Regression for matrix-valued data via Kronecker products factorization

Yin-Jen Chen

Minh Tang

Department of Statistics, North Carolina State University

May 1, 2024

Abstract

We study the matrix-variate regression problem $Y_i = \sum_k \beta_{1k} X_i \beta_{2k}^\top + E_i$ for i = 1, 2..., nin the high dimensional regime wherein the response Y_i are matrices whose dimensions $p_1 \times p_2$ outgrow both the sample size n and the dimensions $q_1 \times q_2$ of the predictor variables X_i i.e., $q_1, q_2 \ll n \ll p_1, p_2$. We propose an estimation algorithm, termed KRO-PRO-FAC, for estimating the parameters $\{\beta_{1k}\} \subset \Re^{p_1 \times q_1}$ and $\{\beta_{2k}\} \subset \Re^{p_2 \times q_2}$ that utilizes the Kronecker product factorization and rearrangement operations from Van Loan and Pitsianis (1993). The KRO-PRO-FAC algorithm is computationally efficient as it does not require estimating the covariance between the entries of the $\{Y_i\}$. We establish perturbation bounds between $\hat{\beta}_{1k} - \beta_{1k}$ and $\hat{\beta}_{2k} - \beta_{2k}$ in spectral norm for the setting where either the rows of E_i or the columns of E_i are independent sub-Gaussian random vectors. Numerical studies on simulated and real data indicate that our procedure is competitive, in terms of both estimation error and predictive accuracy, compared to other existing methods.

Keywords: matrix regression, Kronecker product, low-rank approximation, matrix perturbations

1 Introduction

Regression is one of the most important and widely studied inference tasks in statistics and machine learning. Traditional applications of regression mainly focus on settings where the response variables Y are either scalars or, more generally, $Y \in \Re^d$ for some "small" d. With the recent advancements in computation and storage technology, it is now common to encounter scenarios where the responses are (large) matrices. Examples include data from multivariate bioassay study (Vølund, 1980), electroencephalography (Li and Zhang, 2017), images denoising (Kamm and Nagy, 1998; Nagy, 1996; Cai et al., 2019), and factors models in econometrics (Chen et al., 2019; Wang et al., 2019).

These type of data naturally leads to the simple and intuitive notion of matrix-variate regression wherein, given a collection of predictor and response tuples $\{(X_i, Y_i)\}_{i=1}^n$ with $X_i \in \Re^{q_1 \times q_2}$ and $Y_i \in \Re^{p_1 \times p_2}$, one typically assumes that the Y_i are related to the X_i through the linear model

$$\operatorname{vec}(Y_i) = \nu \operatorname{vec}(X_i) + \operatorname{vec}(E_i), \quad i = 1, ..., n$$
(1)

where $\nu \in \Re^{p_1 p_2 \times q_1 q_2}$ are the unknown regression coefficients, $E_i \in \Re^{p_1 \times p_2}$ are the unobserved noise matrices, and **vec** denote the vectorization operator that concatenates the column vectors of the input matrix. We note that any linear model for $\{(X_i, Y_i)\}$ can be written in the form of Eq. (1) In the high-dimensional regime wherein the dimensions of response variables $(p_1.p_2)$ grow much faster than the sample size n, i.e., $p_i/n \to \infty$, the regression coefficient ν is overparameterized and consistent estimation of ν is generally unfeasible unless one impose some structural assumptions on Eq. (1) so as to reduce the effective number of parameters in ν .

Two of the most widely studied and adopted regularity conditions for ν is that it is low-rank and/or sparse; see e.g., Yuan et al. (2007); Obozinski et al. (2011); Negahban and Wainwright (2011); Chen et al. (2012); Bunea et al. (2012); Bing and Wegkamp (2019); Zheng et al. (2019); Zou et al. (2020) and the references therein. In particular Negahban and Wainwright (2011) noted that low-rank constraints are analogous to imposing sparsity on the data without explicitly specifying any basis.

Despite the popularity of these sparse and/or low-rank assumptions, they do not lead to a significant reduction in complexity of ν when the feature vectors $\{X_i\}$ are low-dimensional but the response $\{Y_i\}$ are high-dimensional. More specifically, suppose $p_1 = p_2 = p$, $q_1 = q_2 = q$, $p \gg q$ and $p/n \to \infty$. If we only assume that ν is low-rank so that $rk(\nu) = d \ll n$ then we still need to estimate on the order of $O(d(p^2 + q^2))$ parameters for any matrix factorization of ν (such as SVD) and is computationally infeasible as it is equivalent to the estimation of covariance matrices in high-dimensional univariate linear regression with p predictor variables and n scalar responses. In contrast if we assume sparsity on ν then, denoting the number of non-zero entries in ν by s, we will in general need $n = \omega(s)$ to estimate ν consistently. As $s/p^2 \ll p^{-1}$, this implies that almost all of the entries in the responses $\{Y_i\}$ are ignorable. This is a rather strong assumption that should be justified on a case-by-case basis.

In this paper we consider a more refined variant of the low-rank assumption on ν by assuming that it admits a representation/ approximation in terms of a (sum of) Kronecker products of smaller matrices. More specifically, we shall assume that ν is of the form

$$\nu = \sum_{k=1}^d \beta_{2k} \otimes \beta_{1k}$$

for some collection of $p_1 \times q_1$ matrices $\{\beta_{11}, \ldots, \beta_{1d}\}$ and $p_2 \times q_2$ matrices $\{\beta_{21}, \ldots, \beta_{2d}\}$. The number of effective parameters in ν is then $d(p_1q_1 + p_2q_2)$ which is substantially smaller than $O(d(p_1p_2 + q_1q_2))$ for $p_1 \gg q_1, p_2 \gg q_2$. See Beylkin and Mohlenkamp (2002); De Lathauwer et al. (2000); Tyrtyshnikov (2004) for futher discussion of Kronecker product factorization and its use in large-scale matrix approximations.

Finally the linear model in Eq. (1) with the Kronecker product structure for ν is equivalent to the *bi-linear model*

$$Y_{i} = \sum_{k=1}^{d} \beta_{1k} X_{i} \beta_{2k}^{\top} + E_{i}.$$
 (2)

Under this perspective the $\{\beta_{1k}\}$ (resp. $\{\beta_{2k}\}$) can be interpreted as the row effects (resp. column effects) of X_i on the response Y_i . The special case of d = 1 was considered previously in Ding and Cook (2016) wherein the authors studied estimation of β_{11} and β_{21} using two-step MLEs; see Section 2 for further discussions. In a related vein, Chen et al. (2019); Wang et al. (2019); Chen and Fan (2021) considered factor models for Y_i of the form $Y_i = \beta_1 X_i \beta_2^{\top}$ but, in contrast to the current paper, they assume that the X_i are either unknown or unobserved. They then propose to estimate β_1 and β_2 via two-step PCA Finally, the bi-linear modeling of $\{Y_i\}$ also arise in the context of image recognition (Crainiceanu et al., 2011; Wang et al., 2016; Ye, 2005; Zhang, 2005). In particular Crainiceanu et al. (2011) proposed the notion of population value decomposition for summarizing images population $\{Y_i\}$ by assuming that $Y_i \approx PV_iD$ where P and D encode "population frame of reference" for all $\{Y_i\}$ while V_i encode "subject-level" features specific to a given Y_i . Their P and D thus serve identical roles to that of $\{\beta_{11}, \beta_{21}\}$ in Eq. (2) (when d = 1).

In this paper we study estimation of $\{\beta_{1k}, \beta_{2k}\}$ for the model in Eq. (2). Inspired by the work of Van Loan and Pitsianis (1993) on the nearest Kronecker product problem, we observe that ν exhibits a low-rank representation after *reshaping and rearranging* the entries of ν . In other words, while $\nu = \sum_{k=1}^{d} \beta_{2k} \otimes \beta_{1k}$ itself need not be low-rank, its rearranged version still admits a low-rank representation or approximation. Leveraging this observation we propose an algorithm, termed KRO-PRO-FAC, for estimating ν with computational complexity of $O(p_1p_2dq_1q_2)$ flops; see Section 2). We next studied the theoretical properties of the KRO-PRO-FAC algorithm and show that it yield, under reasonably mild conditions on the noise E_i and the dimensions p_i (compared with the sample size n), consistent estimates of the $\{\beta_{1k}, \beta_{2k}\}$; see Section 3. Numerical experiments on simulated and real data are presented in Section 4. In particular our procedure is shown to be competitive, in terms of both estimation error and predictive accuracy, to other existing methods.

2 Methodologies

We now introduce some basic notations used throughout this paper. For $p \in \mathbb{N}$, we denote the set $\{1, ..., p\}$ by [p]. Let $\mathcal{O}(\cdot)$, $\mathcal{O}(\cdot)$ and $\Theta(\cdot)$ represent the standard big-O, little-o and big-theta relationships. For two arbitrary real sequences $(a_n)_{n \in \mathbb{N}}$ and $(b_n)_{n \in \mathbb{N}}$, we write $a_n \ll b_n$ if a_n/b_n converges to 0 as $n \to \infty$; $a_n \asymp b_n$ means a_n/b_n has a finite and non-zero limit as $n \to \infty$. For an arbitrary matrix $M = (M_{ij}) \in \Re^{p \times q}$, the Frobenius norm, spectral norm and nuclear norm of M are denoted by $||M||_F$, $||M||_2$ and $||M||_*$, and if M is square then $\mathbf{tr}(M)$ and |M| denote its trace and determinant. The symbol ' \otimes ' represents the Kronecker product between matrices while I_p denote the $p \times p$ identity matrix. The vectorization of a $p \times q$ matrix M is defined as

$$\operatorname{vec}(M) = [a_{11} \dots a_{p1} \dots a_{12} \dots a_{p2} \dots a_{p1} \dots a_{pq}]^T \in \Re^{pq}.$$

and we denote its inverse by $\operatorname{vec}^{-1}(m, p, q)$ where m is a vector in \Re^{pq} .

