}(D), \quad \mathbb{E} [\|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|^2] = \sigma^2 N \|\Delta_{\hat{\mathbf{A}}} L_{\star}\|_F^2 > 0 = \mathbb{E} [\text{Tr} (E^{\top} L_{\star}^{\top} \Delta_{\hat{\mathbf{A}}}^{\top} E)].$$

This inequality hints that if for a set of matrices $\mathcal{A}'(D) \subseteq \mathcal{A}(D)$, there is a uniform result

$$\mathcal{E} := \left\{ \forall \mathbf{A} \in \mathcal{A}'(D) : \|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|^2 \geq 2 \text{Tr} (E^{\top} L_{\star}^{\top} \Delta_{\hat{\mathbf{A}}}^{\top} E) \right\}, \quad (15)$$

with high probability as seen from their means, then the empirical risk minimizer $\hat{\mathbf{A}}$ belongs to the set $\mathcal{A}(D) \setminus \mathcal{A}'(D)$ with the same high probability by a simple Bayesian argument. Hence, the proof of Theorem 4.1 is reduced to proving Equation (15) for a suitable set of matrices.

Fix a $\mathbf{A} \in \mathcal{A}'(D)$ and study the martingale series defined through the differences sequences

$$d_{t,i}^{(n)} = \left((\mathbf{A} - \mathbf{A}^*) x_t^{(n)} \right)_i \left(\xi_t^{(n)} \right)_i / \sigma^2, \quad (16)$$

where the series is first ordered in i , then in t , and finally in n . The sum of the differences is then

$$Y_{\mathbf{A}} = \sum_{n,t,i} d_{t,i}^{(n)} = \sum_{n,t} \left\langle (\mathbf{A} - \mathbf{A}^*) x_t^{(n)}, \xi_t^{(n)} \right\rangle = \frac{1}{\sigma^2} \text{Tr} (E^{\top} L_{\star}^{\top} \Delta_{\mathbf{A}}^{\top} E),$$

and the quadratic predictable variation of the series is

$$W_{\mathbf{A}} = \sum_{n,t,i} \mathbb{E}_{(\xi_t^{(n)})_i} \left[\left(d_{t,i}^{(n)} \right)^2 \right] = \frac{1}{\sigma^2} \sum_{n,t} \left\| (\mathbf{A}^* - \mathbf{A}) x_t^{(n)} \right\|^2 = \frac{1}{\sigma^2} \|\Delta_{\mathbf{A}} L_{\star} E\|^2.$$

The condition that is asked in Equation (15) is then that the sum of the differences $Y_{\mathbf{A}}$ is large compared to the quadratic predictable variation $W_{\mathbf{A}}$, i.e., $\mathcal{E} = \{ \forall \mathbf{A} \in \mathcal{A}'(D) : W_{\mathbf{A}} \leq 2Y_{\mathbf{A}} \}$.

In order to prove probabilistic statements on \mathcal{E} , we use Freedman's inequality (Freedman, 1975; Dzhaparidze & Van Zanten, 2001) which gives control on $Y_{\mathbf{A}}$ and $W_{\mathbf{A}}^R$ for a particular \mathbf{A} :

$$\mathbb{P}(Y_{\mathbf{A}} \geq r_Y, W_{\mathbf{A}}^R \leq r_W) \leq \exp\left(-\frac{r_Y^2/2}{r_W + Rr_Y}\right), \quad (17)$$

where $W_{\mathbf{A}}^R = W_{\mathbf{A}} + \sum_{n,t,i} \mathbb{1}_{d_{t,i}^{(n)} > R} \left(d_{t,i}^{(n)}\right)^2$ and $r_Y, r_W, R > 0$ are arbitrary constants. As the noise is sub-Gaussian, it is possible to upper bound $W_{\mathbf{A}}^R$ with $W_{\mathbf{A}}$:

$$\sum_{n,t,i} \left(d_{t,i}^{(n)}\right)^2 \leq \sup_{n,t,i} \left(\xi_t^{(n)}\right)_i \cdot \sum_{n,t} \left\| (\mathbf{A} - \mathbf{A}^*) x_t^{(n)} \right\|^2 \xrightarrow{\text{w.h.p}} \forall \mathbf{A} : W_{\mathbf{A}}^R \leq (1 + 2c'^2 \sigma^2 \ln dTN) W_{\mathbf{A}}.$$

Further, assume that there are uniform upper and lower bounds on $Y_{\mathbf{A}}$ and $W_{\mathbf{A}}$, respectively:

$$\exists 0 < \alpha_L < \alpha_U \text{ such that } \forall \mathbf{A} \in \mathcal{A}'(D) : Y_{\mathbf{A}} \leq \alpha_U \text{ and } \alpha_L \leq W_{\mathbf{A}} \leq W_{\mathbf{A}}^R. \quad (18)$$

Then, letting $\gamma = 2(1 + 2c'^2 \sigma^2 \ln dTN)$ and k' such that $R^{k'} \alpha_L \geq \gamma \alpha_U$, we have

$$\mathbb{P}(W_{\mathbf{A}} \leq 2\gamma Y_{\mathbf{A}}) \leq \mathbb{P}(W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}) \leq \bigcup_{k=1}^{k'} \mathbb{P}(W_{\mathbf{A}}^R \leq R^k \alpha_L, \gamma Y_{\mathbf{A}} \geq R^{k-1} \alpha_L).$$

Each of the terms in the union can be controlled by the Freedman's inequality in Equation (17) with the choices of $r_Y = \alpha_L R^{k-1}$ and $r_W = \alpha_L R^k$:

$$\mathbb{P}(W_{\mathbf{A}} \leq 2Y_{\mathbf{A}}) \leq \sum_{k=1}^{k'} \exp(-\alpha_L R^{k-2}/4) \leq \exp\left(-\frac{\alpha_L}{4R} + \ln \ln \frac{2e\alpha_U}{R\alpha_L}\right). \quad (19)$$

As can be seen from Equation (19), the probability of the event $\{W_{\mathbf{A}} \leq 2Y_{\mathbf{A}}\}$ is largely controlled with the lower bound α_L as the ratio α_U/α_L only matters logarithmically. This is crucial as the two bounds differ with the condition number κ of the linear operator L_* , which can scale with T .

Therefore, it is possible to control the event \mathcal{E} with a union bound over an ϵ -net of $\mathcal{A}'(D)$. In particular, α_L needs to be uniformly bounded below such that α_L/R is of scale $\ln |\mathcal{N}_{\epsilon}(\mathcal{A}'(D))|$. And, this is achieved by Hanson-Wright inequality (Hanson & Wright, 1971) which allows us to derive the needed uniform lower and upper bounds in Equation (18) with high probability as long as $\mathcal{A}'(D)$ is composed of matrices that satisfy

$$\frac{\|\Delta_{\mathbf{A}}\|_F^2}{\|\Delta_{\mathbf{A}}\|_{\text{op}}^2} \geq \ln |\mathcal{N}_{\epsilon}(\mathcal{A}'(D))|, \quad \text{where } \epsilon \sim \frac{\text{polysqrt}(p, d, N, T, e^{\kappa})}{1 + c^2 \ln \frac{1}{\delta}}.$$

Here, we pick up the dependency on κ as ϵ scales with κ . This is needed to bound the worst-case errors while transitioning from point-wise bounds on the ϵ -net $\mathcal{N}_{\epsilon}(\mathcal{A}'(D))$ to the whole set $\mathcal{A}'(D)$.

Finally, since there is a uniform bound on $\|\Delta_{\mathbf{A}}\|_{\text{op}} \leq 2D$ implied by Assumption 3.4, setting

$$\mathcal{A}'(D) = \left\{ \mathbf{A} \mid \|\Delta_{\mathbf{A}}\|_F^2 \geq CD^2 \frac{pdr}{N} \text{polylog}(p, d, N, T, \kappa) \right\},$$

for some $C = \mathcal{O}(\frac{1}{\delta})$ is sufficient to deduce $\hat{\mathbf{A}} \in \mathcal{A}(D) \setminus \mathcal{A}'(D)$ with probability $1 - \delta$. The proof of Theorem 4.1 is then complete as $\|\Delta_{\mathbf{A}}\|_F^2 = \|\mathbf{M}_{\mathbf{A}} - \mathbf{M}_{\mathbf{A}^*}\|_F^2 \geq (T - p)\|\mathbf{A} - \mathbf{A}^*\|_F^2$.

7 CONCLUSION

In this work, we extend non-asymptotic linear system identification theory and derive upper bounds on the sample complexity of learning long-context linear autoregressive models. Our bounds improve upon the existing arguments specific to first-order systems by employing a uniform concentration argument over prediction differences. We further establish improved statistical rates when learning under a low-rank assumption. Finally, we show that even with long or unbounded generative contexts, *misspecification* still allows the estimation of the matrices with a reduced sample complexity and for stable systems.

While this work makes significant progress for non-asymptotic linear system identification theory, several technical questions remain open for further investigation. Can the OLS operator norm be coarsely controlled to derive rates for unconstrained OLS in the $NT = \Omega(pdr)$ regime? Is it possible to find efficient algorithms that would benefit from low-rank assumptions? Lastly, can misspecification be beneficial for marginally stable systems?

REFERENCES

- Pierre Alquier, Karine Bertin, Paul Doukhan, and Rémy Garnier. High-dimensional VAR with low-rank transition. *Statistics and Computing*, 30(4):1139–1153, March 2020. doi: 10.1007/s11222-020-09929-7. URL <https://doi.org/10.1007/s11222-020-09929-7>.
- Francesco Bacchiocchi, Gianmarco Genalti, Davide Maran, Marco Mussi, Marcello Restelli, Nicola Gatti, and Alberto Maria Metelli. Autoregressive bandits, 2024.
- Sumanta Basu, Xianqi Li, and George Michailidis. Low rank and structured modeling of high-dimensional vector autoregressions. *IEEE Transactions on Signal Processing*, 67(5):1207–1222, 2019. doi: 10.1109/TSP.2018.2887401.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, aug 2013. ISSN 0162-8828. doi: 10.1109/TPAMI.2013.50. URL <https://doi.org/10.1109/TPAMI.2013.50>.
- Thomas Blumensath and Mike E. Davies. Iterative hard thresholding for compressed sensing. *Applied and Computational Harmonic Analysis*, 27(3):265–274, 2009. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2009.04.002>. URL <https://www.sciencedirect.com/science/article/pii/S1063520309000384>.
- Etienne Boursier, Mikhail Konobeev, and Nicolas Flammarion. Trace norm regularization for multi-task learning with scarce data. In *Conference on Learning Theory*, pp. 1303–1327. PMLR, 2022.
- George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time Series Analysis: Forecasting and Control*. John Wiley & Sons, 2015.
- Peter J Brockwell and Richard A Davis. *Time Series: Theory and Methods*. Springer science & business media, 1991.
- Florentina Bunea, Yiyuan She, and Marten H. Wegkamp. Optimal selection of reduced rank estimators of high-dimensional matrices. *The Annals of Statistics*, 39(2):1282–1309, 2011. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/29783674>.
- Samuel Burer and Renato D.C. Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, February 2003. ISSN 1436-4646. doi: 10.1007/s10107-002-0352-8. URL <http://dx.doi.org/10.1007/s10107-002-0352-8>.
- Samuel Burer and Renato D.C. Monteiro. Local minima and convergence in low-rank semidefinite programming. *Mathematical Programming*, 103(3):427–444, December 2004. ISSN 1436-4646. doi: 10.1007/s10107-004-0564-1. URL <http://dx.doi.org/10.1007/s10107-004-0564-1>.
- Emmanuel J. Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359, 2011. doi: 10.1109/TIT.2011.2111771.
- Liam Collins, Aryan Mokhtari, Sewoong Oh, and Sanjay Shakkottai. MAML and ANIL provably learn representations. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4238–4310. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/collins22a.html>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

