

Supplementary Materials for “Likelihood-based Dimension Folding on Tensor Data”

Ning Wang, Xin Zhang

Department of Statistics, Florida State University, Tallahassee, FL, 32306

and Bing Li

Department of Statistics, Pennsylvania State University, University Park, PA 16802

S.1 Primary biliary cirrhosis data analysis

In this section, we analyze the data from a Mayo Clinical trial on primary biliary cirrhosis (PBC). This data set is also used in Xue & Yin (2014) and Sheng & Yuan (2019). The data set contains laboratory results for 312 patients. Several biomarkers such as bilirubin, albumin and prothrombin time are adopted to diagnose PBC. We study this data set in the same way as Xue & Yin (2014) did. We treat the time baseline as one fold of the covariates and the multivariate predictors repeatedly measured over time as another fold of the covariates and thus we have predictors in a matrix form. We focus our attention on the measurements of bilirubin, albumin level, and prothrombin time at time points 6 months, 1 year, 2 years, and 3 years and thus the predictors are in the form of 3×4 matrix. The response is the survival time of a patient which is continuous. We partition the response variable into 8 slices, and treat each slice as one class, then apply FLAD and FELAD to get the estimation of the dimension folding subspace. To make a comparison with Xue & Yin (2014) and Sheng & Yuan (2019), we use $d_1 = d_2 = 1$.

The estimated basis matrices Γ_1 and Γ_2 are $(0.151, -0.989, 0.007)^T$ and $(0.274, 0.240, -0.009, 0.931)^T$ by FLAD, and $(0.025, -0.999, 0.023)^T$ and $(0.094, 0.469, -0.761, 0.438)^T$ by FELAD. To test the sig-

		$\Gamma_{1,11}$	$\Gamma_{1,21}$	$\Gamma_{1,31}$	$\Gamma_{2,11}$	$\Gamma_{2,21}$	$\Gamma_{2,31}$	$\Gamma_{2,41}$
FLAD	Est.	0.151	-0.989	0.007	0.274	0.240	-0.009	0.931
	pvalue	0.002	0.006	0.904	0.484	0.216	0.736	0
FELAD	Est.	0.025	-0.999	0.023	0.094	0.469	-0.761	0.438
	pvalue	0.712	0.020	0.964	0.906	0.232	0.064	0.420
FSIR	Est.	0.018	-0.227	-0.974	0.623	0.204	0.727	0.208
	pvalue	0.002	0	0	0.606	0.596	0.598	0.616
FDR	Est.	0.057	-0.152	-0.987	-0.086	-0.241	-0.153	0.955
	pvalue	0.490	0	0.594	0.980	0.052	0.896	0.266
DCOV	Est.	0.989	-0.144	0.036	0.390	0.401	0.495	0.665
	pvalue	0	0	0.056	0	0	0	0
FMAVE	Est.	0.119	-0.989	-0.087	0.491	0.1811	0.159	0.837
	pvalue	0	0	0.366	0.170	0.200	0.586	0.146

Table S1: Bootstrapped p-values of coefficients estimated by different methods with bootstrap sample size 1000.

nificance of the coefficients, we compute the bootstrapped p-value for each element of Γ_1 and Γ_2 . Table S1 indicates that, for FLAD, the bilirubin and albumin levels at year 3 significantly affects the length of the survival time at level 0.05, and for FELAD, the albumin level significantly affects the length of the survival time at level 0.05. By Table S1, all the methods indicate that albumin significantly affects the length of survival time at level 0.05, and the results of FLAD, FSIR, DCOV and FMAVE show that bilirubin significantly affects the survival time. DCOV is the most stable method: the subspace estimation does not change so much with respect to different bootstrapped samples.