2.1 Dual Kronecker products structure

We first discuss the special case of Eq. (2) with d = 1, i.e., given a collection of matrix-variate predictors $\{X_i\}_{i=1}^n \subset \Re^{q_1 \times q_2}$ and matrix-variate responses $\{Y_i\}_{i=1}^n \subset \Re^{p_1 \times p_2}$, we consider the bilinear model of the form

$$Y_i = \beta_1 X_i \beta_2^\top + E_i, \quad i \in [n]$$
(3)

where $\beta_1 \in \Re^{p_1 \times q_1}$ and $\beta_2 \in \Re^{p_2 \times q_2}$ are the unknown regression coefficients and E_i are unobserved noise matrices. Under this model each column (resp. row) of Y_i is a noisy perturbation of some linear combination of the columns of β_1 (resp. rows of β_2^{\top}). Recall that Eq. (3) can be rewritten in vectors form as

$$\operatorname{vec}^{\top}(Y_i) = \operatorname{vec}^{\top}(X_i)(\beta_2^{\top} \otimes \beta_1^{\top}) + \operatorname{vec}^{\top}(E_i), \quad i \in [n]$$

$$\tag{4}$$

and thus, by collecting all the $\mathbf{vec}^{\top}(Y_i)$ into a matrix and letting $\nu = \beta_2 \otimes \beta_1$, leads to the linear regression model

$$\mathcal{Y} = \mathcal{X}\nu^T + \mathcal{E} \tag{5}$$

where \mathcal{Y} and \mathcal{E} are $n \times p_1 p_2$ matrices whose rows are the $\operatorname{vec}^{\top}(Y_i)$ and $\operatorname{vec}^{\top}(E_i)$ respectively while the design matrix $\mathcal{X} \in \Re^{n \times q_1 q_2}$ has rows $\operatorname{vec}^{\top}(X_i)$. Several variants of formulation in Eq. (4) and Eq. (5) have been discussed in the literature; see e.g., Zhao and Leng (2014) for the case of $p_1 = p_2 = 1$ and Kong et al. (2019) for the case of $q_1 = q_2 = 1$ and ν being low-rank. Here Eq. (5) assumes that the mean coefficient ν admits a Kronecker product representation of β_2 and β_1 . Without the Kronecker product factorization, ν can be easily overparameterized with $\mathcal{O}(p_1 p_2 q_1 q_2)$ elements to be estimated. While ν is identifiable, is parametrization in terms of β_1 and β_2 is only identifiable up to a constant, i.e., $\beta_2 \otimes \beta_1 = c\beta_2 \otimes c^{-1}\beta_1$ for any non-zero constant c.

As we allude to in the introduction, the bi-linear model in Eq. (3) had been studied previously in Ding and Cook (2016) and we now describe the pertinent details of this work in the context of the current paper. Denote the covariance matrix of E_i as $\Sigma_{\text{vec}(E)}$. Ding and Cook (2016) then assume that $\Sigma_{\text{vec}(E)}$ can be decomposed as $\Sigma_{\text{vec}(E)} = \Sigma_2 \otimes \Sigma_1$; here Σ_1 and Σ_2 represent the covariance matrix for the rows and and columns of Y_i respectively, i.e.,

$$\Sigma_{1} = \mathbb{E}\{(Y_{1} - \mathbb{E}(Y_{1}))(Y_{1} - \mathbb{E}(Y_{1}))^{T}\}, \quad \Sigma_{2} = \mathbb{E}\{(Y_{1} - \mathbb{E}(Y_{1}))^{T}(Y_{1} - \mathbb{E}(Y_{1}))\}$$
(6)

The above structure for $\Sigma_{\text{vec}(E)}$ is quite natural for longitudinal data wherein each subject is measured repeatedly over two different domains. For example the rows of Y_i can record measurements over time while the columns of Y_i record different covariates. Given Eq. (6), the number of parameters in $\Sigma_{\text{vec}(E)}$ is then reduced drasticically from $\mathcal{O}(p_1^2 p_2^2)$ to $\mathcal{O}(p_1^2 + p_2^2)$. This in turn allows the dimension p_1 and p_2 to possibly outgrow the sample size n, i.e., $n \ll \min\{p_1, p_2\}$.

Ding and Cook (2016) then consider MLE estimation of $\beta_1, \beta_2, \Sigma_1, \Sigma_2$ by further assuming that the E_i follows the matrix normal distribution, i.e., $\text{vec}(E_i) \sim \mathcal{N}(0, \Sigma_{\text{vec}(E)})$; for more on the matrix normal distribution see De Waal (1985); Gupta and Nagar (1999) and the references therein. With the above Kronecker product structure for the mean and covariance of Y_i , the log-likelihood of $\{Y_i\}$ given $\{X_i\}$ is (ignoring unimportant constants)

$$2\ell(\{Y_i\};\theta) = -np_1 \ln |\Sigma_1| - np_2 \ln |\Sigma_2| - \sum_{i=1}^n \operatorname{tr} \left\{ \Sigma_2^{-1} (Y_i - \beta_1 X_i \beta_2^{\top})^T \Sigma_1^{-1} (Y_i - \beta_1 X_i \beta_2^{\top}) \right\}$$
(7)

where $\theta = (\beta_1, \beta_2, \Sigma_1, \Sigma_2)$. Let $\hat{\theta}$ denote the MLE of θ from Eq. (7). As there are no closed-form expression $\hat{\theta}$, Ding and Cook (2016) proposed a two-stage iterative algorithm for finding $\hat{\theta}$ that is motivated by earlier work of Dutilleul (1999). More specifically the algorithm sequentially updates either the row parameters (β_1, Σ_1) or the column parameters (β_2, Σ_2) , with the remaining parameters hold fixed. While the dual Kronecker product structure and the resulting MLE procedure provides a convenient way to model both the mean and covariance of the rows (and columns) simultaneously, there are nevertheless two major concerns regarding this approach. Firstly the MLE procedure is guaranteed to converge only to a *stationary* point but not necessarily a *global* optimum. Secondly, the update for (β_1, Σ_1) (resp. (β_2, Σ_2)) require inverting Σ_2 (resp. Σ_1) and thus each updates involve possibly $\mathcal{O}(n(p_1^2p_2 + p_1p_2^2))$ flops, which is computationally prohibitive for moderate and/or large values of p_1 and p_2 . In light of the above drawbacks we propose in Section 2.2 a more computationally efficient procedure which estimates only β_1 and β_2 but not Σ_1, Σ_2 or Cov[vec(Y_i)].

2.2 Kronecker products factorization and low-rank approximation

If we assume a high-dimensional setting where the sample size n is small or comparable to the dimensions $\min\{p_1, p_2\}$ of the response then it is generally the case that we can not estimate $\operatorname{Cov}(Y_i)$ to any reasonable degree of accuracy. One simple and intuitive remedy to this issue is to ignore the structure in $\operatorname{Cov}(Y_i)$ and instead focus our effort on estimating ν .

Our starting point is the observation that although the OLS estimate $\tilde{\nu} = [(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}\mathcal{Y}]^{\top}$ is a simple and elegant estimate of ν , it does not share the same Kronecker product structure as that for $\nu = (\beta_2 \otimes \beta_1)$. It is thus natural to consider projecting $\tilde{\nu}$ onto the set formed by Kronecker products of matrices with appropriate dimensions. In particular let p_1, q_1, p_2, q_2 be positive integers and M be a matrix of dimensions $p_1p_2 \times q_1q_2$. The nearest Kronecker product approximation to M with respect to the dimensions $\{p_i, q_i\}$ is defined as

$$\underset{\beta_1 \in \Re^{p_1 \times q_1}, \beta_2 \in \Re^{p_2 \times q_2}}{\operatorname{smin}} \| M - \beta_2 \otimes \beta_1 \|_F^2.$$
(8)

Van Loan and Pitsianis (1993) showed that Eq. (8) has a closed-form solution given by the truncated SVD of a *rearranged* version of M. More specifically first partition M into smaller matrices $M_{ij} \in \Re^{p_1 \times q_1}$ for $1 \le i \le p_2$ and $1 \le j \le q_2$, i.e.,

$$M = \begin{bmatrix} M_{11} & M_{12} & \cdots & M_{1q_2} \\ M_{21} & M_{22} & \cdots & M_{2q_2} \\ \vdots & \vdots & \ddots & \vdots \\ M_{p_21} & M_{p_22} & \cdots & M_{p_2q_2}. \end{bmatrix}$$
(9)

Next define the rearrangement operation $R(\cdot)$: $\Re^{p_1p_2 \times q_1q_2} \to \Re^{p_2q_2 \times p_1q_1}$ by

$$R(M) = \begin{bmatrix} A_1 \\ A_2 \\ \cdots \\ A_{q_2} \end{bmatrix}, \qquad A_j = \begin{bmatrix} \operatorname{vec}(M_{1j})^\top \\ \operatorname{vec}(M_{2j})^\top \\ \cdots \\ \operatorname{vec}(M_{p_2j})^\top \end{bmatrix}.$$
(10)

We emphasize that R(M) and M generally have different dimensions. In particular, if $p_1 \simeq p_2$, $q_1 \simeq q_2$ and $q_i \ll p_i$ then M is a *tall* matrix but the dimensions of R(M) are comparable. The solution of Eq. (8) is then equivalent to finding the closest rank-1 representation of R(M), i.e.,

$$\min_{\beta_1,\beta_2} \|M - \beta_2 \otimes \beta_1\|_F = \min_{\beta_1,\beta_2} \|R(M) - \operatorname{vec}(\beta_2)\operatorname{vec}(\beta_1)^\top\|_F$$
(11)

and thus, by the Eckart-Young-Mirsky theorem (Eckart and Young, 1936), we can take $\beta_1 = \sigma_1^{1/2} \operatorname{vec}^{-1}(\mathcal{V}_1, p_1, q_1)$ and $\beta_2 = \sigma_1^{1/2} \operatorname{vec}^{-1}(\mathcal{U}_1, p_2, q_2)$ where σ_1, \mathcal{U}_1 and \mathcal{V}_1 are the leading singular values and (left and right) singular vectors of R(M), respectively. See Van Loan and Pitsianis (1993) for more details. In summary if we assume a Kronecker product structure for the regression coefficient ν in Eq. (5) then our estimate for ν, β_1 and β_2 is given by

Step 1: Let $\tilde{\nu} \leftarrow [(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}\mathcal{Y}]^{\top} \in \Re^{p_1p_2 \times q_1q_2}$ be the OLS estimate of ν and let $R(\tilde{\nu})$ be its Pitsianis-Van Loan rearrangement; see Eq. (10).