- John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, pp. 272–279, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi: 10.1145/1390156.1390191. URL <https://doi.org/10.1145/1390156.1390191>.
- Kacha Dzharidze and JH Van Zanten. On bernstein-type inequalities for martingales. *Stochastic processes and their applications*, 93(1):109–117, 2001.
- Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander Toshev, Vaishaal Shankar, Joshua M Susskind, and Armand Joulin. Scalable pre-training of large autoregressive image models, 2024.
- Mathieu Even. Stochastic gradient descent under markovian sampling schemes. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.
- Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Finite time identification in unstable linear systems, 2018.
- Dylan Foster, Tuhin Sarkar, and Alexander Rakhlin. Learning nonlinear dynamical systems from a single trajectory. In Alexandre M. Bayen, Ali Jadbabaie, George Pappas, Pablo A. Parrilo, Benjamin Recht, Claire Tomlin, and Melanie Zeilinger (eds.), *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of Machine Learning Research*, pp. 851–861. PMLR, 10–11 Jun 2020. URL <https://proceedings.mlr.press/v120/foster20a.html>.
- Simon Foucart and Srinivas Subramanian. Iterative hard thresholding for low-rank recovery from rank-one projections, 2018.
- David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pp. 100–118, 1975.
- Christophe Giraud, François Roueff, and Andres Sanchez-Perez. Aggregation of predictors for nonstationary sub-linear processes and online adaptive forecasting of time varying autoregressive processes. *The Annals of Statistics*, 43(6):2412–2450, 2015. ISSN 00905364. URL <http://www.jstor.org/stable/43818856>.
- N. K. Gupta, R. K. Mehra, and W. E. Hall. Application of optimal input synthesis to aircraft parameter identification. *Journal of Dynamic Systems, Measurement, and Control*, 98(2):139–145, June 1976. ISSN 1528-9028. doi: 10.1115/1.3427000. URL <http://dx.doi.org/10.1115/1.3427000>.
- James D Hamilton. *Time Series Analysis*. Princeton university press, 2020.
- D. L. Hanson and F. T. Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971. ISSN 00034851. URL <http://www.jstor.org/stable/2240253>.
- Moritz Hardt, Tengyu Ma, and Benjamin Recht. Gradient descent learns linear dynamical systems. *J. Mach. Learn. Res.*, 19(1):1025–1068, jan 2018. ISSN 1532-4435.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- Elad Hazan, Karan Singh, and Cyril Zhang. Learning linear dynamical systems via spectral filtering. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances*

- in *Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 6702–6712, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/165a59f7cf3b5c4396ba65953d679f17-Abstract.html>.
- T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis; Machine Intelligence*, 44(09):5149–5169, sep 2022. ISSN 1939-3539. doi: 10.1109/TPAMI.2021.3079209.
- Daniel Hsu, Sham M. Kakade, and Tong Zhang. Random design analysis of ridge regression. In Shie Mannor, Nathan Srebro, and Robert C. Williamson (eds.), *Proceedings of the 25th Annual Conference on Learning Theory*, volume 23 of *Proceedings of Machine Learning Research*, pp. 9.1–9.24, Edinburgh, Scotland, 25–27 Jun 2012. PMLR. URL <https://proceedings.mlr.press/v23/hsul2.html>.
- Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 427–435, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/jaggi13.html>.
- Prateek Jain, Suhas S. Kowshik, Dheeraj Nagaraj, and Praneeth Netrapalli. Streaming linear system identification with reverse experience replay. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 30140–30152, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/fd2c5e4680d9a01dba3aada5ece22270-Abstract.html>.
- Yassir Jedra and Alexandre Proutiere. Sample complexity lower bounds for linear system identification. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 2676–2681, 2019. doi: 10.1109/CDC40024.2019.9029303.
- Vladimir R Kostic, Pietro Novelli, Andreas Maurer, Carlo Ciliberto, Lorenzo Rosasco, and massimiliano pontil. Learning dynamical systems via koopman operator regression in reproducing kernel hilbert spaces. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=Tz11knIPVfp>.
- Lennart Ljung. *System Identification*, pp. 163–173. Birkhäuser Boston, Boston, MA, 1998. ISBN 978-1-4612-1768-8. doi: 10.1007/978-1-4612-1768-8_11. URL https://doi.org/10.1007/978-1-4612-1768-8_11.
- Horia Mania, Michael I. Jordan, and Benjamin Recht. Active learning for nonlinear system identification with guarantees. *J. Mach. Learn. Res.*, 23(1), jan 2022. ISSN 1532-4435.
- Thomas M. Moerland, Joost Broekens, Aske Plaat, and Catholijn M. Jonker. Model-based reinforcement learning: A survey, 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Mark Rudelson and Roman Vershynin. Hanson-Wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18(none):1 – 9, 2013. doi: 10.1214/ECP.v18-2865. URL <https://doi.org/10.1214/ECP.v18-2865>.

- Tuhin Sarkar and Alexander Rakhlin. Near optimal finite time identification of arbitrary linear dynamical systems. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5610–5618. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/sarkar19a.html>.
- Max Simchowitz, Horia Mania, Stephen Tu, Michael I. Jordan, and Benjamin Recht. Learning without mixing: Towards a sharp analysis of linear system identification. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet (eds.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 439–473. PMLR, 06–09 Jul 2018. URL <https://proceedings.mlr.press/v75/simchowitz18a.html>.
- Alex Simpkins. System identification: Theory for the user, 2nd edition (Ijung, I.; 1999). *IEEE Robotics and Automation Magazine*, 19(2):95–96, 1999. doi: 10.1109/MRA.2012.2192817.
- Nilesh Tripuraneni, Chi Jin, and Michael Jordan. Provable meta-learning of linear representations. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10434–10443. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/tripuraneni21a.html>.
- Anastasios Tsiamis, Ingvar Ziemann, Nikolai Matni, and George J Pappas. Statistical learning theory for control: A finite-sample perspective. *IEEE Control Systems Magazine*, 43(6):67–97, 2023.
- Stephen Tu, Roy Frostig, and Mahdi Soltanolkotabi. Learning from many trajectories, 2023.
- Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009. ISBN 1441412697.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi: 10.1017/9781108627771.
- F. T. Wright. A Bound on Tail Probabilities for Quadratic Forms in Independent Random Variables Whose Distributions are not Necessarily Symmetric. *The Annals of Probability*, 1(6):1068 – 1070, 1973. doi: 10.1214/aop/1176996815. URL <https://doi.org/10.1214/aop/1176996815>.
- Oğuz Kaan Yüksel, Etienne Boursier, and Nicolas Flammarion. First-order ANIL provably learns representations despite overparametrisation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=if2vRbS8Ew>.
- Ingvar Ziemann and Stephen Tu. Learning with little mixing. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

ORGANIZATION OF THE APPENDIX

The appendix is organized as follows,

- In Appendix A, we provide preliminary tools needed for our analyses: Hanson-Wright and Freedman inequalities, supremum of sub-Gaussian processes and proof of Proposition 3.5.
- In Appendix B, we prove Theorems 4.1 to 4.3 jointly under Theorem B.5.
- In Appendix C, we provide numerical experiments to verify our theoretical findings.

A PRELIMINARY TOOLS

A.1 HANSON-WRIGHT INEQUALITY

We use Hanson-Wright inequality (Hanson & Wright, 1971; Wright, 1973; Rudelson & Vershynin, 2013) to show concentration of certain second-order terms.

Theorem A.1. (*Hanson-Wright*) Let $Z = (Z_1, \dots, Z_n) \in \mathbb{R}^n$ be a random vector with independent components Z_i which satisfy $\mathbb{E}[Z_i] = 0$ and $\|Z_i\|_{\psi_2} \leq K$. Let P be an $n \times n$ matrix. Then, for every $r \geq 0$,

$$\mathbb{P}(|Z^\top P Z - \mathbb{E}[Z^\top P Z]| > r) \leq 2 \exp\left(-C_{HW} \min\left(\frac{r^2}{K^4 \|P\|_F^2}, \frac{r}{K^2 \|P\|_{\text{op}}}\right)\right).$$

The bound can be turned into a one-sided bound by dropping the constant 2.

Remark A.2. For data regime considered in this paper, $K = c\sigma$ in Theorem A.1.

A.2 FREEDMAN'S INEQUALITY

We use an extension of Freedman's inequality (Freedman, 1975) to non-bounded differences by Dzhaparidze & Van Zanten (2001) to show concentration of certain second-order terms. For the sake of completeness, we provide the original Freedman's inequality. We also remark that it is possible to use the original Freedman's inequality in our proofs to deal with *any bounded noise*.

Theorem A.3. (*Freedman's inequality*) Let Y_0, \dots, Y_n be a real-valued martingale series that is adapted to the filtration $\mathcal{F}_0, \dots, \mathcal{F}_n$ where $Y_0 = 0$. Let d_1, \dots, d_n be the difference sequence induced, i.e.,

$$d_i = Y_i - Y_{i-1} \quad \text{for } i = 1, \dots, n.$$

Assume that d_i is upper bounded by some R , i.e., $|d_i| \leq R$ for all i . Let W_i be the quadratic variation of the martingale series, i.e.,

$$W_i = \sum_{j=1}^i \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \quad \text{for } i = 1, \dots, n.$$

Then, for any $r, W > 0$,

$$\mathbb{P}(\exists k \geq 0 : Y_k \geq r \text{ and } W_k \leq W) \leq \exp\left(-\frac{r^2/2}{W + Rr}\right).$$

Theorem A.4. (*Freedman's inequality with non-bounded differences*) Let Y_0, \dots, Y_n be a real-valued martingale series that is adapted to the filtration $\mathcal{F}_0, \dots, \mathcal{F}_n$ where $Y_0 = 0$. Let d_1, \dots, d_n be the difference sequence induced, i.e.,

$$d_i = Y_i - Y_{i-1} \quad \text{for } i = 1, \dots, n.$$

Let W_i^R be the quadratic variation of the martingale series plus an error term for large differences,

$$W_i^R = \sum_{j=1}^i \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] + d_i^2 \mathbb{1}_{\{|d_i| > R\}} \quad \text{for } i = 1, \dots, n. \quad (20)$$

We set $W_i = W_i^0$ for ease of notation. Then, for any $r, R, W > 0$,

$$\mathbb{P}(\exists k \geq 0 : Y_k \geq r \text{ and } W_k^R \leq W) \leq \exp\left(-\frac{r^2/2}{W + Rr}\right).$$

We extend these theorems in Lemma A.5 to compare the quadratic variation with the martingale series itself. This is useful in our proofs to show certain events necessarily implied by empirical risk minimization do not occur with high probability.

Lemma A.5. *Let \mathcal{E} be the following event*

$$\mathcal{E} = \{W_n^R \geq \alpha_L\} \cap \{Y_n \leq \alpha_U\},$$

where $0 \leq \alpha_L \leq \alpha_U$, $R > 1$ are constants and W_n^R and Y_n are defined as in Theorem A.4. Then, for any $\gamma > 0$, we have the following concentration inequality

$$\mathbb{P}(\{W_n^R \leq \gamma Y_n\} \cap \mathcal{E}) \leq \exp\left(-\frac{\alpha_L}{2R\gamma} + \ln\left(\ln\left(\frac{\alpha_U}{R\alpha_L}\right) + 1\right)\right).$$

Proof. Without loss of generality, we assume $\gamma = 1$ as it is possible to scale the martingale series with $\frac{1}{\gamma}$ and set $R \leftarrow \frac{1}{\gamma}R$, $\alpha_L \leftarrow \frac{1}{\gamma^2}\alpha_L$ and $\alpha_U \leftarrow \frac{1}{\gamma}\alpha_U$.