In Figure S1, we show the smoothing spline plots for the logarithm of the response versus the projected predictor, and scatter plots for the fitted value versus residuals. The results of FLAD and FELAD are

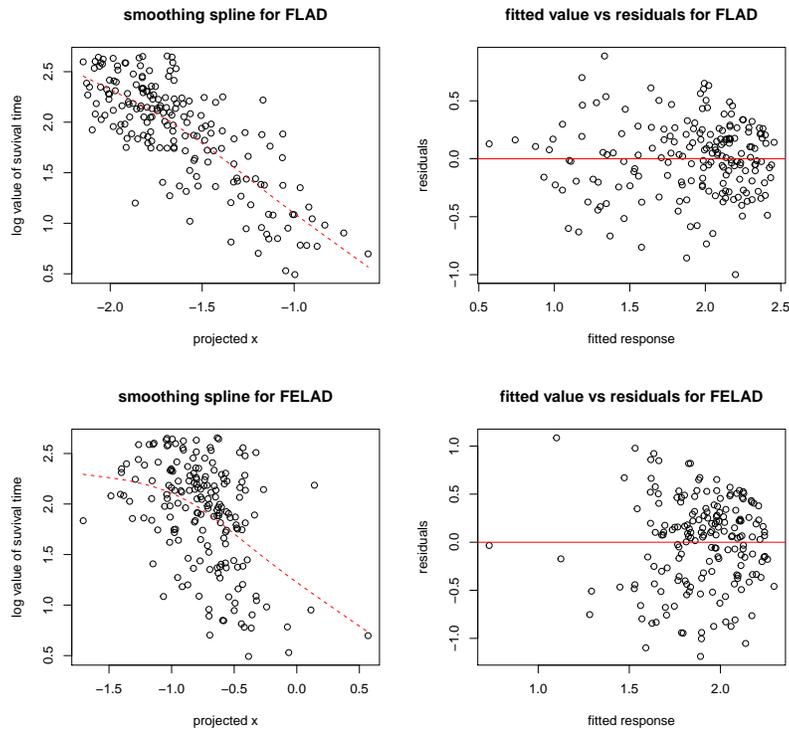


Figure S1: Left panels: Smoothing splines for the logarithm of response versus the projected predictor. Right panels: plots for the fitted values versus the residuals. The smoothing parameter is 0.076 for FLAD, and 0.033 for FELAD.

clear. It shows that the survival time has a negative relation with the projected predictor. Also, our analysis indicates that the survival time has a negative relationship with bilirubin level and a positive relationship with the albumin level, which is consistent with the medical outcome. The results are similar to those given in Xue & Yin (2014).

S.2 Proofs

S.2.1 Some Technical Lemmas

We first provide some lemmas that will be used in the proofs. The first three lemmas can be found in Appendix of Cook & Forzani (2009). Let \mathbf{B} be a symmetric positive definite matrix, and at $(\boldsymbol{\alpha}, \boldsymbol{\alpha}_0) \in \mathbb{R}^{p \times p}$ be a full rank matrix with $\boldsymbol{\alpha}^T \boldsymbol{\alpha}_0 = 0$.

Lemma 1.

$$\boldsymbol{\alpha}(\boldsymbol{\alpha}^T \mathbf{B} \boldsymbol{\alpha})^{-1} \boldsymbol{\alpha}^T + \mathbf{B}^{-1} \boldsymbol{\alpha}_0 (\boldsymbol{\alpha}_0^T \mathbf{B}^{-1} \boldsymbol{\alpha}_0)^{-1} \boldsymbol{\alpha}_0^T \mathbf{B}^{-1} = \mathbf{B}^{-1}. \quad (1)$$

As a consequence of Lemma 1, $\mathbf{I}_p - \mathbf{P}_{\boldsymbol{\alpha}(\mathbf{B})}^T = \mathbf{P}_{\boldsymbol{\alpha}_0(\mathbf{B}^{-1})}$. Additionally, if $(\boldsymbol{\alpha}, \boldsymbol{\alpha}_0)$ is an orthogonal matrix, then we have the following two lemmas.