Step 2: Compute the SVD $R(\tilde{\nu}) = \sum_{k=1}^r \hat{\sigma}_k \hat{\mathcal{U}}_k \hat{\mathcal{V}}_k^{\top}$ with $\hat{\sigma}_1 \geq \cdots \geq \hat{\sigma}_r$ and $r \leq \min\{p_1q_1, p_2q_2\}$.

Step 3: Let $\operatorname{vec}(\hat{\beta}_2) = \hat{\sigma}_1^{1/2} \hat{\mathcal{U}}_1$ and $\operatorname{vec}(\hat{\beta}_1) = \hat{\sigma}_1^{1/2} \hat{\mathcal{V}}_1$

Step 4: Output the estimate $\hat{\nu} \leftarrow \hat{\beta}_2 \otimes \hat{\beta}_1$ for ν .

We emphasize that, despite the close connection between Kronecker products approximation and low-rank approximations described in Eq. (11), the assumption of a Kronecker factorization for ν is quite different from the assumption that ν is low-rank. Indeed, the rank of ν can be as large as q_1q_2 (assuming $p_1p_2 \ge q_1q_2$) even when $R(\nu)$ is a rank-1 matrix. This difference distinguishes our work from those which introduce penalty terms to induce low-rank structure on ν directly; see e.g., Kong et al. (2019); Wang et al. (2021); Feng et al. (2021) for recent examples of this latter approach. We now consider a simple simulation study to further illustrate this distinction.

Example. We set the dimensions of Y_i and X_i as $p_1 = p_2 = q_1 = q_2 = 10$. We then generate n = 3000 samples of the $\{(X_i, Y_i)\}$ pair according to the model $Y_i = \beta_1 X_i \beta_2^\top + E_i$ where the $\{X_i\}$ are iid random vectors with $\operatorname{vec}(X_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the $\operatorname{vec}(E_i)$ are also iid random vectors with $\operatorname{vec}(X_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the $\operatorname{vec}(E_i)$ are also iid random vectors with $\operatorname{vec}(E_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Given the $\{(X_i, Y_i)\}$ we first compute the OLS estimate $\tilde{\nu}$ and its rearranged version $R(\tilde{\nu})$. Next define, for a matrix M and an integer $k \geq 1$, the function $f_k(M) = (\sum_{i=1}^k \sigma_i(M))/||M||_*$ corresponding to the (normalized) sum of the k largest singular values of M; here $||M||_*$ denotes the nuclear norm of M. We then compute, for each $k \in [100]$, the quantity $f_k(\tilde{\nu})$ and $f_k(R(\tilde{\nu}))$. Finally we repeat the above steps for 1000 Monte Carlo replicates. We note that the β_1 and β_2 are fixed constants and do not vary with the Monte Carlo replicates.

Figure 1 plots the (normalized) cumulative sum of the first k singular values of $\tilde{\nu}$ and $R(\tilde{\nu})$ for k varying in $\{1, 2, ..., 100\}$; note that $f_k(\tilde{\nu}) = f_k(R(\tilde{\nu})) = 1$ when k = 100. From Figure 1 we see that the largest singular value of $R(\tilde{\nu})$ accounts for, on average, roughly 87% of $||R(\tilde{\nu})||_*$ and thus a rank-1 approximation of $R(\tilde{\nu})$ is expected to preserve most of the information in $R(\tilde{\nu})$ while also removing the noise from the small singular values in $R(\tilde{\nu})$. In contrast the largest singular value of $\tilde{\nu}$ only explains 5% of $||\tilde{\nu}||_*$ and thus computing $\hat{\nu}$ using low-rank approximations to $\tilde{\nu}$ is possibly problematic.

2.3 KRO-PRO-FAC algorithm

A natural extension of the optimization problem in Eq. (8) is to approximate a matrix ν using a sum of d Kronecker products which, by the above discussions, can be related to the sum of d



Figure 1: cumulative singular value (averaged over 1000 replications) for (a) the OLS estimate $\tilde{\nu}$ and (b) the rearranged estimate $R(\tilde{\nu})$

rank-1 matrices via

$$\underset{\{(\beta_{1k},\beta_{2k})\}}{\arg\min} \|\nu - \sum_{k=1}^{d} \beta_{2k} \otimes \beta_{1k}\|_{F} = \underset{\{(\beta_{1k},\beta_{2k})\}}{\arg\min} \|R(\nu) - \sum_{k=1}^{d} \operatorname{vec}(\beta_{2k}) \operatorname{vec}(\beta_{1k})^{\top}\|_{F}.$$
(12)

A solution of Eq. (12) is then once again given by the truncated SVD of $R(\nu)$. Eq. (12) furthermore suggests a more general version of the regression problem in Eq. (3), namely that

$$Y_{i} = \sum_{k=1}^{d} \beta_{1k} X_{i} \beta_{2k}^{\top} + E_{i}, \quad i \in [n]$$
(13)

with $d \ll n \ll \min\{p_1q_1, p_2q_2\}$. Eq. (13) can be rewritten as

$$\mathcal{Y} = \mathcal{X}\nu^T + \mathcal{E}, \quad \nu = \sum_{k=1}^d \beta_{2k} \otimes \beta_{1k}.$$
 (14)

Here we refer to d in Eq. (14), as the Kronecker product rank of ν . For ease of exposition (and without loss of generality) we shall assume that the $\{\beta_{1k}, \beta_{2k}\}$ are orthogonal, i.e., $\operatorname{vec}(\beta_{1s})^{\top}\operatorname{vec}(\beta_{1t}) = \operatorname{vec}(\beta_{2s})^{\top}\operatorname{vec}(\beta_{2t}) = 0$ for all $s \neq t$ and $\|\beta_{1s}\|_F = \|\beta_{2s}\|_F$. Our estimate for ν and $\{\beta_{1k}, \beta_{2k}\}$ proceeds in an analogous manner to that described in Section 2.2. In particular we first compute the OLS estimate $\tilde{\nu} = [(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}\mathcal{Y}]^{\top}$, then rearrange $\tilde{\nu}$ to obtain $R(\tilde{\nu})$, and finally compute the truncated SVD of $R(\tilde{\nu})$ to keep only the *d* largest singular values and singular vectors. We termed this procedure as the KRO-PRO-FAC (Kronecker product factorization) estimate of ν . See Algorithm 1 for a more formal descriptions. As the Kronecker product rank *d* of ν is generally unknown, we estimate it using the ratio of singular values as described in Lam and Yao (2012) and Ahn and Horenstein (2013), i.e., we estimate *d* by

$$\hat{d} = \underset{j \in \{1, \dots, \bar{d}\}}{\arg \max \hat{\sigma}_j / \hat{\sigma}_{j+1}}$$
(15)

where d is a pre-specified constant and $\hat{\sigma}_k$'s are the singular values of $R(\tilde{\nu})$ in a descending order.

The computational complexity of the KRO-PRO-FAC algorithm is $\mathcal{O}(p_1p_2q_1q_2\min\{p_1q_1, p_2q_2\})$ with the main computational bottleneck being the SVD of $R(\tilde{\nu})$. If d is either known or is estimated to be much smaller than the dimensions of $R(\tilde{\nu})$ then the cost of the SVD step reduces to $\mathcal{O}(p_1p_2q_1q_2d)$ flops by using either Lanczos bidiagonalization and/or randomized SVD, see e.g., Halko et al. (2011); Musco and Musco (2015); Tropp et al. (2017) and the references therein. Hence the complexity for the full algorithm itself drops to $\mathcal{O}(np_1p_2q_1q_2)$. In contrast, any algorithm that involves estimating the covariance matrices for the rows and/or columns will requires at least $\mathcal{O}(n(p_1^2p_2 + p_1p_2^2))$ flops which is an enormous computational burden for large values of p_1 and/or p_2 .

Remark 1. We note that even if ν does not have the form as specified in Eq. (13) it can nevertheless be well-approximated by a sum of Kronecker products. Kronecker products provide a computational efficient building block for approximating large matrices in numerical linear algebra application. See Beylkin and Mohlenkamp (2002); De Lathauwer et al. (2000); Tyrtyshnikov (2004) for some general theory and see (Kamm and Nagy, 1998; Nagy, 1996; Werner et al., 2008; Greenewald and Hero, 2015) for specific examples in image restoration and covariance estimation. We emphasize that if M is a $p_1p_2 \times q_1q_2$ matrix with $p_1p_2 \gg q_1q_2$ then a rank d SVD of M will require computing left singular vectors of length p_1p_2 while its Kronecker product factorization only require computing factors of dimensions $p_1 \times q_1$ and $p_2 \times q_2$.