Let $\mathcal{G} = \{\alpha_L, R\alpha_L, \dots, R^k\alpha_L\}$ where k is the smallest positive integer such that

$$R^k\alpha_L \geq \alpha_U.$$

Then, by a union bound,

$$\begin{aligned} \mathbb{P}(W_n^R \leq Y_n \cap \mathcal{E}) &\leq \mathbb{P}(\cup_{i=1}^k (\{W_n^R \leq R^i\alpha_L, Y_n \geq R^{i-1}\alpha_L\} \cap \mathcal{E})) \\ &\leq \mathbb{P}(\cup_{i=1}^k (\{W_n^R \leq R^i\alpha_L, Y_n \geq R^{i-1}\alpha_L\})) \\ &\leq \sum_{i=1}^k \mathbb{P}(\{W_n^R \leq R^i\alpha_L, Y_n \geq R^{i-1}\alpha_L\}). \end{aligned}$$

By applying Theorem A.4 with $r = R^{i-1}\alpha_L$ and $W = R^i\alpha_L$, we obtain

$$\mathbb{P}(W_n^R \leq R^i\alpha_L, Y_n \geq R^{i-1}\alpha_L) \leq \exp(-\alpha_L R^{i-2}/4),$$

for each $i = 1, \dots, k$. The result follows by noting that

$$\begin{aligned} \sum_{i=1}^k \mathbb{P}(\{W_n^R \leq R^i\alpha_L, Y_n \geq R^{i-1}\alpha_L\}) &\leq \sum_{i=1}^k \exp(-\alpha_L R^{i-2}/4) \\ &\leq \exp\left(-\frac{\alpha_L}{4R} + \ln k\right) \\ &\leq \exp\left(-\frac{\alpha_L}{4R} + \ln\left(\ln\left(\frac{\alpha_U}{R\alpha_L}\right) + 1\right)\right). \end{aligned}$$

□

A.3 SUPREMUM OF THE NOISE

We need the following lemma to control the supremum of the noise in our proofs.

Lemma A.6. *Let X_1, \dots, X_n be i.i.d. mean zero and σ^2 -sub-Gaussian random variables (in the sense provided in Assumption 3.1). Then, there exist a universal constant c' such that for any $t > 0$,*

$$\mathbb{P}\left(\sup_{i=1, \dots, n} X_i > c'\sigma\sqrt{2\ln n} + t\right) \leq 2\exp\left(-\frac{t^2}{2c'^2\sigma^2}\right).$$

Proof. By the sub-Gaussian property, we have a universal constant c' such that

$$\mathbb{P}(X_i \geq r) \leq \exp\left(-\frac{r^2}{2c'^2\sigma^2}\right).$$

Then, by the union bound,

$$\mathbb{P}\left(\sup_{i=1, \dots, n} X_i \geq r\right) \leq \cup_i \mathbb{P}(X_i \geq r) \leq n\exp\left(-\frac{r^2}{2c'^2\sigma^2}\right).$$

The result follows by setting $r = c'\sigma\sqrt{2\ln n} + t$.

□

Corollary A.7. For any $\delta > 0$, there exists a universal constant $c'(\delta) = \mathcal{O}\left(\sqrt{\ln \frac{1}{\delta}}\right)$ such that

$$\sup_{t,n} \|\xi_t^{(n)}\|_\infty \leq c'(\delta) \sigma \sqrt{2 \ln dTN},$$

with probability $1 - \delta$.

Proof. Each component $(\xi_t^{(n)})_i$ are i.i.d. of each other and sub-Gaussian with parameter σ . Therefore, by Lemma A.6, we have

$$\mathbb{P}\left(\sup_{t,n} \|\xi_t^{(n)}\|_\infty > c' \sigma \sqrt{2 \ln dTN} + r\right) \leq dNT \exp\left(-\frac{r^2}{2c'^2 \sigma^2}\right).$$

Select $r = c' \sigma \sqrt{2 (\ln dTN + \ln \frac{1}{\delta})}$ to obtain the desired confidence level of δ . Note that

$$r \leq c' \sigma \sqrt{2} \left(\sqrt{\ln dTN} + \sqrt{\ln \frac{1}{\delta}} \right),$$

and the constant $c'(\delta)$ need to satisfy

$$c'(\delta) \leq 1 + \sqrt{\frac{\ln \frac{1}{\delta}}{\ln dTN}} \leq 1 + \sqrt{\frac{1}{\ln 2}} \sqrt{\ln \frac{1}{\delta}},$$

as $dTN > p \geq 2$. Thus, $c'(\delta)$ can be picked such that it is a universal constant in δ . \square

A.4 PROOF OF PROPOSITION 3.5

Proof. Using Lemma B.12, we have $\|M_{A^*}\|_{\text{op}} \leq \sqrt{p} \|A^*\|_{\text{op}}$, directly leading to (ii).

For (i), we have

$$\|M_{A^*}\|_{\text{op}} \leq \sum_{i=1}^p \|M_{A^{*,(i)}}\|_{\text{op}}, \quad \text{where } A^{*,(i)} = \left(\underbrace{0_d, 0_d, \dots, 0_d}_{i-1 \text{ times}}, A^{*,i}, 0_d, \dots, 0_d \right).$$

Then, it is easy to see that

$$\|M_{A^{*,(i)}}\|_{\text{op}} \leq \|A_i^*\|_{\text{op}}.$$

\square

B PROOF OF THEOREMS 4.1 TO 4.3

Before proving the main theorems, we recall certain definitions from the main body of the paper:

Definition B.1. For any $A \in \mathbb{R}^{d \times pd}$, let $\Delta_A = \Delta_{A,p'}$ be defined as follows

$$\Delta_{A,i} = (M_{A_i} - M_{A_i^*}),$$

where $A_i = (A_1, \dots, A_i, 0_d, \dots, 0_d)$, $A_i^* = (A_1^*, \dots, A_i^*, 0_d, \dots, 0_d)$.

Let $\xi^{(i)} \in \mathbb{R}^{Td}$ be the whole noise concatenated in time, i.e.,

$$\xi^{(i)} = \left(\xi_1^{(i)}, \dots, \xi_T^{(i)} \right),$$

and let $E \in \mathbb{R}^{Td \times N}$ be the matrix that collects the noise for all sequences, i.e.,

$$E = \left(\xi^{(1)}, \dots, \xi^{(n)} \right).$$

Proposition B.2. With the definitions of Definition B.1, we have the following properties:

$$\begin{aligned} \sum_{n,t} \langle (A_i - A_i^*) X_t^{(n)}, \xi_t^{(n)} \rangle &= \text{Tr}(E^\top \Delta_{A,i} L_* E), \\ \sum_{n,t} \left\| (A_i - A_i^*) X_t^{(n)} \right\|^2 &= \|(M_{A_i} - M_{A_i^*}) L_* E\|_F^2 = \|\Delta_{A,i} L_* E\|_F^2. \end{aligned}$$

Definition B.3. Let $\mathcal{A}_{r,p}(D)$ and $\mathcal{S}_{r,p}(C, D)$ be the search and solution set for constants $C, D \geq 1$:

$$\mathcal{A}_{r,p'}(D) = \left\{ \mathbf{A} \in \mathbb{R}^{d \times pd} \mid \|\Delta_{\mathbf{A}}\|_{\text{op}} \leq D, \text{rank}(\mathbf{A}_i) \leq r, A_{p'+1} = \dots = A_p = 0 \right\},$$

$$\mathcal{S}_{r,p'}(C, D) = \left\{ \mathbf{A} \in \mathcal{A}(D) \mid \|\mathbf{A} - \mathbf{A}_{p'}^*\|_F^2 \leq CD^2 \eta^2 \frac{p' dr (\ln \tau)^2}{N(T - p')} \right\},$$

where η is a constant that captures an additional factor for the misspecified setting,

$$\eta = \begin{cases} 1 & \text{if } p' = p, \\ \max \left\{ 1, 1 + \left\| \left(M_{\mathbf{A}^*} - M_{\mathbf{A}_{p'}^*} \right) L_{\star} \right\|_{\text{op}} \right\} & \text{if } p' < p, \end{cases}$$

and τ is the following term:

$$\tau = e\sigma_{\text{cond}}(L_{\star}) \sqrt{p' d N T} \sqrt{\frac{T}{T - p'}}.$$

Let $\mathcal{G}_{r,p'}(C, D)$ be defined as follows,

$$\mathcal{G}_{r,p'}(C, D) = \left\{ \mathbf{A} \in \mathcal{A}_{r,p}(D) \mid \frac{\|\Delta_{\mathbf{A}}\|_F^2}{\|\Delta_{\mathbf{A}}\|_{\text{op}}^2} \leq C \eta^2 \frac{p' dr (\ln \tau)^2}{N} \right\}.$$

We set $\mathcal{A}_{r,p'} = \mathcal{A}_{r,p'}(\infty)$ and $\mathcal{G}(C)_{r,p'} = \mathcal{G}(C, \infty)$. Lastly, we drop the subscript r, p' when the statement is valid for all r, p' .

Remark B.4. For strictly stable systems with $\|M_{\mathbf{A}^*}\|_{\text{op}} < 1$ and $\|M_{\mathbf{A}_{p'}^*}\|_{\text{op}} < 1$, the factor η is controlled by Corollary B.10. However, for marginally stable systems or explosive systems, there is no a prior good upper bound on η , implying that the misspecification results only applies to strictly stable systems.

B.1 THEOREM STATEMENT

In this subsection, we state Theorem B.5 that generalizes the statements in Theorems 4.1 to 4.3. We give a proof that reduces Theorem B.5 to a uniform concentration result in Theorem B.6. The proof of Theorem B.6 is deferred to Appendix B.5. Lastly, Corollary B.7 gives a corollary that removes the constraints on the diameter of the search set.

Theorem B.5. Let Assumptions 3.1 and 3.4 hold. Furthermore, let Assumption 3.7 for $r < d$ and Assumption 3.9 for $p' < p$ hold. Let $\hat{\mathbf{A}}$ be the following estimator:

$$\hat{\mathbf{A}} = \underset{\mathbf{A} \in \mathcal{A}(D)}{\text{argmin}} \mathcal{L}(\mathbf{A}).$$

Then, for any small $\delta > 0$, there exist $C(\delta) = \mathcal{O}(\ln(1/\delta))$ such that

$$\mathbb{P}(\hat{\mathbf{A}} \in \mathcal{S}(C(\delta), D)) \geq 1 - \delta.$$

Proof. Let $\mathcal{E}_{\mathbf{A}}$ be the following event

$$\mathcal{E}_{\mathbf{A}} = \{\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E)\}.$$

By Corollary B.13, $\mathcal{G}(C, D) \subset \mathcal{S}(C, D)$ and thus, for any random choice of \mathbf{A} ,

$$\mathbb{P}(\{\mathbf{A} \in \mathcal{S}(C(\delta), D)\}) \geq \mathbb{P}(\{\mathbf{A} \in \mathcal{G}(C, D)\}) = 1 - \mathbb{P}(\{\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)\}).$$

For the choice of $\hat{\mathbf{A}}$, $\mathbb{P}(\mathcal{E}_{\hat{\mathbf{A}}}) = 1$ by Corollary B.15 and

$$\mathbb{P}(\{\hat{\mathbf{A}} \in \mathcal{S}(C(\delta), D)\}) \geq 1 - \mathbb{P}(\{\hat{\mathbf{A}} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)\} \mid \mathcal{E}_{\hat{\mathbf{A}}}).$$

By Bayes rule, we have

$$\begin{aligned} \mathbb{P}(\{\hat{\mathbf{A}} \in \mathcal{S}(C(\delta), D)\}) &\geq 1 - \mathbb{P}(\mathcal{E}_{\hat{\mathbf{A}}} \mid \{\hat{\mathbf{A}} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)\}) \mathbb{P}(\{\hat{\mathbf{A}} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)\}) \\ &\geq 1 - \mathbb{P}(\mathcal{E}_{\hat{\mathbf{A}}} \mid \{\hat{\mathbf{A}} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)\}). \end{aligned}$$

Then, the proof is complete by applying Theorem B.6 to the right-hand side. \square

Theorem B.6. *Let all the assumptions of Theorem B.5 hold. Then, for any small $\delta > 0$, there exist a constant $C(\delta) = \mathcal{O}(\ln(1/\delta))$ such that*

$$\mathbb{P}(\exists \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C(\delta), D) : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E)) \leq \delta. \quad (21)$$