Lemma 2.

$$|\boldsymbol{\alpha}_0^T \mathbf{B} \boldsymbol{\alpha}_0| = |\mathbf{B}| |\boldsymbol{\alpha}^T \mathbf{B}^{-1} \boldsymbol{\alpha}|. \quad (2)$$

Lemma 3.

$$\begin{aligned} (\boldsymbol{\alpha}_0^T \mathbf{B}^{-1} \boldsymbol{\alpha}_0)^{-1} &= \boldsymbol{\alpha}_0^T \mathbf{B} \boldsymbol{\alpha}_0 - \boldsymbol{\alpha}_0^T \mathbf{B} \boldsymbol{\alpha} (\boldsymbol{\alpha}^T \mathbf{B} \boldsymbol{\alpha})^{-1} \boldsymbol{\alpha}^T \mathbf{B} \boldsymbol{\alpha}_0, \\ -(\boldsymbol{\alpha}_0^T \mathbf{B}^{-1} \boldsymbol{\alpha}_0)^{-1} (\boldsymbol{\alpha}_0^T \mathbf{B}^{-1} \boldsymbol{\alpha}) &= (\boldsymbol{\alpha}_0^T \mathbf{B} \boldsymbol{\alpha}) (\boldsymbol{\alpha}^T \mathbf{B} \boldsymbol{\alpha})^{-1}. \end{aligned} \quad (3)$$

We will also use the following two lemmas about matrix derivatives.

Lemma 4. *Suppose $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{X}$ are full rank matrix, \mathbf{D} is full rank symmetric matrix. Then we have:*

$$\frac{d(\text{tr}((\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})\mathbf{D}^{-1}(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})^T))}{d\mathbf{X}} = 2\mathbf{A}^T(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})\mathbf{D}\mathbf{B}^T. \quad (4)$$

Lemma 5. *Suppose \mathbf{A}, \mathbf{B} are full rank matrices, and \mathbf{X} is symmetric matrix. Then we have:*

$$\frac{d(\text{tr}(\mathbf{A}\mathbf{X}^{-1}\mathbf{B}))}{d\mathbf{X}} = -(\mathbf{X}^{-1}\mathbf{B}\mathbf{A}\mathbf{X}^{-1})^T. \quad (5)$$

S.2.2 Proof of Proposition 1

Proof. Let Γ_m be the basis matrix for \mathcal{S}_m , $\Gamma = \bigotimes_{m=M}^1 \Gamma_m$, (Γ, Γ_0) be a full rank matrix with $\Gamma^T \Gamma_0 = 0$, and $\mathbf{P}_{\Gamma(\Sigma_k)} = \Gamma(\Gamma^T \Sigma_k \Gamma)^{-1} \Gamma^T \Sigma_k$. It is easy to show that $Y \mid \mathbf{X} \sim Y \mid (\bigotimes_{m=M}^1 \mathbf{P}_{\mathcal{S}_m}) \text{vec}(\mathbf{X})$ is equivalent to $\mathbf{X} \mid (\Gamma^T \mathbf{X}, Y) \sim \mathbf{X} \mid \Gamma^T \mathbf{X}$. By definition, $\text{cov}(\text{vec}(\mathbf{X}) \mid \Gamma^T \mathbf{X}, Y = k) = (\mathbf{I}_{\prod_{m=1}^M p_m} - \mathbf{P}_{\Gamma(\Sigma_k)}^T) \Sigma_k$ and $\text{E}(\text{vec}(\mathbf{X}) \mid \Gamma^T \text{vec}(\mathbf{X}), Y = k) = \text{vec}(\boldsymbol{\mu}) + \boldsymbol{\eta} \boldsymbol{\omega}_k + \mathbf{P}_{\Gamma(\Sigma_k)}^T (\text{vec}(\mathbf{X}) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\eta} \boldsymbol{\omega}_k)$, where $\boldsymbol{\omega}_k$ is the vector in $\mathbb{R}^{\dim(\mathcal{M})}$, and $\boldsymbol{\eta}$ is a basis matrix for \mathcal{M} . It is easy to show that $\text{var}(\text{vec}(\mathbf{X}) \mid \Gamma^T \mathbf{X}, Y = k)$ and $\text{E}(\text{vec}(\mathbf{X}) \mid \Gamma^T \text{vec}(\mathbf{X}), Y = k)$ are constant if and only if $(\mathbf{I}_{\prod_{m=1}^M p_m} - \mathbf{P}_{\Gamma(\Sigma_k)}^T) \Sigma_k$ and $\mathbf{P}_{\Gamma(\Sigma_k)}^T$ are constant and $\mathbf{P}_{\Gamma(\Sigma_k)}^T \boldsymbol{\eta} = \boldsymbol{\eta}$.