Remark 2. We note that Kronecker products factorization also featured prominently in the work of Cai et al. (2019) but their research question is subtantially different from that considered in the

current paper. In particular our setting is that of linear regression where the goal is to estimate the factorization (β_1, β_2) of ν given both the responses $\{Y_i\}$ and feature vectors $\{X_i\}$, i.e., our estimation of (β_1, β_2) is a *supervised learning* problem. In contrast Cai et al. (2019) uses Kronecker product approximation to perform dimension reduction of the $\{Y_i\}$ without observing any $\{X_i\}$, i.e., they are considering an *unsupervised* learning problem.

Algorithm 1: KRO-PRO-FAC algorithmInput: \mathcal{Y}, \mathcal{X} and $(p_2, q_2), (p_1, q_1), \bar{d}$ Output: $(\hat{\beta}_{1k}, \hat{\beta}_{2k})$ 1 Compute the OLS estimate $\tilde{\nu} \leftarrow [(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}\mathcal{Y}]^{\top}$.2 Rearrange $\tilde{\nu}$ to get $R(\tilde{\nu})$ by Eq. (10).3 Perform SVD on $R(\tilde{\nu})$, i.e., $R(\tilde{\nu}) = \sum_{k=1}^{\bar{d}} \hat{\sigma}_k \hat{\mathcal{U}}_k \hat{\mathcal{V}}_k^{\top}$ with $\hat{\sigma}_1 \geq \hat{\sigma}_2 \geq \cdots \geq \hat{\sigma}_{\bar{d}}$.4 Estimate d by $\hat{d} = \arg \max_{j \in \{1, \dots, \bar{d}\}} \hat{\sigma}_j / \hat{\sigma}_{j+1}$ 5 Set $\operatorname{vec}(\hat{\beta}_{2k}) = \hat{\sigma}_k^{1/2} \hat{\mathcal{U}}_k$ and $\operatorname{vec}(\hat{\beta}_{1k}) = \hat{\sigma}_k^{1/2} \hat{\mathcal{V}}_k$.6 Output $\hat{\nu} \leftarrow \sum_{k=1}^{\hat{d}} \hat{\beta}_{2k} \otimes \hat{\beta}_{1k}$

3 Theoretical Results

We now study large-sample and/or asymptotic results for the estimates of $\{\beta_{1k}, \beta_{2k}\}$ obtained by the KRO-PRO-FAC algorithm. Recall that, from our earlier discussions in Section 2, the rearranged OLS estimate $R(\tilde{\nu})$ can be viewed as a sum of rank-1 matrices $R(\nu) = \sum_{k=1}^{d} \operatorname{vec}(\beta_{2k}) \operatorname{vec}(\beta_{1k})^{\top}$ additively perturbed by the noise matrix $R(\tilde{\mathcal{E}})$ where $\tilde{\mathcal{E}} = [(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}\mathcal{E}]^{\top}$. Therefore, if $\|\tilde{\mathcal{E}}\|$ is sufficiently small compared to $\|R(\nu)\|$, then we can apply classical results in matrix perturbation theory such as the sin- Θ theorem (Wedin, 1972) to show that the leading singular vectors of $R(\tilde{\nu})$ are "close" to the $\operatorname{vec}(\beta_{1k})$ and $\operatorname{vec}(\beta_{2k})$.

We now make the above description precise. Let $R(\nu)$ be a rank d matrix for some fixed constant d not depending on p_1, p_2 and n. Denote the SVD of $R(\nu)$ by $R(\nu) = \mathcal{U}\mathcal{D}\mathcal{V}^{\top}$ where $\mathcal{D} = diag(\sigma_k)$ is a $d \times d$ diagonal matrix of singular values, $\mathcal{V} = (\mathcal{V}_1, \ldots, \mathcal{V}_d)$ is a $p_1q_1 \times d$ orthonormal matrix of right singular vectors and $\mathcal{U} = (\mathcal{U}_1, \ldots, \mathcal{U}_d)$ is a $p_2q_2 \times d$ orthonormal matrix of left singular vectors. Next let $\hat{\mathcal{U}}\hat{\mathcal{D}}\hat{\mathcal{V}}^{\top}$ denote the *truncated* SVD corresponding to the d largest singular values and singular vectors of $R(\tilde{\nu})$. We first make an assumption on the relationship between the matrix dimensions p_1 , p_2 , q_1 , q_2 and the sample size n as well as the growth rate for the singular values of $R(\nu)$.

Condition 1. Let p_1, p_2, q_1, q_2 and n satisfy

$$\frac{q_1}{q_2} = \Theta(1), \quad \frac{p_1}{p_2} = \Theta(1), \quad q_1 q_2 \ll n, \quad \ln p_i = o(n), \ i = 1, 2.$$

Furthermore, for sufficiently large p_1, p_2 , assume that the singular values of $R(\nu)$ satisfy

$$\sigma_k = \mathcal{O}(p_1), \ i = 1, 2, \dots, d$$

Condition 1 implies that $R(\nu)$ have bounded condition number.

We next recall the notion of a sub-Gaussian random vector

Definition 1. Let Z be a mean zero random variable. Then Z is said to be sub-Gaussian with variance proxy σ^2 if, for all t > 0 we have

$$\mathbb{P}(|Z| > t) \le 2\exp\left(-\frac{t^2}{2\sigma^2}\right).$$
(16)

In other words, the tail probability of Z behaves similarly to that of a Gaussian distribution with variance σ^2 . A mean zero random vector $\mathbf{Z} \in \Re^p$ is then said to be a sub-Gaussian random vector with covariance proxy Σ if $w^{\top} \mathbf{Z}$ is sub-Gaussian with variance proxy $w^{\top} \Sigma w$ for all $w \in \Re^p$. See Section 2.5 and Section 3.4 of Vershynin (2018) for further discussion and characterizations of sub-Gaussian random vectors.

Now let $\{\xi_1, \ldots, \xi_n\}$ be iid mean zero sub-Gaussian random vectors in $\Re^{p_1p_2}$ with covariance proxy \mathcal{I} where \mathcal{I} is the $p_1p_2 \times p_1p_2$ identity matrix. We shall assume that the noise matrices E_i are of the form

$$\operatorname{vec}(E_i) = \Sigma_{\operatorname{vec}(E)}^{1/2} \xi_i, \qquad i \in [n]$$
(17)

for some $p_1p_2 \times p_1p_2$ positive definite matrix $\Sigma_{\text{vec}(E)}^{1/2}$ satisfying the following condition.

Condition 2. $\Sigma_{\text{vec}(E)}$ is a block diagonal matrix, i.e., $\Sigma_{\text{vec}(E)} = \text{diag}(\Sigma_1, \Sigma_2, \dots, \Sigma_{p_2})$ where each diagonal block is of size $p_1 \times p_1$. Furthermore there exists a positive constant C independent of p_1 , p_2 and n such that

$$\max_{k \in [p_2]} \max_{s \in [p_1]} \Sigma_k(s, s) \le \mathcal{C}.$$
(18)

where $\Sigma_k(s,t)$ is the (s,t) entry of Σ_k .

Remark 3. We note that the block diagonal structure posited in Assumption 2 is different from and arguably more flexible than assuming a Kronecker product structure for $\Sigma_{\text{vec}(E)}$. More specifically an arbitrary $\Sigma_{\text{vec}(E)}$ has $O(p_1^2p_2^2)$ parameters. If $\Sigma_{\text{vec}(E)}$ can be factored into the Kronecker product of a $p_1 \times p_1$ matrix and a $p_2 \times p_2$ matrix then the number of parameters is reduced drastically to $O(p_1^2 + p_2^2)$ parameters. It was noted in Barratt (2018) that $O(p_1^2 + p_2^2)$ parameters is potentially too few as it preclude the use of some common matrix-variate Gaussian distribution to model $\Sigma_{\text{vec}(E)}$. In contrast, under Assumption (2), $\Sigma_{\text{vec}(E)}$ has $O(p_2p_1^2)$ parameters. If $p_1 \times p_2 \times p$ then the above three scenarios correspond to $O(p^4), O(p^2)$ and $O(p^3)$ parameters, respectively. Finally we note that Assumption 2 is satisfied whenever the columns of E_i are uncorrelated. A similar condition can be formulated for the case when the rows of E_i are uncorrelated. These conditions are milder than assuming that the entries of E_i are mutually independent as is done in Obozinski et al. (2011); Negahban and Wainwright (2011); Bunea et al. (2012); Bing and Wegkamp (2019); Zou et al. (2020).

With the above assumptions in place, we now state our theoretical results for bounding the estimation error between $\hat{\mathcal{U}}$ (resp. $\hat{\mathcal{V}}$) and \mathcal{U} (resp. \mathcal{V}). These errors are stated in terms of the sin- Θ distance between linear subspaces, i.e., given two orthonormal matrices \mathcal{W}_1 and \mathcal{W}_2 the sin- Θ distance between the linear subspaces spanned by \mathcal{W}_1 and \mathcal{W}_2 is defined as

$$\|\sin\Theta\left(\mathcal{W}_1,\mathcal{W}_2\right)\| = \sqrt{1 - \sigma_{\min}^2(\mathcal{W}_1,\mathcal{W}_2)}.$$
(19)

where $\sigma_{\min}(\mathcal{W}_1, \mathcal{W}_2)$ is the minimum singular value of $\mathcal{W}_1^{\top} \mathcal{W}_2$.