Corollary B.7. *Let all the assumptions of Theorem B.5 hold and consider the following estimator:*

$$\hat{\mathbf{A}}_{\text{OLS}} = \underset{\mathbf{A} \in \mathcal{A}(\infty)}{\text{argmin}} \mathcal{L}(\mathbf{A}).$$

Then, for any small $\delta > 0$, there exist $C(\delta) = \mathcal{O}(\ln(1/\delta))$ such that

$$\mathbb{P}(\hat{\mathbf{A}}_{\text{OLS}} \in \mathcal{S}(C(\delta), D)) \geq 1 - \delta,$$

given that NT satisfies the following condition:

$$N(T - p') \geq C\eta^2 p'^2 dr (\ln \tau)^2 \ln \ln D. \quad (22)$$

Proof. Assume that D is sufficiently large such that

$$\mathcal{A}(D)^\circ \subset \mathcal{G}(C, D),$$

i.e., the interior of $\mathcal{A}(D)$ contains $\mathcal{G}(C, D)$. We have the following relation:

$$\mathcal{A}(\infty) = \{\mathbf{A}' = \alpha \mathbf{A} \mid \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D), \alpha \geq 1 \in \mathbb{R}\}.$$

Then, by Theorem B.6, we have

$$\begin{aligned} \mathbb{P}(\exists \mathbf{A} \in \mathcal{A}(\infty) \setminus \mathcal{G}(C(\delta), D) : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E)) \\ = \mathbb{P}(\exists \alpha \geq 1 \in \mathbb{R}, \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C(\delta), D) : \|\Delta_{\alpha \mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\alpha \mathbf{A}} L_{\star} E)) \\ = \mathbb{P}(\exists \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C(\delta), D) : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E)) \leq \delta, \end{aligned}$$

as $\forall \alpha \geq 1$, we have the following:

$$\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E) \implies \|\Delta_{\alpha \mathbf{A}} L_{\star} E\|_F^2 \leq 2\eta \text{Tr}(E^\top \Delta_{\alpha \mathbf{A}} L_{\star} E).$$

Thus, the result is complete by applying the same argument as in Theorem B.5 where $\mathcal{A}(D)$ is replaced by $\mathcal{A}(\infty)$.

We only need to provide a D such that $\mathcal{A}(D)^\circ \subset \mathcal{G}(C, D)$. It is easier to ensure the inclusion $\mathcal{A}(D)^\circ \subset \mathcal{S}(C, D)$. For any $\mathbf{A} \in \mathcal{S}(C, D)$, we have

$$\|\Delta_{\mathbf{A}}\|_{\text{op}}^2 \leq p' \|\mathbf{A}\|_{\text{op}}^2 \leq p' \|\mathbf{A}\|_F^2 \leq CD^2 \eta^2 \frac{p' dr (\ln \tau)^2}{N(T - p')},$$

from Lemma B.12. Therefore, we need to find a D such that

$$D^2 \geq CD^2 \eta^2 \frac{T}{T - p} \frac{p' dr (\ln \tau)^2}{N} \ln \ln D, \quad (23)$$

where we make the $\ln \ln D$ factor in C explicit. See Appendix B.5 for the details why this constant is needed. Lastly, Equation (23) is satisfied for large enough N and T that verifies the condition in Equation (22). \square

B.2 TECHNICAL LEMMAS

In this subsection, we present simple technical results on L_{\star} , $M_{\mathbf{A}}$ and $\Delta_{\mathbf{A}}$ that are used in the proof of Theorem B.5.

Lemma B.8. *L_{\star} and $M_{\mathbf{A}^{\star}}$ satisfy the following relations:*

$$L_{\star} = M_{\mathbf{A}^{\star}} L_{\star} + I, \quad M_{\mathbf{A}^{\star}} = (L_{\star} - I) L_{\star}^{-1}, \quad L_{\star} = (I - M_{\mathbf{A}^{\star}})^{-1}.$$

Proof. The first relation follows from a direct computation. For the second, note that L_\star is invertible since it is a lower triangular matrix with non-zero diagonals. Lastly,

$$\begin{aligned} L_\star &= I + M_{\mathbf{A}^\star} L_\star \\ &= I + M_{\mathbf{A}^\star} + M_{\mathbf{A}^\star}^2 L_\star \\ &= \dots \\ &= I + M_{\mathbf{A}^\star} + M_{\mathbf{A}^\star}^2 + \dots + M_{\mathbf{A}^\star}^{T-1} \\ &= (I - M_{\mathbf{A}^\star})^{-1}, \end{aligned}$$

where we have used the fact that $M_{\mathbf{A}^\star}^T = 0_{Td}$. \square

Lemma B.9. Assume that $\|M_{\mathbf{A}^\star}\|_{\text{op}} < 1$. Then, the operator norm and minimum singular value of L_\star are bounded as follows,

$$\frac{1}{1 + \|M_{\mathbf{A}^\star}\|_{\text{op}}} \leq \|L_\star\|_{\text{op}} \leq \frac{1}{1 - \|M_{\mathbf{A}^\star}\|_{\text{op}}}, \quad \frac{1}{2} \leq \sigma_{\min}(L_\star).$$

Proof. By Weyl's inequality for singular values on the identity $L_\star = M_{\mathbf{A}^\star} L_\star + I$ from Lemma B.8,

$$\begin{aligned} \|L_\star\|_{\text{op}} &\leq \|I\|_{\text{op}} + \|M_{\mathbf{A}^\star} L_\star\|_{\text{op}} \leq 1 + \|M_{\mathbf{A}^\star}\|_{\text{op}} \|L_\star\|_{\text{op}}, \\ \|L_\star\|_{\text{op}} &\geq \|I\|_{\text{op}} - \|M_{\mathbf{A}^\star} L_\star\|_{\text{op}} \geq 1 - \|M_{\mathbf{A}^\star}\|_{\text{op}} \|L_\star\|_{\text{op}}. \end{aligned}$$

This implies the desired inequalities for $\|L_\star\|_{\text{op}}$. For the lower bound on minimal singular value, use Lemma B.8,

$$\sigma_{\min}(L_\star) = \sigma_{\min}((I - M_{\mathbf{A}^\star})^{-1}) = \frac{1}{\|I - M_{\mathbf{A}^\star}\|_{\text{op}}} \geq \frac{1}{1 + \|M_{\mathbf{A}^\star}\|_{\text{op}}} \geq \frac{1}{2}.$$

Corollary B.10. Assume that $\|M_{\mathbf{A}^\star}\|_{\text{op}} < 1$ and $\|M_{\mathbf{A}_{p'}^\star}\|_{\text{op}} < 1$. Then, we have

$$\eta \leq \frac{2}{1 - \|M_{\mathbf{A}^\star}\|_{\text{op}}}.$$

Proof. Applying Lemma B.9,

$$\eta \leq \|M_{\mathbf{A}^\star} - M_{\mathbf{A}_{p'}^\star}\|_{\text{op}} \|L_\star\|_{\text{op}} \leq \frac{1}{1 - \|M_{\mathbf{A}^\star}\|_{\text{op}}} \left(\|M_{\mathbf{A}^\star}\|_{\text{op}} + \|M_{\mathbf{A}_{p'}^\star}\|_{\text{op}} \right) \leq \frac{2}{1 - \|M_{\mathbf{A}^\star}\|_{\text{op}}}.$$

Lemma B.11. Assume that $\|M_{\mathbf{A}^\star}\|_{\text{op}} \leq D$. Then, the operator norm and minimum singular value of L_\star are bounded as follows,

$$\|L_\star\|_{\text{op}} \leq \frac{D^T - 1}{D - 1}, \quad \frac{1}{D + 1} \leq \sigma_{\min}(L_\star).$$

Proof. By Weyl's inequality for singular values on the identity $L_\star = I + M_{\mathbf{A}^\star} + \dots + M_{\mathbf{A}^\star}^{T-1}$ from Lemma B.8,

$$\|L_\star\|_{\text{op}} \leq \|I\|_{\text{op}} + \sum_{t=1}^{T-1} \|M_{\mathbf{A}^\star}^t\|_{\text{op}} \leq \sum_{t=0}^{T-1} D^t \leq \frac{D^T - 1}{D - 1}.$$

For the lower bound on minimal singular value, use Lemma B.8,

$$\sigma_{\min}(L_\star) = \sigma_{\min}((I - M_{\mathbf{A}^\star})^{-1}) = \frac{1}{\|I - M_{\mathbf{A}^\star}\|_{\text{op}}} \geq \frac{1}{1 + \|M_{\mathbf{A}^\star}\|_{\text{op}}} \geq \frac{1}{D + 1}.$$

Lemma B.12. For any $\mathbf{A} \in \mathbb{R}^{d \times pd}$,

$$\begin{aligned} \|\mathbf{A}\|_{\text{op}} &\leq \|\mathbf{M}_{\mathbf{A}}\|_{\text{op}} \leq \sqrt{p'} \|\mathbf{A}\|_{\text{op}}, \\ \frac{1}{T} \|\Delta_{\mathbf{A}}\|_F^2 &\leq \|\mathbf{A} - \mathbf{A}_{p'}^*\|_F^2 \leq \frac{1}{T - p'} \|\Delta_{\mathbf{A}}\|_F^2. \end{aligned}$$

Proof. Let $u = (u_1, \dots, u_T) \in \mathbb{R}^{Td}$ be an arbitrary vector with $\|u\|_2^2 = 1$. Then, setting $u_{-a} = 0$ for any $a \geq 0$,

$$\begin{aligned} \|M_{\mathbf{A}}u\|_2^2 &= \sum_{i=1}^T \|(M_{\mathbf{A}}u)_i\|_2^2 = \sum_{i=1}^T \left\| \sum_{k=1}^p A_k u_{i-k} \right\|_2^2 \\ &= \sum_{i=1}^T \|\mathbf{A}_{:p'} u_{i-p':i-1}\|_2^2 \leq \sum_{i=1}^T \|\mathbf{A}_{:p'}\|_{\text{op}}^2 \|u_{i-p':i-1}\|_2^2 \\ &\leq \|\mathbf{A}\|_{\text{op}}^2 \sum_{i=1}^T p' \|u_i\|_2^2 = p' \|\mathbf{A}\|_{\text{op}}^2. \end{aligned}$$

The left-hand side of the first inequality follows by picking $u_{p'+1:T} = 0$ and $u_{1:p'}$ as the maximal singular vector of $\mathbf{A}_{:p'}$ with unit length. The second inequality follows by a simple computation. \square

Corollary B.13. For any $\mathbf{A} \in \mathcal{A}(D)$,

$$\|\mathbf{A} - \mathbf{A}_{p'}^*\|_F^2 \leq \frac{D^2}{T - p'} \frac{\|\Delta_{\mathbf{A}}\|_F^2}{\|\Delta_{\mathbf{A}}\|_{\text{op}}^2}.$$

Proof. By definition of $\mathcal{A}(D)$,

$$\frac{D^2}{T - p'} \frac{\|\Delta_{\mathbf{A}}\|_F^2}{\|\Delta_{\mathbf{A}}\|_{\text{op}}^2} \geq \frac{1}{T - p'} \|\Delta_{\mathbf{A}}\|_F^2,$$

and the result follows by Lemma B.12. \square

Proposition B.14. The empirical risk minimizer $\hat{\mathbf{A}}$, i.e.,

$$\hat{\mathbf{A}} \in \operatorname{argmin}_{\mathbf{A} \in \mathcal{A}(D)} \mathcal{L}(\mathbf{A}), \quad (24)$$

implies $\mathcal{L}(\hat{\mathbf{A}}) \leq \mathcal{L}(\mathbf{A}_{p'}^*)$, which can be rewritten as follows:

$$\|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|_F^2 \leq 2 \operatorname{Tr} \left(E^{\top} L_{\star}^{\top} \Delta_{\hat{\mathbf{A}}}^{\top} (I - M_{\mathbf{A}_{p'}^*}) L_{\star} E \right). \quad (25)$$