By Lemma 1, $(\mathbf{I}_{\prod_{m=1}^M p_m} - \mathbf{P}_{\Gamma(\Sigma_k)}^T) \Sigma_k = \Gamma_0 (\Gamma_0^T \Sigma_k^{-1} \Gamma_0)^{-1} \Gamma_0^T = C_1$ for some constant C_1 . Because $\mathbf{P}_{\Gamma(\Sigma_k)}^T$, $k = 1, \dots, K$, are constant, we know that $\mathbf{I} - \mathbf{P}_{\Gamma(\Sigma_k)}^T = \Gamma_0 (\Gamma_0^T \Sigma_k^{-1} \Gamma_0)^{-1} \Gamma_0^T \Sigma_k^{-1} = C_2$ for some constant C_2 . So $\Gamma_0 \Sigma_k^{-1}$ is a constant, which implies that $\mathbf{Q}_{\bigotimes_{m=M}^1 \mathcal{S}_m} \Sigma_k^{-1}$ does not change with k . Because $\mathbf{P}_{\Gamma(\Sigma_k)}^T$ is a constant, $\mathbf{P}_{\Gamma(\Sigma_k)}^T \boldsymbol{\eta} = \mathbf{P}_{\Gamma(\Sigma)}^T \boldsymbol{\eta} = \boldsymbol{\eta}$ for any k . And $\mathbf{P}_{\Gamma(\Sigma)}^T \boldsymbol{\eta} = \boldsymbol{\eta}$ if and only if $\mathbf{P}_{\Gamma(\Sigma)} (\Sigma^{-1} \boldsymbol{\eta}) = \Sigma^{-1} \boldsymbol{\eta}$. So we have $\Sigma^{-1} \boldsymbol{\eta} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$, which is equivalent to $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$. \square

S.2.3 Proof of Proposition 2

Proof. From the proof of Proposition 1, condition “ $\mathbf{Q}_{\bigotimes_{m=M}^1 \mathcal{S}_m} \Sigma_k^{-1}$ is a constant” of Proposition 1 is equivalent to $\mathbf{P}_{\Gamma(\Sigma_k)} = \mathbf{P}_{\Gamma(\Sigma)}$ and $\Sigma_k (\mathbf{I}_{\prod_{m=1}^M p_m} - \mathbf{P}_{\Gamma(\Sigma_k)}) = \Sigma (\mathbf{I}_{\prod_{m=1}^M p_m} - \mathbf{P}_{\Gamma(\Sigma)})$, which is equivalent to $\Sigma_k - \Sigma = \mathbf{P}_{\Gamma(\Sigma)}^T (\Sigma_k - \Sigma) \mathbf{P}_{\Gamma(\Sigma)}$. Also $\Sigma_k - \Sigma = \mathbf{P}_{\Gamma(\Sigma)}^T (\Sigma_k - \Sigma) \mathbf{P}_{\Gamma(\Sigma)}$ is equivalent to $\Sigma^{-1} (\Sigma - \Sigma_k) \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$. See Cook & Forzani (2009) for the proof of the equivalences.

By Proposition 1, we have $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$ and $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$, so $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$. Also, from Proposition 1 and the previous equivalences, we have $\Sigma^{-1} (\Sigma_k - \Sigma) \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$ and $\Sigma^{-1} (\Sigma_k - \Sigma) \subseteq \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$, and hence $\Sigma^{-1} (\Sigma_k - \Sigma) \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$. If we can show $\bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m = \bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m)$, then Proposition 2 is

proved.