Theorem 1. Let $\{(X_i, Y_i)\}$ satisfy the linear model in Eq. (13) for some fixed $d \ge 1$ not depending on n and suppose that Condition 1 and 2 holds. Then there exists a constant C > 0 such that, with probability at least $1 - n^{-3}$, the following holds simultaneously,

$$|\hat{\sigma}_k - \sigma_k| \le C q_1 q_2 \max_{k \in p_2} \max_{s \in [p_1]} \Sigma_k(s, s) \frac{\sqrt{p_1} + \sqrt{p_2} + \ln p_1 + \ln p_2}{\sqrt{n}}, \quad (20)$$

 $\max\{\|\sin\Theta(\mathcal{U},\hat{\mathcal{U}})\|, \|\sin\Theta(\mathcal{V},\hat{\mathcal{V}})\|\} \le C q_1 q_2 \max_{k \in p_2} \max_{s \in [p_1]} \Sigma_k(s,s) \frac{\sqrt{p_1} + \sqrt{p_2} + \ln p_1 + \ln p_2}{\sqrt{np_1p_2}}.$ (21)

Theorem 1 implies the following upper bound for the error of $\hat{\nu} = \sum_{k=1}^{d} \hat{\beta}_{2k} \otimes \hat{\beta}_{1k}$ as an estimate for $\nu = \sum_{k=1}^{d} \beta_{2k} \otimes \beta_{1k}$. In particular the *relative* error of $\hat{\nu} - \nu$ converges to 0 as p_1, p_2 and n diverge and thus $\hat{\nu}$ is a consistent estimate for ν .

Corollary 1. Suppose $p_1 = p_2 = p$ and consider the setting in Theorem 1. Then there exists a constant C > 0 such that with probability at least $1 - n^{-3}$,

$$\frac{\|\hat{\nu} - \nu\|_F}{\|\nu\|_F} \le \frac{C}{\sqrt{np}}.\tag{22}$$

4 Numerical experiments

We evaluate the numerical performance of the KRO-PRO-FAC algorithm through a few simulation studies and real data analysis.

4.1 Simulation studies

For the simulation experiments we set the dimensions of Y_i and X_i to be $p_1 = p_2 = 500$ and $q_1 = q_2 = 2$ while the sample size n is chosen in {200, 400, 1000, 2000, 3000}. For ease of exposition we only consider the special case of Eq. (13) with d = 1, and thus $\nu = \beta_2 \otimes \beta_1$ where β_1 and β_2 are 500×2 matrices. We first generate $\operatorname{vec}(\beta_1)$ from the standard multivariate normal distribution on \Re^{1000} and similarly for $\operatorname{vec}(\beta_2)$; note that neither β_1 nor β_2 are expected to be sparse and furthermore the estimation of these β_k when n = 200 or n = 400 falls within the setting of regression with high-dimensional responses. We then generate X_1, X_2, \ldots, X_n where the $\operatorname{vec}(X_i)$ are iid standard multivariate normals in \Re^4 .

Given the $\{X_i\}$ we then consider the following 4 different models for the random noises $\{E_i\}$. The first three models corresponds to $\operatorname{vec}(E_i) \in \Re^{25000}$ that are multivariate normals while the last model corresponds to $\operatorname{vec}(E_i)$ with entries independently sampled from Student's t distribution with 5 degrees of freedom. The entries of E_i for Model 4 have heavier tails compared to that for Models 1–3.

- **Model** 1: Identity covariance: $\Sigma_{\text{vec}(E)} = \mathbf{I}_{p_1 p_2 \times p_1 p_2}$ and $\text{vec}(E_i)$'s are generated independently from the standard multivariate normal distribution.
- **Model** 2: Banded covariance: $\Sigma_{\text{vec}(E)} = \mathbf{L}\mathbf{L}^{\top}$ where \mathbf{L} is a lower triangular banded matrix in $\Re^{p_1p_2 \times p_2p_2}$ with $\mathbf{L}_{ij} = 0$ for i < j or i j > b. The bandwidth b is set to 5 and the diagonal elements are generated from $\mathcal{N}(3, 1)$ and the non-zero off-diagonal elements are generated from $\mathcal{N}(0, 1)$. \mathbf{L} is fixed over the 100 replications.

- **Model** 3: AR(1): $\Sigma_{\text{vec}(E)} = (\rho^{|i-j|})_{p_1p_2 \times p_1p_2}$ with $\rho = 0.9$. Here we generate $\text{vec}(E_1), \ldots, \text{vec}(E_n)$ based on Matlab codes from arima.
- **Model** 4: Heavy-tailed: $\Sigma_{\text{vec}(E)}$ is proportional to $\mathbf{I}_{p_1p_2 \times p_1p_2}$ and the entries of $\text{vec}(E_i)$'s are random samples from the Student's t-distribution with 5 degrees of freedom.

For each choice of the noise model for E_i we then generate $\{Y_i\}$ according to Eq. (3) and then estimate $\hat{\nu}$ based on the $\{X_i, Y_i\}$ using the KRO-PRO-FAC algorithm. For illustrative comparisons we considered, in addition to the default described in Algorithm 1, two other variants which performs rank regularization of either the responses or the OLS estimate. More specifically the first variant uses, instead of the observed Y_i , its truncated rank $-\alpha$ SVD $Y_i^{(\alpha)}$ for estimating ν . We termed this variant as KRO-PRO-FAC (α) and note that it is motivated by the fact that while Y_i is, with probability 1 full rank, $\mathbb{E}[Y_i] = \beta_1 X_i \beta_2^{\top}$ is low-rank for all i and thus a rank-regularized version of the $\{Y_i\}$ might lead to better estimate of ν . The second variant also performs rank regularization, but on the OLS estimate $\tilde{\nu}$ as opposed to the responses $\{Y_i\}$. Letting $\tilde{\nu}^{(\gamma)}$ be the truncated rank $-\gamma$ SVD of $\tilde{\nu}$ we then perform the remaining steps of Algorithm 1 with $\tilde{\nu}^{(\gamma)}$ in place of $\tilde{\nu}$. We termed this variant as rdu-rank-KRO (γ) and note that it is motivated by the notion of reduced-rank-regression in Izenman (1975). For this simulation we chose $\alpha = \gamma = 2$.

Finally we also estimate ν using the MLE based procedure described in Ding and Cook (2016). Recall that this MLE based approach posits both a Kronecer product structure for both the regression coefficient ν and the covariance matrix of $\text{vec}(E_i)$. We use the implementation from is based on R codes from MatrixEnv and denote the resulting estimates as dual-KRO-MLE. Table 1 summarizes some key differences between the 4 methods described above. For numerical comparisons we evaluate the relative errors $\|\hat{\nu} - \nu\|_F / \|\nu\|_F$ for each methods and *averaged* these over 100 independent Monte Carlo replicates. The results are presented in Table 2 through Table 5 for the four noise models described above.

For Model 1 we see from Table 2 that both the KRO-PRO-FAC and KRO-PRO-FAC (α) method have the smallest estimation error. The dual-KRO-MLE estimate is substantially less accurate compared to that of KRO-PRO-FAC and KRO-PRO-FAC (α) especially when the sample size is small, e.g., n = 200 or n = 400. This is in a sense expected as the entries of E_i are iid and thus there are few if any benefits in estimating and/or incorporating the covariance structure of $\{Y_i\}$. Finally, the rdu-rank-KRO (γ) method has the highest estimation error and this observation

method	data	ν estimation	Kronecker structure on	
KRO-PRO-FAC	Y_i	OLS	mean	
KRO-PRO-FAC (α)	rank $-\alpha Y_i^{(\alpha)}$	OLS	mean	
rdu-rank-KRO (γ)	Y_i	rank $-\gamma$ OLS $\tilde{\nu}^{(\gamma)}$	mean	
dual-KRO-MLE	Y_i	column & row separate estimates	mean & covariance	

Table 1: Method Comparison

Table 2: Average relative estimation error (%) under the identity covariance (Model 1)

	sample size (n)				
$\ \nu - \nu\ _F / \ \nu\ _F$	200	400	1000	2000	3000
KRO-PRO-FAC	0.339	0.237	0.151	0.106	0.087
KRO-PRO-FAC ($\alpha = 2$)	0.339	0.238	0.151	0.106	0.087
rdu-rank-KRO ($\gamma = 2$)	63.998	63.998	63.997	63.997	63.997
dual-KRO-MLE	76.938	52.702	20.721	8.377	6.563

also extends to the results for Model 2 through 4 as presented in Table 3 through Table 5 This is once again expected as, recalling the earlier discussions in Example 2.2, the Kronecker structure in the regression coefficient ν is fundamentally different from assuming ν to be low-rank. In other words imposing rank constraints on $\tilde{\nu}$ only leads to information loss due to model misspecification.