Proof. By Lemma B.8,

$$\begin{aligned} \mathcal{L}(\mathbf{A}) &= \|(M_{\mathbf{A}} - I) L_{\star} E\|_2^2 = \|(M_{\mathbf{A}} - M_{\mathbf{A}^*}) L_{\star} - I\|_F^2 \|E\|_F^2 \\ &= \left\| \left[(M_{\mathbf{A}} - M_{\mathbf{A}_{p'}^*}) L_{\star} + (M_{\mathbf{A}_{p'}^*} - M_{\mathbf{A}^*}) L_{\star} - I \right] E \right\|_F^2 \\ &= \left\| (M_{\mathbf{A}} - M_{\mathbf{A}_{p'}^*}) L_{\star} E \right\|_F^2 + \left\| \left[(M_{\mathbf{A}_{p'}^*} - M_{\mathbf{A}^*}) L_{\star} - I \right] E \right\|_F^2 \\ &\quad + 2 \operatorname{Tr} \left(E^{\top} L_{\star}^{\top} (M_{\mathbf{A}} - M_{\mathbf{A}_{p'}^*})^{\top} \left[(M_{\mathbf{A}_{p'}^*} - M_{\mathbf{A}^*}) L_{\star} - I \right] E \right) \\ &= \left\| (M_{\mathbf{A}} - M_{\mathbf{A}_{p'}^*}) L_{\star} E \right\|_F^2 + \left\| (M_{\mathbf{A}_{p'}^*} - I) L_{\star} E \right\|_F^2 + 2 \operatorname{Tr} \left(E^{\top} L_{\star}^{\top} (M_{\mathbf{A}} - M_{\mathbf{A}_{p'}^*})^{\top} (M_{\mathbf{A}_{p'}^*} - I) L_{\star} E \right). \end{aligned}$$

Then, for $\hat{\mathbf{A}}$ that satisfy $\mathcal{L}(\hat{\mathbf{A}}) \leq \mathcal{L}(\mathbf{A}_{p'}^*)$, we have

$$\mathcal{L}(\hat{\mathbf{A}}) - \mathcal{L}(\mathbf{A}_{p'}^*) = \left\| (M_{\hat{\mathbf{A}}} - M_{\mathbf{A}_{p'}^*}) L_{\star} E \right\|_F^2 + 2 \operatorname{Tr} \left(E^{\top} L_{\star}^{\top} (M_{\hat{\mathbf{A}}} - M_{\mathbf{A}_{p'}^*})^{\top} (M_{\mathbf{A}_{p'}^*} - I) L_{\star} E \right) \leq 0,$$

which implies the desired result. \square

Corollary B.15. *Observe that for $p' = p$, Proposition B.14 reads*

$$\|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|_2^2 \leq 2 \operatorname{Tr} (E^{\top} \Delta_{\hat{\mathbf{A}}} L_{\star} E) .$$

For $p' < p$, one can write the following relaxed condition for any $\hat{\mathbf{A}}$:

$$\begin{aligned} \|\Delta_{\hat{\mathbf{A}}} L_{\star} E\|_2^2 &\leq 2 \left\| \left(I - M_{\mathbf{A}^{\star}} \right) L_{\star} \right\|_{\text{op}} \operatorname{Tr} (E^{\top} \Delta_{\hat{\mathbf{A}}} L_{\star} E) \\ &= 2 \left\| I_{Td} + \left(M_{\mathbf{A}^{\star}} - M_{\mathbf{A}^{\star}'} \right) L_{\star} \right\|_{\text{op}} \operatorname{Tr} (E^{\top} \Delta_{\hat{\mathbf{A}}} L_{\star} E) \\ &\leq 2 \left(1 + \left\| \left(M_{\mathbf{A}^{\star}} - M_{\mathbf{A}^{\star}'} \right) L_{\star} \right\|_{\text{op}} \right) \operatorname{Tr} (E^{\top} \Delta_{\hat{\mathbf{A}}} L_{\star} E) \\ &= 2\eta \operatorname{Tr} (E^{\top} \Delta_{\hat{\mathbf{A}}} L_{\star} E) . \end{aligned}$$

B.3 LOWER AND UPPER ISOMETRIES

In Theorem B.16, we present a uniform lower bound on $\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2$. In order to establish this lower bound, we first start with a point-wise lower bounds in Lemma B.19 that relies on Hanson-Wright inequality for bounding the deviation of quadratic forms of sub-Gaussian vectors. Then, we use Lemmas B.21 and B.22 with a discretization argument in Theorem B.16 to establish uniform isometries. Finally, with Corollary B.17, we have a uniform control over the range of both $\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2$ and $\operatorname{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E)$.

Theorem B.16. *Let $\delta > 0$ be small and fixed. Then, there exist a constant $1 \leq C(\delta) = \ln(\frac{1}{\delta})$ such that the following holds uniformly for all $\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)$ and $C \geq C(\delta)$:*

$$\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \geq \frac{\sigma^2}{8} \sigma_{\min}(L_{\star})^2 N \|\Delta_{\mathbf{A}}\|_F^2 , \quad (26)$$

with probability at least $1 - \delta$.

Proof. By Lemmas B.19 and B.20, with probability at least $1 - \delta_1 - \delta_2$, the following holds:

$$\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \geq \sigma^2 (1 - c^2 \nu_1) \sigma_{\min}(L_{\star})^2 N \|\Delta_{\mathbf{A}}\|_F^2 , \quad (27)$$

for any arbitrary $\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)$ where

$$\delta_1 = \exp(-C_{HW} C \nu_1^2 p' d r \ln \tau) , \quad \delta_2 = \exp(-C_{HW} C \nu_2^2 p' d r \ln \tau) .$$

Let $\mathcal{B}(C, D)$ be the normalized $\mathcal{A}(D) \setminus \mathcal{G}(C, D)$,

$$\mathcal{B}(C, D) = \left\{ \frac{\mathbf{A}}{\|\mathbf{A}\|_F} \mid \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D) \right\} .$$

Then, since the conditions are homogeneous, Equation (27) holds for any $\mathbf{A} \in \mathcal{B}(C, D)$ with probability $1 - \delta_1 - \delta_2$.

Let $\mathcal{N}_{\epsilon}(D)$ be ϵ -net over the set $\mathcal{B}(C, D)$. Hence, with probability at least

$$1 - \delta_0 = 1 - |\mathcal{N}_{\epsilon}(D)|(\delta_1 + \delta_2) ,$$

the condition Equation (27) holds $\forall \mathbf{A} \in \mathcal{N}_{\epsilon}(D)$. Moreover, by Lemmas B.21 to B.23, we have

$$\|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \geq \frac{1}{2} \sigma^2 (1 - c^2 \nu_1) \sigma_{\min}(L_{\star})^2 N \|\Delta_{\mathbf{A}}\|_F^2 - \sigma^2 (1 + c^2 \nu_3) \epsilon^2 \|L_{\star}\|_{\text{op}}^2 p' d N T$$

$$\operatorname{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E) \leq \sigma^2 c^2 \nu_2 \|L_{\star}\|_{\text{op}} \sqrt{C p' d r N \ln \tau} \|\Delta_{\mathbf{A}}\|_F + \sigma^2 (1 + c^2 \nu_3) \epsilon \|L_{\star}\|_{\text{op}} \sqrt{p' d N T} ,$$

$\forall \mathbf{A} \in \mathcal{B}(C, D)$ with probability at least $1 - \delta_0 - \delta_3$ where

$$\delta_3 = \exp(-C_{HW} \nu_3 d N T) , .$$

Recall that

$$\|\Delta_{\mathbf{A}}\|_F^2 \geq (T - p) \|\mathbf{A}\|_F^2 = T - p ,$$

for any $\mathbf{A} \in \mathcal{B}(C, D)$ due to the normalization.

Setting $\nu_1 = \frac{1}{4c^2}$, $\nu_2 = \frac{1}{4c^2}$ and ϵ such that

$$\epsilon = \frac{1}{2(1+c^2\nu_3)} \cdot \sqrt{\frac{T-p'}{T}} \min \left\{ \frac{1}{\sigma_{\text{cond}}(L_\star)} \sqrt{\frac{r}{dNT}}, \sqrt{\frac{1}{p'd}} \right\},$$

and recalling that $C \geq 1$ and $\ln \tau \geq 1$,

$$\begin{aligned} \frac{1}{2} \sigma^2 (1 - c^2 \nu_1) \sigma_{\min}(L_\star)^2 N \|\Delta_{\mathbf{A}}\|_F^2 - \sigma^2 (1 + c^2 \nu_3) \epsilon^2 \|L_\star\|_{\text{op}}^2 p' d N T &\geq \frac{\sigma^2}{8} \sigma_{\min}(L_\star)^2 N \|\Delta_{\mathbf{A}}\|_F^2, \\ \sigma^2 c^2 \nu_2 \|L_\star\|_{\text{op}} \sqrt{C p' d r N \ln \tau} \|\Delta_{\mathbf{A}}\|_F + \sigma^2 (1 + c^2 \nu_3) \epsilon \|L_\star\|_{\text{op}} \sqrt{p' d N T} &\leq \frac{\sigma^2}{2} \|L_\star\|_{\text{op}} \sqrt{C p' d r N \ln \tau} \|\Delta_{\mathbf{A}}\|_F, \end{aligned}$$

$\forall \mathbf{A} \in \mathcal{B}(C, D)$ with probability $1 - \delta_0 - \delta_3$.

By homogeneity, this implies that $\forall \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)$, with probability at least $1 - \delta_0 - \delta_3$,

$$\begin{aligned} \|\Delta_{\mathbf{A}} L_\star E\|_F^2 &\geq \frac{\sigma^2}{8} \sigma_{\min}(L_\star)^2 N \|\Delta_{\mathbf{A}}\|_F^2, \\ \text{Tr}(E^\top \Delta_{\mathbf{A}} L_\star E) &\leq \frac{\sigma^2}{2} \|L_\star\|_{\text{op}} \sqrt{C p' d r N \ln \tau} \|\Delta_{\mathbf{A}}\|_F. \end{aligned}$$

Lastly, we can ensure that $\delta_3 < \delta/2$ with the choice of

$$\nu_3(\delta) = \frac{\ln \frac{1}{\delta/2}}{C_{HW}} > \frac{\ln \frac{1}{\delta/2}}{C_{HW} d N T}.$$

Moreover, the ϵ -net size can be bounded as follows:

$$|\mathcal{N}_\epsilon(D)| \leq \left(\frac{9}{\epsilon}\right)^{(p'd+d+1)r} \leq \exp\left(3p'dr \ln \frac{9}{\epsilon}\right) \leq \exp\left(9p'dr \ln \frac{1}{\epsilon}\right).$$

For more details on ϵ -nets on low-rank matrices, see Candès & Plan (2011, Lemma 3.1). Then,

$$\delta_0 = |\mathcal{N}_\epsilon(D)| (\delta_1 + \delta_2) \leq \exp\left(\left(9 - \frac{1}{16c^4} C_{HW} C\right) p' d r \ln \tau + 9p'dr \ln \left(1 + \frac{c^2}{C_{HW}} \ln \frac{1}{\delta/2}\right)\right),$$

where we use that $\frac{1}{\epsilon} < \tau \left(1 + \frac{c^2}{C_{HW}} \ln \frac{1}{\delta/2}\right)$. Thus, $\delta_0 < \delta/2$ can be made with the choice of

$$\begin{aligned} C(\delta) &= \frac{16c^4}{C_{HW}} \left(9 + \ln \frac{1}{\delta/2} + 9 \ln \left(1 + \frac{c^2}{C_{HW}} \ln \frac{1}{\delta/2}\right)\right) \\ &> \frac{16c^4}{C_{HW}} \left(9 + \frac{\ln \frac{1}{\delta/2}}{p' d r \ln \tau} + 9 \frac{\ln \left(1 + \frac{c^2}{C_{HW}} \ln \frac{1}{\delta/2}\right)}{\ln \tau}\right). \end{aligned}$$

Here, $C(\delta)$ is a constant that is independent of p, p', d, r, N, T such that $C(\delta) = \mathcal{O}(\ln(1/\delta))$. \square

Corollary B.17. *For any small $\delta > 0$, there exists a constant $1 \leq C(\delta) = \mathcal{O}(\ln(1/\delta))$ such that the following holds uniformly for all $\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)$ and $C \geq C(\delta)$:*

$$\begin{aligned} \inf_{\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)} \|\Delta_{\mathbf{A}} L_\star E\|_F^2 &\geq \frac{\sigma^2}{8} \sigma_{\min}(L_\star)^2 C \eta^2 D^2 p' d r (\ln \tau)^2, \\ \sup_{\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)} \text{Tr}(E^\top \Delta_{\mathbf{A}} L_\star E) &\leq \frac{\sigma^2}{2} \|L_\star\|_{\text{op}} C \eta D p' d r (\ln \tau)^{3/2}, \end{aligned} \tag{28}$$

Proof. Plug in the results from Theorem B.16 and use the definition of the set $\mathcal{A}(D) \setminus \mathcal{G}(C, D)$. \square

Definition B.18. For applying Theorem A.1 in our setup, consider the following objects:

$$\begin{aligned}\tilde{E} &= (\xi^{(1)\top}, \dots, \xi^{(N)\top})^\top \in \mathbb{R}^{NTd}, \\ \tilde{\Delta}_A &= \text{diag}(\Delta_A) \in \mathbb{R}^{NTd} \times \mathbb{R}^{NTd},\end{aligned}$$

where $\text{diag}(P)$ puts P in the diagonal blocks of a larger diagonal matrix.