We first show $\bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m) \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$. For $m = 1, \dots, M$, and any $\delta_m \in \mathcal{S}_m \cap \tilde{\mathcal{S}}_m$, we have $\bigotimes_{m=M}^1 \delta_m \in \bigotimes_{m=M}^1 \mathcal{S}_m$ and $\bigotimes_{m=M}^1 \delta_m \in \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$. So $\bigotimes_{m=M}^1 \delta_m \in \bigotimes_{m=M}^1 \mathcal{S} \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}$. Then we know that $\bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m) \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$.

Next, we show $\bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m \subseteq \bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m)$. For any vector $\nu \in \bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$, there exist $\alpha_m \in \mathcal{S}_m$ and $\beta_m \in \tilde{\mathcal{S}}_m$, $m = 1, \dots, M$, such that $\nu = \bigotimes_{m=M}^1 \alpha_m = \bigotimes_{m=M}^1 \beta_m$. It can be shown that $\bigotimes_{m=M}^1 \alpha_m = \bigotimes_{m=M}^1 \beta_m$ if and only if $\alpha_m = c_m \beta_m$ for some constants c_m . This implies that $\alpha_m \in \tilde{\mathcal{S}}_m$ and hence $\alpha_m \in \mathcal{S}_m \cap \tilde{\mathcal{S}}_m$. Therefore, $\nu = \bigotimes_{m=M}^1 \alpha_m \in \bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m)$ and $\bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m \subseteq \bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m)$. \square

S.2.4 Proof of Proposition 3

Proof. Let $\bigotimes_{m=M}^1 \Gamma_m$ be the basis matrix for $\bigotimes_{m=M}^1 \mathcal{S}_m$. If \mathcal{S}_m is a mode- m dimension folding envelope subspace, it must be a dimension folding subspace. From Proposition 1, \mathcal{S}_m is a mode- m dimension folding subspace if (a) $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$ and (b) $\mathbf{Q}_{\bigotimes_{m=M}^1 \mathcal{S}_m} \Sigma_k^{-1}$ is a constant with respect to k . The envelope covariance structure naturally satisfies (b). So $\Sigma^{-1} \mathcal{M} \subseteq \bigotimes_{m=M}^1 \mathcal{S}_m$ and $\Sigma_k = (\bigotimes_{m=M}^1 \mathbf{P}_{\mathcal{S}_m}) \Sigma_k (\bigotimes_{m=M}^1 \mathbf{P}_{\mathcal{S}_m}) + \mathbf{Q}_{\bigotimes_{m=M}^1 \mathcal{S}_m} \Sigma \mathbf{Q}_{\bigotimes_{m=M}^1 \mathcal{S}_m}$ make \mathcal{S}_m to be a mode- m dimension folding envelope subspace. \square

S.2.5 Proof of Proposition 4

Proof. From the proof of Proposition 2 in Wang et al. (2019), we have $\bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m$ satisfies the two conditions in our Proposition 3. We also proved that $\bigotimes_{m=M}^1 \mathcal{S}_m \cap \bigotimes_{m=M}^1 \tilde{\mathcal{S}}_m = \bigotimes_{m=M}^1 (\mathcal{S}_m \cap \tilde{\mathcal{S}}_m)$ in our Proposition 2. So the intersection of two mode- m dimension folding envelope subspace is a mode- m dimension folding envelope subspace. \square