For Model 2 we see from Table 3 that the KRO-PRO-FAC algorithm has the smallest estimation error with the KRO-PRO-FAC (α) variant being slightly worse. The estimate obtained from the

Table 3: Average relative estimation error (%) under the banded covariance with a bandwidth $b = 5 \pmod{2}$

	sample size (n)					
$\ \nu - \nu\ _F / \ \nu\ _F$	200	400	1000	2000	3000	
KRO-PRO-FAC	1.508	1.059	0.666	0.472	0.385	
KRO-PRO-FAC ($\alpha = 2$)	1.552	1.116	0.746	0.576	0.506	
rdu-rank-KRO ($\gamma = 2$)	63.963	63.997	63.999	63.998	63.998	
dual-KRO-MLE	38.934	19.116	7.211	2.292	2.418	

	sample size (n)					
$\ \nu - \nu\ _F / \ \nu\ _F$	200	400	1000	2000	3000	
KRO-PRO-FAC	0.351	0.247	0.157	0.110	0.090	
KRO-PRO-FAC ($\alpha = 2$)	0.512	0.440	0.391	0.373	0.371	
rdu-rank-KRO $(\gamma=2)$	63.998	63.998	63.997	63.997	63.997	
dual-KRO-MLE	0.250	0.177	0.113	0.079	0.064	

Table 4: Average relative estimation error (%) under the AR(1) setting with $\rho = 0.9$ (Model 3)

Table 5: Average relative estimation error (%) under the Student's t-distribution with 5 degrees of freedom (Model 4)

	sample size (n)				
$\ \nu - \nu\ _F / \ \nu\ _F$	200	400	1000	2000	3000
KRO-PRO-FAC	0.439	0.308	0.194	0.137	0.112
KRO-PRO-FAC ($\alpha = 2$)	0.440	0.308	0.195	0.137	0.113
rdu-rank-KRO $(\gamma=2)$	63.998	63.998	63.997	63.997	63.997
dual-KRO-MLE	69.049	39.734	12.677	7.141	5.010

dual-KRO-MLE algorithm is noticably worse compared to both the KRO-PRO-FAC and KRO-PRO-FAC (α) and furthermore appeared to be sensitive to the sample size n, i.e., its estimation error is much larger than its competitors when n = 200 or n = 400. We note that for this Model 2, 90% of the non-zero correlations in $\Sigma_{vec(E)}$ have absolute value less than 0.5, which suggests either weak or mild dependence among rows and columns in $vec(Y_i)$.

For Model 3 we see from Table 4 that the dual-KRO-MLE algorithm yields the most accurate estimates with errors that are slightly smaller than that of KRO-PRO-FAC and KRO-PRO-FAC (α) methods. There is thus value in joint modeling of the mean ν and the covariance structure for the $\{Y_i\}$. Note, however, that the KRO-PRO-FAC algorithm is much less computationally demanding compared to dual-KRO-MLE.

Finally, for Model 4 we see from Table 5 that the KRO-PRO-FAC algorithm outperforms all of its competitors. In particular it is slightly better than KRO-PRO-FAC(α) and is much better than dual-KRO-MLE. These results are similar to that in Table 2 and one possible explanation for this similarity is that both models induce the same covariance structure for $\{Y_i\}$.

4.2 Real data analysis

We now apply the KRO-PRO-FAC algorithm to the electroencephacology (EEG) dataset from the UC Irvine Machine Learning Repository. The data arises from a study of EEG measurements related to alcoholoism in which there are 122 subjects from either the alcoholic group (77 subjects) or the control group (45 subjects). For each subject a series of voltage measurements is made at 256 different time points from 64 different regions of the scalp, i.e., the EEG response for the *i*th subject in the *j*th group (with j = 1 and j = 2 denoting the alcoholic and control) can be viewed as a matrix Y_{ij} with 256 rows and 64 columns. A key research question for this dataset is to identify which of the 64 brain channel accounts for most of the differences in voltages measurements between the two groups.

To answer the above inquiry we partition the data according to the subject grouping and fit a bi-linear model of the form Eq (13) to each group. As the EEG dataset contains no other covariates, this lead to a model of the form

$$\operatorname{vec}(Y_{ij}) = \left(\sum_{k=1}^{d^{(j)}} \beta_{2k}^{(j)} \otimes \beta_{1k}^{(j)}\right) + \operatorname{vec}(E_{ij}), \quad i \in [n_j], \ j = 1, 2$$
(23)

where $\beta_{2k}^{(j)} \in \Re^{64 \times 1}$ and $\beta_{1k}^{(j)} \in \Re^{256 \times 1}$. In other words, the mean response $\nu^{(j)} = \mathbb{E}[\operatorname{vec}(Y_{ij})]$ for the *j*th group is a sum of $d^{(j)}$ Kronecker products and thus $\nu^{(1)} - \nu^{(2)}$ is the effect of alcoholism (when compared to the control group) on the voltage measurements. We emphasize that the number of Kronecker factors $d^{(j)}$ are possibly different between the two groups. We apply the KRO-PRO-FAC algorithm to these $\{Y_{ij}\}$ with $d^{(1)} = 2$ and $d^{(2)} = 3$ chosen via the singular value ratio criterion as described in Eq. (15). Let $\hat{\nu}^{(1)} = \sum_{k=1}^{2} \hat{\beta}_{2k}^{(1)} \otimes \hat{\beta}_{1k}^{(1)}$ and $\hat{\nu}^{(2)} = \sum_{k=1}^{2} \hat{\beta}_{2k}^{(2)} \otimes \hat{\beta}_{1k}^{(2)}$ be the resulting estimates of $\nu^{(1)}$ and $\nu^{(2)}$.

Given these $\hat{\nu}^{(1)}$ and $\hat{\nu}^{(2)}$ we then follow the same post-processing steps described in Ding and Cook (2016) for multiple testing among the brain locations. Firsly, we isolate the alcoholism effects of each channel by averaging out the time effects $\hat{\nu}^{(1)} - \hat{\nu}^{(2)}$, ie., we take the *column* means of the $\operatorname{vec}^{-1}(\hat{\nu}^{(1)} - \hat{\nu}^{(2)}, 256, 64)$ where $\operatorname{vec}^{-1}(\cdot, 256, 64)$ yields a matrix of dimensions 256×64 . This yields in a vector $\hat{\theta} \in \Re^{64}$ which we then conduct multiple t-test for the null hypothesis that $\mathbb{H}_0: \theta_i = 0$ and compute the resulting p-values. Finally we apply the Benjamini–Yekutieli procedure (Benjamini and Yekutieli, 2001) to adjust these *p*-values.

The left panel of Figure 2 reports these adjusted p-values (on a \log_{10} scale). For comparisons we also repeat the same post-processing analysis but replaced the estimates $\hat{\nu}^{(1)}$ and $\hat{\nu}^{(2)}$ by the the OLS estimate $\bar{Y}^{(1)} = n_1^{-1} \sum_{i \in n_1} Y_{i1}$ and $\bar{Y}^{(2)} = n_2^{-1} \sum_{i \in n_2} Y_{i2}$ and present the adjusted *p*-values for these OLS estimates in the right panel of Figure 2. Figure 2 indicates that, for a significant level of 0.05, the KRO-PRO-FAC estimates lead to the detectation of 20 possibly relevant channels while the OLS estimates detect only 3 possibly relevant channels. We note that Li and Zhang (2017); Ding and Cook (2016) also analyzed the same data set and their estimates detect 24 and 26 possibly relevant channels, respectively. Our detections using the KRO-PRO-FAC estimates are thus comparable with those from Li and Zhang (2017); Ding and Cook (2016); indeed they all detected the regions from 21 to 25, from 44 to 52 and from 57 and 62. The main benefit of using the KRO-PRO-FAC estimates is that they can be commputed efficiently and do not depend on knowing or estimating $\text{Cov}[\text{vec}(Y_{ij})]$.

5 Conclusion

In this paper we consider matrix regression $Y_i = \sum_k \beta_{1k} X_i \beta_{2k}^\top + E_i$ where the responses Y_i are highdimensional matrices and propose a computationally efficient procedure for estimating $\{\beta_{1k}, \beta_{2k}\}$



Figure 2: Benjamini–Yekutieli adjusted p-values (on a scale of $-\log_{10}$) for 64 brain channels obtained from the (a) Kronecker products estimates in Eq. (23) with d = 2 and d = 3 for the alcoholic and control group (b) OLS estimates

based on the nearest Kronecker products approximation to the OLS estimate $\hat{\nu}$ of $\nu = \sum_k \beta_{2k} \otimes \beta_{1k}$. We now mention three potential directions for future research. The empirical results in Section 4.1 show that the KRO-PRO-FAC procedure has smallest estimation error when the noise entries for E_i are independent and is slightly worse than the dual-KRO-MLE procedure of Ding and Cook (2016) when the noise of E_i are highly correlated. As the dual-KRO-MLE is somewhat computationally demanding, it will be valuable to refine our KRO-PRO-FAC procedure for handling highly dependent rows and columns without compromising its computation efficiency. Secondly, the performance of the low-rank variant KRO-PRO-FAC (α) is also quite competitive but its theoretical property is currently unaddressed. Finally, for many type of matrix data such as those arising in image analysis, the ordering of the rows and columns for these matrices are based on latent but important features. For example, pixels' intensities in an image usually exhibit some continuity in both vertical and horizontal directions. How to meaningfully extract these latent features and incorporate them into the matrix regression problem is an open and interesting research question.

References

- Ahn, S. C. and Horenstein, A. R. (2013) Eigenvalue ratio test for the number of factors. *Econo*metrica, 81, 1203–1227.
- Barratt, S. (2018) A matrix Gaussian distribution. arXiv preprint arXiv:1804.11010.
- Benjamini, Y. and Yekutieli, D. (2001) The control of the false discovery rate in multiple testing under dependency. Annals of Statistics, 1165–1188.
- Beylkin, G. and Mohlenkamp, M. J. (2002) Numerical operator calculus in higher dimensions. Proceedings of the National Academy of Sciences, 99, 10246–10251.
- Bhatia, R. (2013) Matrix analysis. Springer.
- Bing, X. and Wegkamp, M. H. (2019) Adaptive estimation of the rank of the coefficient matrix in high-dimensional multivariate response regression models. *The Annals of Statistics*, 47, 3157– 3184.