Lemma B.19. For any $A \in \mathcal{A} \setminus \mathcal{G}(C)$ and $\nu \in (0, 1)$, with probability at least

$$1 - \exp(-C_{HW} C \nu^2 p' dr \ln \tau),$$

we have the following

$$\|\Delta_A L_\star E\|_F^2 \geq \sigma^2 (1 - c^2 \nu) \sigma_{\min}(L_\star)^2 N \|\Delta_A\|_F^2.$$

Proof. First, observe that

$$\|\Delta_A L_\star E\|_F^2 \geq \sigma_{\min}(L_\star)^2 \|\Delta_A E\|_F^2.$$

Applying Theorem A.1 with $P = \tilde{\Delta}_A^\top \tilde{\Delta}_A$ and $r = c^2 \sigma^2 \nu \|\tilde{\Delta}_A\|_F^2$,

$$\begin{aligned}\mathbb{P}(\tilde{E}^\top \tilde{\Delta}_A^\top \tilde{\Delta}_A \tilde{E} - \mathbb{E}[\tilde{E}^\top \tilde{\Delta}_A^\top \tilde{\Delta}_A \tilde{E}] \geq c^2 \sigma^2 \nu \|\tilde{\Delta}_A\|_F^2) \\ \leq \exp\left(-C_{HW} \min\left\{\nu^2 \frac{\|\tilde{\Delta}_A\|_F^4}{\|\tilde{\Delta}_A^\top \tilde{\Delta}_A\|_F^2}, \nu \frac{\|\tilde{\Delta}_A\|_F^2}{\|\tilde{\Delta}_A^\top \tilde{\Delta}_A\|_{\text{op}}}\right\}\right).\end{aligned}$$

Observe that $\tilde{E}^\top \tilde{\Delta}_A^\top \tilde{\Delta}_A \tilde{E} = \text{Tr}(E^\top \Delta_A^\top \Delta_A E) = \|\Delta_A E\|_F^2$ and

$$\mathbb{E}[\tilde{E}^\top \tilde{\Delta}_A^\top \tilde{\Delta}_A \tilde{E}] = \mathbb{E}[\text{Tr}(\tilde{E} \tilde{E}^\top \tilde{\Delta}_A^\top \tilde{\Delta}_A)] = \sigma^2 \|\tilde{\Delta}_A\|_F^2 = \sigma^2 N \|\Delta_A\|_F^2.$$

Furthermore, $\|\tilde{\Delta}_A\|_F^4 = N^2 \|\Delta_A\|_F^4$, $\|\tilde{\Delta}_A^\top \tilde{\Delta}_A\|_F^2 = N \|\Delta_A^\top \Delta_A\|_F^2$ and $\|\tilde{\Delta}_A^\top \tilde{\Delta}_A\|_{\text{op}} = \|\Delta_A^\top \Delta_A\|_{\text{op}} = \|\Delta_A\|_{\text{op}}^2$. Plugging these into the bound,

$$\mathbb{P}(\|\Delta_A E\|_F^2 - \sigma^2 N \|\Delta_A\|_F^2 \geq c^2 \sigma^2 \nu N \|\Delta_A\|_F^2) \leq \exp\left(-C_{HW} \min\left\{\nu^2 N \frac{\|\Delta_A\|_F^4}{\|\Delta_A^\top \Delta_A\|_F^2}, \nu N \frac{\|\Delta_A\|_F^2}{\|\Delta_A\|_{\text{op}}^2}\right\}\right).$$

Then, using $\|\Delta_A^\top \Delta_A\|_F^2 \leq \|\Delta_A\|_F^2 \|\Delta_A\|_{\text{op}}^2$ and $\nu < 1$,

$$\mathbb{P}(\|\Delta_A E\|_F^2 \geq \sigma^2 (1 - c^2 \nu) N \|\Delta_A\|_F^2) \geq 1 - \exp\left(-C_{HW} \nu^2 N \frac{\|\Delta_A\|_F^2}{\|\Delta_A\|_{\text{op}}^2}\right).$$

The result follows from the definition of set $\mathcal{A} \setminus \mathcal{G}(C)$. \square

Lemma B.20. For any $A \in \mathcal{A} \setminus \mathcal{G}(C)$ and $\nu \in (0, 1)$, with probability at least

$$1 - \exp(-C_{HW} C \nu^2 p' dr \ln \tau),$$

we have the following

$$\text{Tr}(E^\top \Delta_A L_\star E) \leq c^2 \sigma^2 \nu \|L_\star\|_{\text{op}} \sqrt{C p' dr N \ln \tau} \|\Delta_A\|_F.$$

Proof. First, by the properties of trace and Frobenius norm, we have

$$\text{Tr}(E^\top \Delta_A L_\star E) \leq \|L_\star\|_{\text{op}} \text{Tr}(E^\top \Delta_A E).$$

Applying Theorem A.1 with $P = \tilde{\Delta}_A$ and $r = c^2 \sigma^2 \nu \sqrt{C p' dr \ln \tau} \|\tilde{\Delta}_A\|_F$,

$$\begin{aligned}\mathbb{P}(\tilde{E}^\top \tilde{\Delta}_A \tilde{E} - \mathbb{E}[\tilde{E}^\top \tilde{\Delta}_A \tilde{E}] \geq c^2 \sigma^2 \nu \sqrt{C p' dr \ln \tau} \|\tilde{\Delta}_A\|_F) \\ \leq \exp\left(-C_{HW} \min\left\{\nu^2 C p' dr \ln \tau, \nu \sqrt{C p' dr \ln \tau} \frac{\|\tilde{\Delta}_A\|_F}{\|\tilde{\Delta}_A\|_{\text{op}}}\right\}\right).\end{aligned}$$

Noting $\mathbb{E}[\tilde{E}^\top \tilde{\Delta}_A \tilde{E}] = 0$ and rewriting,

$$\begin{aligned}\mathbb{P}(\text{Tr}(E^\top \Delta_A E) \leq c^2 \sigma^2 \nu \sqrt{C p' dr N \ln \tau} \|\Delta_A\|_F) \geq \\ 1 - \exp\left(-C_{HW} \min\left\{\nu^2 C p' dr \ln \tau, \nu \sqrt{C p' dr N \ln \tau} \frac{\|\Delta_A\|_F}{\|\Delta_A\|_{\text{op}}}\right\}\right).\end{aligned}$$

The result follows from the definition of set $\mathcal{A} \setminus \mathcal{G}(C)$. \square

Lemma B.21. For any $\nu \geq 1$, with probability at least

$$1 - \exp(-C_{HW}\nu dNT),$$

we have the following

$$\|E\|_F^2 - \sigma^2 dNT \leq c^2 \sigma^2 \nu dNT.$$

Proof. Applying Theorem A.1 with $P = I_{dTN}$, $r = c^2 \sigma^2 \nu dNT$,

$$\mathbb{P}(\tilde{E}^\top \tilde{E} - \mathbb{E}[\tilde{E}^\top \tilde{E}] \geq c^2 \sigma^2 \nu dNT) \leq \exp(-C_{HW}\nu dNT).$$

The result follows after a simple computation. \square

Lemma B.22. For any $\mathbf{A}_1, \mathbf{A}_2 \in \mathcal{A}$,

$$\|\Delta_{\mathbf{A}_2} L_\star E\|_F^2 \geq \frac{1}{2} \|\Delta_{\mathbf{A}_1} L_\star E\|_F^2 - p' \|\mathbf{A}_1 - \mathbf{A}_2\|_{\text{op}}^2 \|L_\star\|_{\text{op}}^2 \|E\|_F^2. \quad (29)$$

Proof. By the properties of Frobenius norm and Lemma B.12,

$$\begin{aligned} \|\Delta_{\mathbf{A}_1} L_\star E\|_F^2 &= \|(\Delta_{\mathbf{A}_1} - \Delta_{\mathbf{A}_2} + \Delta_{\mathbf{A}_2}) L_\star E\|_F^2 \\ &\leq 2\|\Delta_{\mathbf{A}_2} L_\star E\|_F^2 + 2\|(\Delta_{\mathbf{A}_1} - \Delta_{\mathbf{A}_2}) L_\star E\|_F^2 \\ &\leq 2\|\Delta_{\mathbf{A}_2} L_\star E\|_F^2 + 2\|\Delta_{\mathbf{A}_1} - \Delta_{\mathbf{A}_2}\|_{\text{op}}^2 \|L_\star\|_{\text{op}}^2 \|E\|_F^2 \\ &\leq 2\|\Delta_{\mathbf{A}_2} L_\star E\|_F^2 + 2p' \|\mathbf{A}_1 - \mathbf{A}_2\|_{\text{op}}^2 \|L_\star\|_{\text{op}}^2 \|E\|_F^2. \end{aligned} \quad (30)$$

The results readily follows by reordering terms. \square

Lemma B.23. For any $\mathbf{A}_1, \mathbf{A}_2 \in \mathcal{A}$,

$$\text{Tr}(E^\top \Delta_{\mathbf{A}_2} L_\star E) \leq \text{Tr}(E^\top \Delta_{\mathbf{A}_1} L_\star E) + \sqrt{p'} \|\mathbf{A}_1 - \mathbf{A}_2\|_{\text{op}} \|L_\star\|_{\text{op}} \|E\|_F^2, \quad (31)$$

Proof. By the properties of trace and Lemma B.12,

$$\begin{aligned} \text{Tr}(E^\top \Delta_{\mathbf{A}_2} L_\star E) &= \text{Tr}(E^\top (\Delta_{\mathbf{A}_2} - \Delta_{\mathbf{A}_1} + \Delta_{\mathbf{A}_1}) L_\star E) \\ &= \text{Tr}(E^\top \Delta_{\mathbf{A}_1} L_\star E) + \text{Tr}(E^\top (\Delta_{\mathbf{A}_2} - \Delta_{\mathbf{A}_1}) L_\star E) \\ &\leq \text{Tr}(E^\top \Delta_{\mathbf{A}_1} L_\star E) + \|\Delta_{\mathbf{A}_2} - \Delta_{\mathbf{A}_1}\|_{\text{op}} \|L_\star\|_{\text{op}} \|E\|_F^2 \\ &\leq \text{Tr}(E^\top \Delta_{\mathbf{A}_1} L_\star E) + \sqrt{p'} \|\mathbf{A}_1 - \mathbf{A}_2\|_{\text{op}} \|L_\star\|_{\text{op}} \|E\|_F^2. \end{aligned} \quad (32)$$

\square

B.4 CONCENTRATION INEQUALITIES

In Remark B.24, we show that the quantities of interest that show up in Equation (21) are related to a martingale series and its predictable quadratic variation. This allow us to use Lemma A.5 in Theorem B.26 to quantify the probability of the event in Equation (21) for finite sets of \mathbf{A} .