S.2.6 Proof of Proposition 5

Proof. Let $\mathbf{\Gamma} = \bigotimes_{m=1}^M \mathbf{\Gamma}_m$ be the basis matrix for $\mathcal{T}_{Y|\mathbf{X}}$, and $(\mathbf{\Gamma}, \mathbf{\Gamma}_0)$ be a full rank matrix with $\mathbf{\Gamma}^T \mathbf{\Gamma}_0 = 0$. By Proposition 1, $\text{vec}(\boldsymbol{\mu} - \boldsymbol{\mu}_k) \subseteq \boldsymbol{\Sigma} \mathcal{T}_{Y|\mathbf{X}}$. So we have $\text{vec}(\boldsymbol{\mu} - \boldsymbol{\mu}_k) = \boldsymbol{\Sigma} \mathbf{\Gamma} \boldsymbol{\nu}_k$ for some $\boldsymbol{\nu}_k \in \mathbb{R}^{\prod_{m=1}^M d_m}$. By definition, we have $\mathbf{\Gamma}_0 \text{vec}(\mathbf{X}) \mid (\mathbf{\Gamma}^T \text{vec}(\mathbf{X}), Y = k) \sim N(\rho_k, \boldsymbol{\Theta}_k)$, with $\rho_k = \mathbf{\Gamma}_0^T \text{vec}(\boldsymbol{\mu}) + \mathbf{\Gamma}_0^T \boldsymbol{\Sigma} \mathbf{\Gamma} \boldsymbol{\nu}_k + (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})(\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} \mathbf{\Gamma}^T (\text{vec}(\mathbf{X}) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\Sigma} \mathbf{\Gamma} \boldsymbol{\nu}_k)$, and $\boldsymbol{\Theta}_k = \mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k \mathbf{\Gamma}_0 - \mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k \mathbf{\Gamma} (\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} \mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma}_0$. By Lemma 3, we have $\boldsymbol{\Theta}_k = (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k^{-1} \mathbf{\Gamma}_0)^{-1}$. Since, by Proposition 1, $\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k^{-1}$ is a constant, we have $\boldsymbol{\Theta}_k = (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0)^{-1} = \mathbf{D}$. Again by Lemma 3 and Proposition 1, we have $(\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}_k^{-1} \mathbf{\Gamma})(\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} = (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma})(\mathbf{\Gamma}^T \boldsymbol{\Sigma} \mathbf{\Gamma}) = \mathbf{H}$. Hence $\rho_k = \mathbf{H} \mathbf{\Gamma}^T \text{vec}(\mathbf{X}) + (\mathbf{\Gamma}_0^T - \mathbf{H} \mathbf{\Gamma}^T) \text{vec}(\boldsymbol{\mu})$. \square

S.2.7 Proof of Proposition 6

Proof. By the equality $f(\mathbf{\Gamma}^T \text{vec}(\mathbf{X}), \mathbf{\Gamma}_0^T \text{vec}(\mathbf{X}) \mid Y = k) = f(\mathbf{\Gamma}^T \text{vec}(\mathbf{X}) \mid Y = k) f(\mathbf{\Gamma}_0^T \text{vec}(\mathbf{X}) \mid \mathbf{\Gamma}^T \text{vec}(\mathbf{X}), Y = k)$ and Proposition 2, the log-likelihood function is

$$\begin{aligned}
L(\mathbf{\Gamma}) &= \sum_{i=1}^n \log f(\mathbf{\Gamma}^T \text{vec}(\mathbf{X}_i), \mathbf{\Gamma}_0^T \text{vec}(\mathbf{X}_i) \mid Y_i = k) \\
&= \sum_{i=1}^n \log f(\mathbf{\Gamma}^T \text{vec}(\mathbf{X}_i) \mid Y_i = k) + \sum_{i=1}^n \log f(\mathbf{\Gamma}_0^T \text{vec}(\mathbf{X}_i) \mid \mathbf{\Gamma}^T \text{vec}(\mathbf{X}_i), Y_i = k) \\
&= -\frac{npq}{2} \log(2\pi) - \frac{n}{2} \log |\mathbf{D}| - \frac{1}{2} \sum_k n_k \log |\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma}| \\
&\quad - \frac{1}{2} \sum_k n_k (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\Sigma} \mathbf{\Gamma} \boldsymbol{\nu}_k)^T \mathbf{\Gamma}^T (\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} \mathbf{\Gamma} (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\Sigma} \mathbf{\Gamma} \boldsymbol{\nu}_k) \\
&\quad - \frac{1}{2} \sum_k n_k (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}))^T \mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu})) \\
&\quad - \frac{1}{2} \sum_k n_k \text{tr} \{ \mathbf{\Gamma}^T \tilde{\boldsymbol{\Sigma}}_k \mathbf{\Gamma} (\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} \} - \frac{1}{2} \sum_k n_k \text{tr} (\mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T \tilde{\boldsymbol{\Sigma}}_k),
\end{aligned}$$