- Bunea, F., She, Y. and Wegkamp, M. H. (2012) Joint variable and rank selection for parsimonious estimation of high-dimensional matrices. *The Annals of Statistics*, 40, 2359–2388.
- Cai, C., Chen, R. and Xiao, H. (2019) Kopa: Automated Kronecker product approximation. arXiv preprint arXiv:1912.02392.
- Chen, E. Y. and Fan, J. (2021) Statistical inference for high-dimensional matrix-variate factor models. *Journal of the American Statistical Association*, 1–18.
- Chen, E. Y., Tsay, R. S. and Chen, R. (2019) Constrained factor models for high-dimensional matrix-variate time series. *Journal of the American Statistical Association*.
- Chen, K., Chan, K.-S. and Stenseth, N. C. (2012) Reduced rank stochastic regression with a sparse singular value decomposition. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 74, 203–221.
- Crainiceanu, C. M., Caffo, B. S., Luo, S., Zipunnikov, V. M. and Punjabi, N. M. (2011) Population value decomposition, a framework for the analysis of image populations. *Journal of the American Statistical Association*, **106**, 775–790.
- De Lathauwer, L., De Moor, B. and Vandewalle, J. (2000) A multilinear singular value decomposition. SIAM Journal on Matrix Analysis and Applications, **21**, 1253–1278.
- De Waal, D. (1985) Matrix-valued distributions. In *Encyclopedia of statistical sciences*, 326–333.
 Wiley Online Library.
- Ding, S. and Cook, R. (2016) Matrix-variate regressions and envelope models. Journal of the Royal Statistical Society: Series B, 80.
- Dutilleul, P. (1999) The MLE algorithm for the matrix normal distribution. Journal of Statistical Computation and Simulation, 64, 105–123.
- Eckart, C. and Young, G. (1936) The approximation of one matrix by another of lower rank. Psychometrika, 1, 211–218.
- Feng, Y., Xiao, L. and Chi, E. C. (2021) Sparse single index models for multivariate responses. Journal of Computational and Graphical Statistics, 30, 115–124.

- Greenewald, K. and Hero, A. O. (2015) Robust Kronecker product PCA for spatio-temporal covariance estimation. *IEEE Transactions on Signal Processing*, **63**, 6368–6378.
- Gupta, A. K. and Nagar, D. K. (1999) Matrix variate distributions. Chapman and Hall/CRC.
- Halko, N., Martinsson, P.-G. and Tropp, J. A. (2011) Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53, 217–288.
- Izenman, A. J. (1975) Reduced-rank regression for the multivariate linear model. Journal of Multivariate Analysis, 5, 248–264.
- Kamm, J. and Nagy, J. G. (1998) Kronecker product and SVD approximations in image restoration. *Linear Algebra and its Applications*, 284, 177–192.
- Kong, D., An, B., Zhang, J. and Zhu, H. (2019) L2rm: Low-rank linear regression models for high-dimensional matrix responses. *Journal of the American Statistical Association*.
- Lam, C. and Yao, Q. (2012) Factor modeling for high-dimensional time series: inference for the number of factors. Annals of Statistics, 694–726.
- Li, L. and Zhang, X. (2017) Parsimonious tensor response regression. Journal of the American Statistical Association, 112, 1131–1146.
- Musco, C. and Musco, C. (2015) Randomized block krylov methods for stronger and faster approximate singular value decomposition. Advances in neural information processing systems, 28.
- Nagy, J. G. (1996) Decomposition of block Toeplitz matrices into a sum of Kronecker products with applications in image restoration. *Tech. rep.*, Southern Methodist University.
- Negahban, S. and Wainwright, M. J. (2011) Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, **39**, 1069–1097.
- Obozinski, G., Wainwright, M. J. and Jordan, M. I. (2011) Support union recovery in highdimensional multivariate regression. *The Annals of Statistics*, **39**, 1–47.

- Tropp, J. A., Yurtsever, A., Udell, M. and Cevher, V. (2017) Practical sketching algorithms for low-rank matrix approximation. SIAM Journal on Matrix Analysis and Applications, 38, 1454–1485.
- Tyrtyshnikov, E. (2004) Kronecker-product approximations for some function-related matrices. Linear Algebra and its Applications, **379**, 423–437.
- Van Loan, C. F. and Pitsianis, N. (1993) Approximation with Kronecker products. In *Linear algebra for large scale and real-time applications*, 293–314. Springer.
- Vershynin, R. (2018) High-Dimensional Probability: An Introduction with Applications in Data Science. Cambridge University Press.
- Vølund, A. (1980) Multivariate bioassay. *Biometrics*, 225–236.
- Wang, D., Liu, X. and Chen, R. (2019) Factor models for matrix-valued high-dimensional time series. *Journal of Econometrics*, 208, 231–248.
- Wang, D., Shen, H. and Truong, Y. (2016) Efficient dimension reduction for high-dimensional matrix-valued data. *Neurocomputing*, **190**, 25–34.
- Wang, D., Zheng, Y., Lian, H. and Li, G. (2021) High-dimensional vector autoregressive time series modeling via tensor decomposition. *Journal of the American Statistical Association*, 1–19.
- Wedin, P.-Å. (1972) Perturbation bounds in connection with singular value decomposition. BIT Numerical Mathematics, 12, 99–111.
- Werner, K., Jansson, M. and Stoica, P. (2008) On estimation of covariance matrices with kronecker product structure. *IEEE Transactions on Signal Processing*, 56, 478–491.
- Ye, J. (2005) Generalized low rank approximations of matrices. Machine Learning, 61, 167–191.
- Yuan, M., Ekici, A., Lu, Z. and Monteiro, R. (2007) Dimension reduction and coefficient estimation in multivariate linear regression. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 69, 329–346.
- Zhang, D. (2005) (2d)² pca: Two-directional two-dimensional pca for efficient face representation and recognition. *Neurocomputing*, **69**, 224–231.

- Zhao, J. and Leng, C. (2014) Structured Lasso for regression with matrix covariates. Statistica Sinica, 799–814.
- Zheng, Z., Bahadori, M. T., Liu, Y. and Lv, J. (2019) Scalable interpretable multi-response regression via seed. J. Mach. Learn. Res., 20, 1–34.
- Zou, C., Ke, Y. and Zhang, W. (2020) Estimation of low rank high-dimensional multivariate linear models for multi-response data. *Journal of the American Statistical Association*, 1–11.

Appendix A: Proofs of Stated Results

A.1 Proof of Theorem 1

For convenience of notation, we take

$$\underbrace{(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}^{\top}}_{q_{1}q_{2}\times n} = \begin{bmatrix} \mathcal{C}_{1}, \mathcal{C}_{2}, \dots, \mathcal{C}_{q_{1}q_{2}} \end{bmatrix}^{\top}, \ \mathcal{C}_{1}, \mathcal{C}_{2}, \dots, \mathcal{C}_{q_{1}q_{2}} \in \Re^{n}$$

$$\underbrace{\mathcal{E}}_{n \times p_{1}p_{2}} = \begin{bmatrix} \operatorname{vec}(E_{1})^{\top} \\ \vdots \\ \operatorname{vec}(E_{n})^{\top} \end{bmatrix} = \begin{bmatrix} \mathcal{D}_{1}, \mathcal{D}_{2}, \dots, \mathcal{D}_{p_{1}p_{2}} \end{bmatrix} = \begin{bmatrix} \mathcal{F}_{1}, \mathcal{F}_{2}, \dots, \mathcal{F}_{p_{1}p_{2}} \end{bmatrix} \Sigma_{\operatorname{vec}(E)}^{1/2}$$
(24)

where $\mathcal{D}_1, \mathcal{D}_2 \cdots \mathcal{D}_{p_1 p_2} \in \Re^n$ and $\mathcal{F}_1, \mathcal{F}_2 \cdots \mathcal{F}_{p_1 p_2} \in \Re^n$ with $\mathcal{F}_s = (\xi_{1s}, \dots, \xi_{ns})^\top$ being independent with independent entries.

Set $\Delta = \left[(\mathcal{X}^{\top} \mathcal{X})^{-1} \mathcal{X}^{\top} \mathcal{E}) \right]^{\top}$ and recall that $R(\nu) = \sum_{k=1}^{d} \operatorname{vec}(\beta_{2k}) \operatorname{vec}(\beta_{1k})^{\top} = \sum_{k=1}^{d} \sigma_k \, \mathcal{U}_k \mathcal{V}_k^{\top}$. We then have

$$R(\hat{\nu}) = \sum_{k=1}^{d} \sigma_k \, \mathcal{U}_k \mathcal{V}_k^{\top} + R(\Delta)$$

By Weyl's inequality (Problem III.6.13 in Bhatia (2013)) and Wedin sin- Θ (Wedin, 1972) theorem we have

$$\left|\hat{\sigma}_{k} - \sigma_{k}\right| \le \left\|R(\Delta)\right\| \tag{26}$$

$$\|\sin\Theta(\mathcal{U}_{1},\hat{\mathcal{U}}_{1})\|_{2}, \|\sin\Theta(\mathcal{V}_{1},\hat{\mathcal{V}}_{1})\|_{2} \leq \frac{\min\{\|\hat{\mathcal{U}}_{1}^{\top}R(\Delta)\|_{2}, \|R(\Delta)\hat{\mathcal{V}}_{1}\|_{2}\}}{\hat{\sigma}_{d}} \leq \frac{\|R(\Delta)\|_{2}}{\hat{\sigma}_{d}}$$
(27)