Remark B.24. Fix $\mathbf{A} \in \mathcal{A}(D)$. Consider the martingale differences sequences

$$d_{t,i}^{(n)} = \left((\mathbf{A} - \mathbf{A}_{p'}^\star) x_t^{(n)} \right)_i \left(\xi_t^{(n)} \right)_i / \sigma^2,$$

where the series is first ordered in i , then in t , and finally in n . Let Y be the sum of the martingale differences, i.e.,

$$Y = \sum_{i,t,n} d_{i,t}^{(n)}.$$

Let $W_{\mathbf{A}}^R$ be the quadratic variation of the series plus an error term as in Theorem A.4, i.e.,

$$W_{\mathbf{A}}^R = \sum_{i,t,n} \mathbb{E}_{(\xi_t^{(n)})_i} \left[\left(d_{i,t}^{(n)} \right)^2 \right] + \sum_{n,t,i} \mathbb{1}_{d_{i,t}^{(n)} > R} \left(d_{i,t}^{(n)} \right)^2,$$

Then, we have the following computations:

$$Y_{\mathbf{A}} = \frac{1}{\sigma^2} \sum_{n,t} \langle (\mathbf{A} - \mathbf{A}_{p'}^*) X_t^{(n)}, \xi_t^{(n)} \rangle = \frac{1}{\sigma^2} \text{Tr} (E^\top \Delta_{\mathbf{A}} L_{\star} E) ,$$

$$W_{\mathbf{A}} := W_{\mathbf{A}}^0 = \frac{1}{\sigma^2} \sum_{n,t} \left\| (\mathbf{A} - \mathbf{A}_{p'}^*) X_t^{(n)} \right\|^2 = \frac{1}{\sigma^2} \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 .$$

Proposition B.25. *Let $R > 0$ be a constant and $\delta > 0$ be small. Then, there exist a constant $C'(\delta) = \mathcal{O}(\ln \frac{1}{\delta})$ such that*

$$\forall \mathbf{A} \in \mathcal{A}(D), \quad C'(\delta) \ln dTN \cdot W_{\mathbf{A}} \geq W_{\mathbf{A}}^R, \quad (33)$$

with probability $1 - \delta$.

Proof. By Corollary A.7, there exist a constant $c'(\delta) = \mathcal{O}(\sqrt{\ln \frac{1}{\delta}})$,

$$\sup_{t,n} \|\xi_t^{(n)}\|_{\infty} \leq c'(\delta) \sigma \sqrt{2 \ln dTN} .$$

Therefore, for any $\mathbf{A} \in \mathcal{A}(D)$, we have

$$\begin{aligned} W_{\mathbf{A}}^R &= W_{\mathbf{A}} + \sum_{n,t,i} \mathbb{1}_{d_{t,i}^{(n)} > R} \left(d_{t,i}^{(n)} \right)^2 \\ &\leq W_{\mathbf{A}} + 2c'(\delta)^2 \ln dTN \sum_{n,t,i} \mathbb{1}_{d_{t,i}^{(n)} > R} \left((\mathbf{A} - \mathbf{A}_{p'}^*) x_t^{(n)} \right)_i^2 \\ &\leq W_{\mathbf{A}} + 2c'(\delta)^2 \ln dTN \sum_{n,t,i} \left((\mathbf{A} - \mathbf{A}_{p'}^*) x_t^{(n)} \right)_i^2 \\ &\leq W_{\mathbf{A}} + 2c'(\delta)^2 \ln dTN \cdot W_{\mathbf{A}} . \end{aligned}$$

Then, by rearranging terms, we have

$$(1 + 2c'(\delta)^2 \ln dTN) W_{\mathbf{A}} \geq W_{\mathbf{A}}^R .$$

□

Theorem B.26. *Let $S \subseteq \mathcal{A}(D)$ be a set and let $\mathcal{E}(S)$ be the following event*

$$\mathcal{E}(S) = \left\{ \inf_{\mathbf{A} \in S} W_{\mathbf{A}} \geq \alpha_L \right\} \cap \left\{ \sup_{\mathbf{A} \in S} Y_{\mathbf{A}} \leq \alpha_U \right\} ,$$

where $\alpha_L, \alpha_U > 0$ are two constants. Then, for any $\gamma > 0, R > 1$ and $\delta > 0$ small, there exist a constant $C'(\delta) = \mathcal{O}(\ln \frac{1}{\delta})$ such that

$$\mathbb{P}(\exists \mathbf{A} \in S : W_{\mathbf{A}} \leq \gamma Y_{\mathbf{A}}) \leq |S| \exp \left(-\frac{\alpha_L}{2C'(\delta)R\gamma \ln dTN} + \ln \left(\ln \left(\frac{\alpha_U}{R\alpha_L} \right) + 1 \right) \right) + \mathbb{P}(\mathcal{E}(S)^C) + \delta .$$

Proof. The statement is trivial for $\alpha_L \leq \alpha_U$ so we consider the case $\alpha_L < \alpha_U$.

For any \mathbf{A} , let $\mathcal{E}_{\mathbf{A}}$ be the following event:

$$\mathcal{E}_{\mathbf{A}} = \{W_{\mathbf{A}}^R \geq \alpha_L\} \cap \{\gamma Y_{\mathbf{A}} \leq \alpha_U\} .$$

Then, by union bound, we have

$$\begin{aligned} \mathbb{P}(\{\exists \mathbf{A} \in S : W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\}) &= \mathbb{P}(\{\exists \mathbf{A} \in S : W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}(S)) + \mathbb{P}(\{\exists \mathbf{A} \in S : W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}(S)^C) \\ &\leq \mathbb{P}(\{\exists \mathbf{A} \in S : W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}(S)) + \mathbb{P}(\mathcal{E}(S)^C) \\ &\leq \sum_{\mathbf{A} \in S} \mathbb{P}(\{W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}(S)) + \mathbb{P}(\mathcal{E}(S)^C) \\ &\leq \sum_{\mathbf{A} \in S} \mathbb{P}(\{W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}_{\mathbf{A}}) + \mathbb{P}(\mathcal{E}(S)^C) . \end{aligned}$$

For any $\mathbf{A} \in \mathcal{S}$,

$$\mathbb{P}(W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}} \cap \mathcal{E}_{\mathbf{A}}) \leq \exp\left(-\frac{\alpha_L}{2R\gamma} + \ln \ln\left(\frac{\alpha_U}{R\alpha_L}\right)\right).$$

by Lemma A.5 which implies that

$$\mathbb{P}(\{\exists \mathbf{A} \in \mathcal{S} : W_{\mathbf{A}}^R \leq \gamma Y_{\mathbf{A}}\}) \leq |\mathcal{S}| \exp\left(-\frac{\alpha_L}{2R\gamma} + \ln\left(\ln\left(\frac{\alpha_U}{R\alpha_L}\right) + 1\right)\right) + \mathbb{P}(\mathcal{E}(\mathcal{S})^C).$$

Finally, let $C'(\delta)$ be the constant from Proposition B.25 and \mathcal{E}_{δ} be the event in Equation (33). Then,

$$\begin{aligned} \mathbb{P}(\{\exists \mathbf{A} \in \mathcal{S} : W_{\mathbf{A}} \leq \gamma Y_{\mathbf{A}}\}) &\leq \mathbb{P}(\{\exists \mathbf{A} \in \mathcal{S} : W_{\mathbf{A}} \leq \gamma Y_{\mathbf{A}}\} \cap \mathcal{E}_{\delta}) + \mathbb{P}(\mathcal{E}_{\delta}^C) \\ &\leq \mathbb{P}(\{\exists \mathbf{A} \in \mathcal{S} : W_{\mathbf{A}}^R \leq C'(\delta)\gamma \ln dTN \cdot Y_{\mathbf{A}}\} \cap \mathcal{E}_{\delta}) + \delta \\ &\leq \mathbb{P}(\{\exists \mathbf{A} \in \mathcal{S} : W_{\mathbf{A}}^R \leq C'(\delta)\gamma \ln dTN \cdot Y_{\mathbf{A}}\}) + \delta. \end{aligned}$$

□

B.5 PROOF OF THEOREM B.6

By Corollary B.17, there exist a constant $C(\delta/4)$ such that for all $C \geq C(\delta/4)$

$$\begin{aligned} \inf_{\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C,D)} \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 &\geq \frac{\sigma^2}{8} \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2, \\ \sup_{\mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C,D)} \text{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E) &\leq \frac{\sigma^2}{2} \|L_{\star}\|_{\text{op}} C \eta D p' dr (\ln \tau)^2, \end{aligned} \quad (34)$$

with probability $1 - \delta/4$.

Let \mathcal{S} be an ϵ -net over $\mathcal{A}(D) \setminus \mathcal{G}(C,D)$. Then, by Theorem B.26, there exist a constant $C'(\delta/4)$ such that

$$\begin{aligned} \mathbb{P}(\exists \mathbf{A} \in \mathcal{S} : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq \gamma \text{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E)) &\leq \\ |\mathcal{S}| \exp\left(-\frac{\alpha_L}{2C'(\delta/4)R\gamma \ln dTN} + \ln\left(\ln\left(\frac{\alpha_U}{R\alpha_L}\right) + 1\right)\right) &+ \frac{2\delta}{4}, \end{aligned}$$

where α_L and α_U are the lower and upper bounds in Equation (34) scaled with σ^2 :

$$\begin{aligned} \alpha_L &= \frac{1}{8} \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2, \\ \alpha_U &= \frac{1}{2} \|L_{\star}\|_{\text{op}} C \eta D p' dr (\ln \tau)^{3/2}. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \mathbb{P}(\exists \mathbf{A} \in \mathcal{S} : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq \gamma \text{Tr}(E^{\top} \Delta_{\mathbf{A}} L_{\star} E)) &\leq \\ |\mathcal{S}| \exp\left(-\frac{\frac{1}{8} \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2}{2C'(\delta/4)R\gamma \ln dTN} + \ln\left(\ln\left(\frac{4\|L_{\star}\|_{\text{op}}}{R\sigma_{\min}(L_{\star})^2 \eta D}\right) + 1\right)\right) &+ \frac{\delta}{2}. \end{aligned} \quad (35)$$

Recall that Lemmas B.21 to B.23 imply

$$\begin{aligned} \|\Delta_{\mathbf{A}_2} L_{\star} E\|_F^2 &\geq \frac{1}{2} \|\Delta_{\mathbf{A}_1} L_{\star} E\|_F^2 - \sigma^2(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}}) p' dNT \epsilon^2 \|L_{\star}\|_{\text{op}}^2, \\ \text{Tr}(E^{\top} \Delta_{\mathbf{A}_2} L_{\star} E) &\leq \text{Tr}(E^{\top} \Delta_{\mathbf{A}_1} L_{\star} E) + \sigma^2(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}}) \sqrt{p'} dNT \epsilon \|L_{\star}\|_{\text{op}}, \end{aligned}$$

with probability at least $1 - \delta/4$. We set ϵ as follows:

$$\epsilon = \min \left\{ \frac{\frac{1}{2} \|L_{\star}\|_{\text{op}} C \eta D p' dr (\ln \tau)^{3/2}}{(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}}) \sqrt{p'} dNT \|L_{\star}\|_{\text{op}}}, \frac{1}{2} \sqrt{\frac{\frac{1}{8} \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2}{(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}}) p' dNT \|L_{\star}\|_{\text{op}}^2}} \right\}.$$

In particular, ϵ is small such that for any $\mathbf{A}_1 \in \mathcal{A}(D) \setminus \mathcal{G}(C, D)$,

$$\exists \mathbf{A}_2 \in \mathcal{S} : \quad \|\Delta_{\mathbf{A}_1} L_{\star} E\|_F^2 \leq \frac{1}{4} \|\Delta_{\mathbf{A}_2} L_{\star} E\|_F^2, \quad \text{Tr}(E^\top \Delta_{\mathbf{A}_1} L_{\star} E) \leq 2 \text{Tr}(E^\top \Delta_{\mathbf{A}_2} L_{\star} E).$$