where $\tilde{\Sigma}_k$ is the sample covariance matrix of \mathbf{X}_k , and $\mathbf{K} = (\mathbf{\Gamma}_0 - \mathbf{\Gamma}\mathbf{H}^T)$. Let $f_k = n_k/n$. The only term in $L(\mathbf{\Gamma})$ that involves $\boldsymbol{\nu}_k$ is

$$T = \sum_k n_k (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\Sigma}\boldsymbol{\Gamma}\boldsymbol{\nu}_k)^T \mathbf{\Gamma}^T (\mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma})^{-1} \mathbf{\Gamma} (\text{vec}(\bar{\mathbf{X}}_y) - \text{vec}(\boldsymbol{\mu}) - \boldsymbol{\Sigma}\boldsymbol{\Gamma}\boldsymbol{\nu}_k).$$

We need to minimize it subject to $\sum_k f_k \boldsymbol{\nu}_k = 0$. Let $\mathbf{B}_k = \mathbf{\Gamma}^T \boldsymbol{\Sigma}_k \mathbf{\Gamma}$, and $\mathbf{Z}_k = \mathbf{\Gamma}^T (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}))$.

For any quantity a_k , let $\bar{a} = \sum_k f_k a_k$. We use the Lagrange multiplier method to minimize $T/n = \sum_k f_k (\mathbf{Z}_k - \bar{\mathbf{B}}\boldsymbol{\nu}_k)^T \mathbf{B}_k^{-1} (\mathbf{Z}_k - \bar{\mathbf{B}}\boldsymbol{\nu}_k) + \lambda \bar{\nu}$. Differentiating it with respect to $\boldsymbol{\nu}_k$, we get

$$-2f_k \mathbf{Z}_k + 2f_k \bar{\mathbf{B}}\boldsymbol{\nu}_k + f_k \mathbf{B}_k \bar{\mathbf{B}}^{-1} \lambda = 0. \quad (6)$$

By summation over the classes k of the above equation, we get $-2\bar{\mathbf{Z}} + \lambda = 0$. Finally, substituting it back into (6), we have $\boldsymbol{\nu}_k = \bar{\mathbf{B}}^{-1} (\mathbf{Z}_k - \mathbf{B}_k \bar{\mathbf{B}}^{-1} \bar{\mathbf{Z}})$. Therefore

$$\begin{aligned} T &= n (\mathbf{\Gamma}^T \text{vec}(\bar{\mathbf{X}}) - \mathbf{\Gamma}^T \text{vec}(\boldsymbol{\mu}))^T \bar{\mathbf{B}}^{-1} (\mathbf{\Gamma}^T \text{vec}(\bar{\mathbf{X}}) - \mathbf{\Gamma}^T \text{vec}(\boldsymbol{\mu})) \\ &= n (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu}))^T \mathbf{\Gamma} \bar{\mathbf{B}}^{-1} \mathbf{\Gamma}^T (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu})). \end{aligned}$$

Notice that, in $L(\mathbf{\Gamma})$, the term

$$\begin{aligned} &\frac{1}{2} \sum_k n_k (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu}))^T \mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T (\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\boldsymbol{\mu})) \\ &= \frac{1}{2} \sum_k n_k (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu}))^T \mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu})) \\ &\quad + \frac{1}{2} \sum_k n_k \text{tr}(\mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T ((\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\bar{\mathbf{X}}))(\text{vec}(\bar{\mathbf{X}}_k) - \text{vec}(\bar{\mathbf{X}}))^T)). \end{aligned}$$