It thus suffices to bound the spectral norm of $R(\Delta)$. First note that $R(\Delta)$ can be written as a block matrix of the form

$$R(\Delta) = \begin{bmatrix} \tilde{\Delta}_{11} & \tilde{\Delta}_{12} & \dots & \tilde{\Delta}_{1q_1} \\ \tilde{\Delta}_{21} & \tilde{\Delta}_{22} & \dots & \tilde{\Delta}_{2q_1} \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{\Delta}_{q_21} & \tilde{\Delta}_{q_22} & \dots & \tilde{\Delta}_{q_2q_1} \end{bmatrix}, \quad \tilde{\Delta}_{k\ell} = \begin{bmatrix} \mathcal{D}_1^{\top} \mathcal{C}_{(k-1)q_1+\ell} & \dots & \mathcal{D}_{p_1}^{\top} \mathcal{C}_{(k-1)q_1+\ell} \\ \mathcal{D}_{p_1+1}^{\top} \mathcal{C}_{(k-1)q_1+\ell} & \dots & \mathcal{D}_{2p_1}^{\top} \mathcal{C}_{(k-1)q_1+\ell} \\ \vdots & \ddots & \vdots \\ \mathcal{D}_{(p_2-1)p_1+1}^{\top} \mathcal{C}_{(k-1)q_2+\ell} & \dots & \mathcal{D}_{p_1p_2}^{\top} \mathcal{C}_{(k-1)q_1+\ell} \end{bmatrix}$$

The matrix $\tilde{\Delta}_{k\ell}$ can be further expressed as

$$\tilde{\Delta}_{k\ell} = \begin{bmatrix} \boldsymbol{\Sigma}_1^{1/2} \zeta_1 & \boldsymbol{\Sigma}_2^{1/2} \zeta_2 & \dots & \boldsymbol{\Sigma}_{p_2}^{1/2} \zeta_{p_1} \end{bmatrix}, \quad \zeta_s = \begin{bmatrix} \mathcal{F}_s, \mathcal{F}_{p_1+s}, \mathcal{F}_{2p_1+s}, \dots, \mathcal{F}_{(p_2-1)p_1+s} \end{bmatrix}^\top \mathcal{C}_{(k-1)q_1+\ell}$$
(28)

Note that for ease of exposition we had suppressed the dependency on k and ℓ in the notation for ζ_s . This should cause minimal confusion as we can fix some arbitrary k and ℓ before proceeding with the subsequent derivations.

We now derive a concentration inequality for $\|\tilde{\Delta}_{k\ell}\|$ using a standard ϵ -net argument.

Step 1: ϵ **net** Let $\epsilon = 1/4$ and choose an ϵ net \mathcal{M} of the sphere \mathcal{S}^{p_2-1} and an ϵ net \mathcal{R} of the sphere \mathcal{S}^{p_1-1} . We have

$$|\mathcal{M}| \le 9^{p_2}, \quad |\mathcal{R}| \le 9^{p_1}$$

The spectral norm of $\tilde{\Delta}_{k\ell}$ can then be bounded as

$$\|\tilde{\Delta}_{k\ell}\|_2 \le 2 \max_{x \in \mathcal{M}, y \in \mathcal{R}} \langle \tilde{\Delta}_{k\ell} x, y \rangle$$

Step 2: Concentration Fix $x \in \mathcal{M}$ and $y \in \mathcal{R}$. We then have

$$\langle \tilde{\Delta}_{k\ell} x, y \rangle = \sum_{i=1}^{p_2} x_i \zeta_i^\top \Sigma_i^{1/2} y = \sum_{i=1}^{p_2} x_i \Big[\sum_{s=1}^{p_1} \sum_{k=1}^{p_1} \Sigma_i^{1/2} (k, s) \zeta_{ik} y_s \Big]$$

Using properties of the Orlicz Ψ_2 -norm (see e.g., Proposition 2.6.1 in Vershynin (2018)) we have

$$\begin{aligned} \|\langle \tilde{\Delta}_{k\ell} x, y \rangle \|_{\psi_{2}}^{2} &\leq \mathcal{C} \sum_{i=1}^{p_{2}} x_{i}^{2} \| \sum_{s=1}^{p_{1}} \sum_{k=1}^{p_{1}} \Sigma_{i}^{1/2}(k, s) \zeta_{ik} y_{s} \|_{\psi_{2}}^{2} \\ &\leq \mathcal{C} \sum_{i=1}^{p_{2}} x_{i}^{2} \sum_{s=1}^{p_{1}} y_{s}^{2} \sum_{k=1}^{p_{1}} \| \Sigma_{i}^{1/2}(k, s) \zeta_{ik} \|_{\psi_{2}}^{2} \\ &\leq \mathcal{C} \max_{i,k} \| \zeta_{ik} \|_{\psi_{2}}^{2} \sum_{i=1}^{p_{2}} x_{i}^{2} \sum_{s=1}^{p_{1}} y_{s}^{2} \sum_{k=1}^{p_{1}} | \Sigma_{i}^{1/2}(k, s) |^{2} \\ &\leq \mathcal{C} \max_{i,k} \| \zeta_{ik} \|_{\psi_{2}}^{2} \max_{i \in [p_{2}]} \max_{s \in [p_{1}]} \Sigma_{i}(s, s) \\ &\leq \mathcal{C} \| (\mathcal{X}^{\top} \mathcal{X})^{-1} \mathcal{X} \|_{2} \max_{i \in [m]} \max_{s \in [r]} \Sigma_{i}(s, s) \end{aligned}$$

where the second to last inequality is because

$$\sum_{k=1}^{p_1} \left| \boldsymbol{\Sigma}_i^{1/2}(k,s) \right|^2 = \sum_{k=1}^{p_1} \boldsymbol{\Sigma}_i^{1/2}(s,k) \boldsymbol{\Sigma}_i^{1/2}(k,s) = \boldsymbol{\Sigma}_i(s,s).$$

Let $K = \|(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}\| \times \max_{i \in [m]} \max_{s \in [r]} \Sigma_i(s, s)$. We therefore have, for all $u \ge 0$, that

$$P(\langle \tilde{\Delta}_{k\ell} x, y \rangle \ge u) \le 2 \exp\left(-cu^2/K^2\right).$$
⁽²⁹⁾

Step 3: Union bound By union over the \mathcal{M} and \mathcal{R} , then with probability $1 - 2 \exp(-u^2)$, we have for any u > 0

$$\|\tilde{\Delta}_{k\ell}\| \leq \mathcal{C}\|(\mathcal{X}^{\top}\mathcal{X})^{-1}\mathcal{X}\|_{2} \max_{i \in [p_{2}]} \max_{s \in [p_{1}]} \Sigma_{i}(s,s)(\sqrt{p_{2}} + \sqrt{p_{1}} + u).$$
(30)

This upper bound is independent of $\{k, \ell\}$ and since $||R(\Delta)|| \leq \sum_{k=1}^{q_1} \sum_{\ell=1}^{q_2} ||\tilde{\Delta}_{k\ell}||$, we obtain the desired results in Theorem 1.

A.2 Proof of Corollary 1

We will continue to use the same notations as that in the proof of Theorem 1. Let $\mathcal{P}_0 = \mathcal{U}\mathcal{U}^{\top}$ and $\mathcal{P}_1 = \mathcal{V}\mathcal{V}^{\top}$. Similarly, let $\hat{\mathcal{P}}_0 = \hat{\mathcal{U}}\hat{\mathcal{U}}^{\top}$ and $\hat{\mathcal{P}}_1 = \hat{\mathcal{V}}\hat{\mathcal{V}}^{\top}$. Note that these matrices are all of rank at most d. As $R(\nu) = \mathcal{P}_0 R(\nu) \mathcal{P}_1$, we have

$$\hat{\mathcal{P}}_0 R(\tilde{\nu}) \hat{\mathcal{P}}_1 - R(\nu) = (\hat{\mathcal{P}}_0 - \mathcal{P}_0) R(\tilde{\nu}) \hat{\mathcal{P}}_1 + \mathcal{P}_0 R(\tilde{\nu}) (\hat{\mathcal{P}}_1 - \mathcal{P}_1) + \mathcal{P}_0 R(\Delta) \mathcal{P}_1.$$

We therefore have

$$\|\hat{\mathcal{P}}_0 R(\tilde{\nu})\hat{\mathcal{P}}_1 - R(\nu)\|_F \le \sqrt{d} \|\hat{\mathcal{P}}_0 - \mathcal{P}_0\|_2 \times \|R(\tilde{\nu})\|_2 + \sqrt{d} \|\hat{\mathcal{P}}_1 - \mathcal{P}_1\|_2 \times \|R(\tilde{\nu})\|_2 + \sqrt{d} \|R(\Delta)\|_2.$$

Now $\|\hat{\mathcal{P}}_0 - \mathcal{P}_0\|_2 \leq 2\|\sin\Theta(\hat{\mathcal{U}},\mathcal{U})\|_2$ and similarly for $\|\hat{\mathcal{P}}_1 - \mathcal{P}_1\|_2$. Then from the conditions in Assumption 1, we have

$$\|\hat{\mathcal{P}}_0 R(\tilde{\nu})\hat{\mathcal{P}}_1 - R(\nu)\|_F \le 2\sqrt{d}\|\sin\Theta(\hat{\mathcal{U}},\mathcal{U})\|_2 \times (\|R(\nu)\|_2 + \|R(\Delta)\|_2) + \|R(\Delta)\|_2 = \mathcal{O}(n^{-1/2}p^{1/2})$$

Finally, as $\hat{\nu}$ and ν are the *inverse* rearrangement of $\hat{\mathcal{P}}_0 R(\tilde{\nu})\hat{\mathcal{P}}_1$ and $R(\nu)$, respectively, we have

$$\frac{\|\hat{\nu} - \nu\|_F}{\|\nu\|_F} = \frac{\|\hat{\mathcal{P}}_0 R(\tilde{\nu})\hat{\mathcal{P}}_1 - R(\nu)\|_F}{\|\nu\|_F} = \mathcal{O}((np)^{-1/2})$$

as desired.