Then, Equation (35) implies that

$$\begin{aligned} & \mathbb{P}(\exists \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D) : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq \gamma \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E)) \leq \\ & |\mathcal{S}| \exp \left(-\frac{\frac{1}{8} \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2}{16 C'(\delta/4) R \gamma \ln dTN} + \ln \left(\ln \left(\frac{4 \|L_{\star}\|_{\text{op}}}{R \sigma_{\min}(L_{\star})^2 \eta D} \right) + 1 \right) \right) + \frac{3\delta}{4}, \end{aligned} \quad (36)$$

by the following computation:

$$\begin{aligned} & \mathbb{P}(\exists \mathbf{A} \in \mathcal{A}(D) \setminus \mathcal{G}(C, D) : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq \gamma \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E)) \\ & \leq \mathbb{P} \left(\exists \mathbf{A} \in \mathcal{S} : \|\Delta_{\mathbf{A}} L_{\star} E\|_F^2 \leq \frac{\gamma}{8} \text{Tr}(E^\top \Delta_{\mathbf{A}} L_{\star} E) \right). \end{aligned}$$

We now have to show that the right-hand side of Equation (36) for $\gamma = 2$ is upper bounded by δ . That is, we need to prove

$$\begin{aligned} & \sigma_{\min}(L_{\star})^2 C \eta^2 D^2 p' dr (\ln \tau)^2 \\ & \geq 256 C'(\delta/4) R \ln dTN \left(\ln |\mathcal{S}| + \ln \left(\ln \left(\frac{4 \|L_{\star}\|_{\text{op}}}{R \sigma_{\min}(L_{\star})^2 \eta D} \right) + 1 \right) + \delta/4 \right). \end{aligned}$$

In order to simplify the expressions, we set $R = 4$, plug in lower bound for $\sigma_{\min}(L_{\star})$ and upper bound for $\|L_{\star}\|_{\text{op}}$ from Lemma B.11, lower bound η with 1, upper bound $\ln dTN$ with $2 \ln \tau$ and derive the following looser condition:

$$C \frac{D^2}{(D+1)^2} p' dr \ln \tau \geq 2048 C'(\delta/4) \left(\ln |\mathcal{S}| + \ln \left(\ln \left(\frac{(D^T - 1)(D+1)^2}{(D-1)D} \right) + 1 \right) + \delta/4 \right). \quad (37)$$

The cardinality of \mathcal{S} is smaller than the cardinality of an ϵ -net \mathcal{S}' covering all of $\mathcal{A}(D) \setminus \mathcal{G}(C, D) \subset \mathcal{A}(D)$. Therefore, we have the following upper bound on $\ln |\mathcal{S}|$:

$$\ln |\mathcal{S}| \leq 9 p' dr \ln \frac{D}{\epsilon} \leq 9 p' dr \ln \left(4\sqrt{2} \left(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}} \right) NT \sigma_{\text{cond}}(L_{\star}) \right), \quad (38)$$

where we grossly upper bound $\frac{1}{\epsilon}$ similar to Equation (37). Then, Equation (37) is satisfied if:

$$\begin{aligned} & C p' dr \ln \tau \geq 2048 C'(\delta/4) \frac{(D+1)^2}{D^2} \left(18 p' dr \ln \tau + 9 \ln \left(4\sqrt{2} \left(1 + c^2 \frac{\ln \frac{4}{\delta}}{C_{HW}} \right) \right) \right) \\ & + 2048 C'(\delta/4) \frac{(D+1)^2}{D^2} \left(\ln \left(\ln \left(\frac{(D^T - 1)(D+1)^2}{(D-1)D} \right) + 1 \right) + \delta/4 \right), \end{aligned} \quad (39)$$

where we plug in the upper bound in Equation (38) and then bound $NT \sigma_{\text{cond}}(L_{\star})$ with τ^2 . Now, note that there exist a constant $C''(D)$ such that

$$\ln \left(\ln \left(\frac{(D^T - 1)(D+1)^2}{(D-1)D} \right) + 1 \right) \leq C''(D) \ln T \leq C''(D) \ln \tau.$$

Therefore, Equation (39) is satisfied for a constant $C = \mathcal{O}(\ln \frac{1}{\delta})$.

Finally, observe that the dependencies of C on c are due to the noise concentration inequalities and are the results of applications of Lemma B.21. Therefore, one can completely remove this dependency by fixing $\delta > \exp(-\nu dNT) + \delta'$ where $\nu = \nu(c)$ is a well-chosen constant and $\delta' > 0$ is small and arbitrary. The $\ln \ln D$ dependency, however, is due to the martingale concentration inequalities in Theorem B.26. Therefore, the dependency of C on D can only be removed by improving the concentration arguments.

C EXPERIMENTS

All experiments in this section are implemented with Python 3 (Van Rossum & Drake, 2009) under PSF license and PyTorch (Paszke et al., 2019) under BSD-3-Clause license. In addition, we use NumPy (Harris et al., 2020) under BSD license.

For all the experiments, \mathbf{A}^* is generated as follows. First, p orthogonal matrices of shape $d \times d$ are sampled. These are then scaled down by $\alpha \cdot p$ where α is arbitrarily set to 0.5. In cases where \mathbf{A} needs to be initialized, we use the same recipe for the student model with p' instead of p and set $\alpha = 1$. For experiments with low-rank ground truth, we set arbitrary $d - r$ singular values to 0 following a SVD decomposition. Each experiment in this section has been run over 3 independent seeds and the average is plotted. As the variance is small and the plots usually overlap, we opt to not plot it for visual clarity.

Theorems 4.1 to 4.3 provide rates on estimation error for empirical minimizers. In the following, we study these rates empirically for various values of p', p, d, N, T and r where $r = d$ for full-rank experiments or $r = 5$ for low-rank experiments and $p' = p$ except it is stated otherwise. We use two quantities, $\beta = NT$, the number of total tokens, $\gamma = pdr$, the number of parameters to estimate, to summarize information in the plots. For Theorems 4.1 and 4.3, $\hat{\mathbf{A}}$ is computed with the OLS estimator and for Theorem 4.2, $\hat{\mathbf{A}}$ is learned with gradient descent with learning rate α on the group-norm regularized loss in Equation (10). The parameter λ and learning rate α are tuned by a grid search.

Figure 1 plots the estimation error for $d \in \{5, 10, 15\}, p \in \{5, 10, 15\}, N \in \{1, 5, 10\}$ and $T \in \{1, 5, 10, 25, 50\} \times pdr/N$. The upper bound in Theorem 4.1 scales with the ratio β/γ up to logarithmic terms as empirically verified by Figure 1. In Figure 2, we verify that there is no individual trend to p and d , which implies that the error depends only on γ . Furthermore, we show the trend in N can be accounted for by incorporating the logarithmic term into β to obtain $\tilde{\beta} = \beta/\ln(1 + \sqrt{N})$.

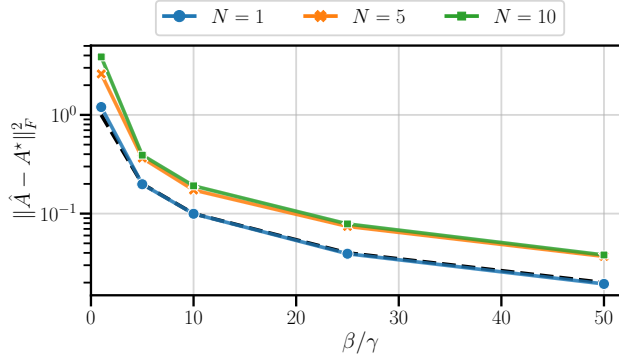


Figure 1: Scaling of estimation error with respect to the ratio $\beta/\gamma = NT/pd^2$ with the OLS estimator. The black dashed line plots γ/β for reference.

Figure 3 plots the estimation error for different degrees of misspecification where the context length is fixed to $p = 15$. The curves for various $p' \in \{5, 10, 15\}$ overlap, which verify the rate $\gamma/\beta = p'd^2/NT$ predicted by Theorem 4.3 holds.

Figure 4 repeats the same plots for low-rank experiments where $d = 15, r = 5$ are fixed and p, N and T are varied as before. Good estimation of \mathbf{A} is not straightforward as λ has to be appropriately tuned. Yet, we see that the group-nuclear norm regularized estimators found with gradient descent after tuning on regularization problem $\lambda \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}\}$ and learning rate $\alpha \in \{10^{-1}, 10^{-2}, 10^{-3}\}$ obtain improved estimation errors than non-regularized OLS estimator. Particularly, the sample efficiency benefits of the group-nuclear norm regularization are amplified in the low-data regime. We leave the analysis of group-nuclear norm regularization as a future work.

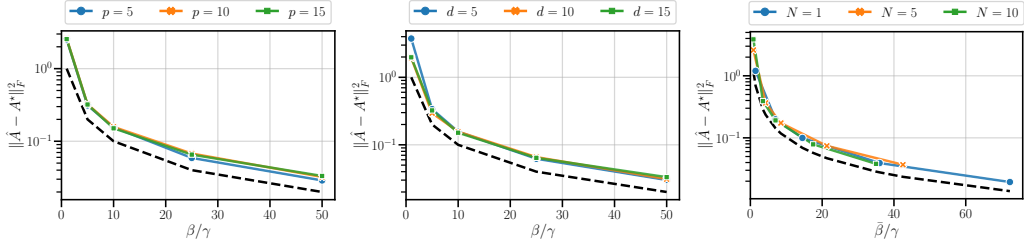


Figure 2: Scaling of estimation error for different values of p , d and N with the OLS estimator. Recall that $\beta = NT$, $\gamma = pd^2$ and $\tilde{\beta} = \beta / \ln(1 + \sqrt{N})$. Black dashed lines are drawn for reference and equals to $\sqrt{\gamma/\beta}$.

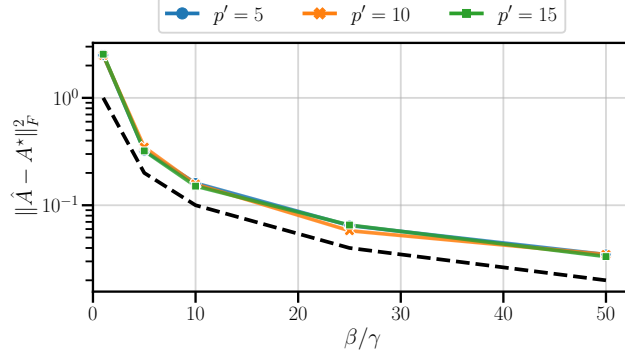


Figure 3: Scaling of estimation error with respect to the ratio $\beta/\gamma = NT/p'd^2$ for different $p' = 5, 10, 15$ with the OLS estimator. The black dashed line plots γ/β for reference.

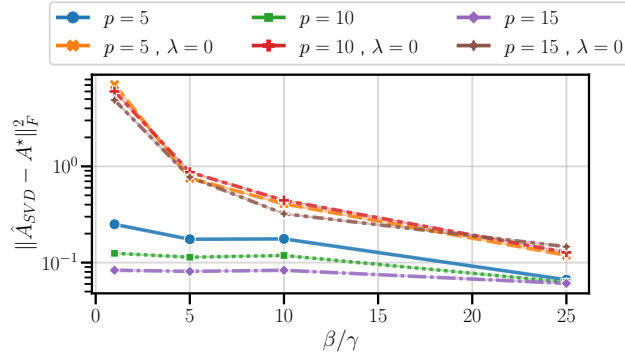


Figure 4: Scaling of estimation error with respect to $\beta/\gamma = \frac{NT}{pdr}$ for different context windows $p = 5, 10, 15$ with the OLS estimator ($\lambda = 0$) and group-nuclear norm regularized estimators.