The only two terms in $L(\mathbf{\Gamma})$ that involve $\text{vec}(\boldsymbol{\mu})$ are one from $\frac{1}{2}T$ and the other one from above equation. Hence we can solve the equation $\frac{\partial L}{\partial \text{vec}(\boldsymbol{\mu})} = 0$, and obtain $\mathbf{\Gamma} \bar{\mathbf{B}}^{-1} \mathbf{\Gamma}^T (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu})) + \mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu})) = 0$, where $\bar{\mathbf{B}} = \mathbf{\Gamma}^T \boldsymbol{\Sigma} \mathbf{\Gamma}$. By Lemma 1, $\mathbf{K}^T = \mathbf{\Gamma}_0^T (\mathbf{I} - \boldsymbol{\Sigma} \mathbf{\Gamma} (\mathbf{\Gamma}^T \boldsymbol{\Sigma} \mathbf{\Gamma})^{-1} \mathbf{\Gamma}^T) = (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0)^{-1} \mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1}$, and $\mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T = \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0 (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0)^{-1} \mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1}$. Again, using Lemma 1, $\boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0 (\mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{\Gamma}_0)^{-1} \mathbf{\Gamma}_0^T \boldsymbol{\Sigma}^{-1} + \mathbf{\Gamma} (\mathbf{\Gamma}^T \boldsymbol{\Sigma} \mathbf{\Gamma})^{-1} \mathbf{\Gamma}^T = \boldsymbol{\Sigma}^{-1}$. Hence we have $n \boldsymbol{\Sigma}^{-1} (\text{vec}(\bar{\mathbf{X}}) - \text{vec}(\boldsymbol{\mu})) = 0$.

Then $L(\Gamma)$ is maximized with respect to $\text{vec}(\boldsymbol{\mu})$ when $\text{vec}(\boldsymbol{\mu}) = \bar{\mathbf{X}}$. Substituting $\bar{\mathbf{X}}$ for $\text{vec}(\boldsymbol{\mu})$, the likelihood function becomes

$$\begin{aligned} L(\Gamma) &= -\frac{npq}{2} \log 2(\pi) - \frac{n}{2} \log |\mathbf{D}| - \frac{1}{2} \sum_k n_k \log |\Gamma^T \boldsymbol{\Sigma}_k \Gamma| \\ &\quad - \frac{1}{2} \sum_k n_k \text{tr} \{ \Gamma^T \tilde{\boldsymbol{\Sigma}}_k \Gamma (\Gamma^T \boldsymbol{\Sigma}_k \Gamma)^{-1} \} - \frac{1}{2} \sum_y n_y \text{tr} (\mathbf{K} \mathbf{D}^{-1} \mathbf{K}^T \tilde{\boldsymbol{\Delta}}_k), \end{aligned}$$

where $\tilde{\boldsymbol{\Delta}}_k = \tilde{\boldsymbol{\Sigma}}_k + (\bar{\mathbf{X}}_k - \bar{\mathbf{X}})(\bar{\mathbf{X}}_k - \bar{\mathbf{X}})^T$. Recall that $\mathbf{K} = (\Gamma_0 - \Gamma \mathbf{H}^T)$. So the only term that involves \mathbf{H} is $\frac{1}{2} \sum_k n_k \text{tr} ((\Gamma_0 - \Gamma \mathbf{H}^T) \mathbf{D}^{-1} (\Gamma_0 - \Gamma \mathbf{H}^T)^T \tilde{\boldsymbol{\Delta}}_k)$.

By solving the equality $\frac{\partial L}{\partial \mathbf{H}} = 0$ and applying Lemma 4, we get $\sum_k n_k \mathbf{D}^{-1} \Gamma_0^T \tilde{\boldsymbol{\Delta}}_k \Gamma + \sum_k n_k \mathbf{D}^{-1} \mathbf{H} \Gamma^T \tilde{\boldsymbol{\Delta}}_k \Gamma = 0$. So $\hat{\mathbf{H}} = (\sum_k n_k \Gamma_0^T \tilde{\boldsymbol{\Delta}}_k \Gamma) (\sum_k n_k \Gamma^T \tilde{\boldsymbol{\Delta}}_k \Gamma)^{-1} = (\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma) (\Gamma^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma)^{-1}$ is the minimizer of $L(\Gamma)$.

Next we solve the equation $\frac{\partial L}{\partial \mathbf{D}} = 0$. By Lemma 5, the solution is

$$\begin{aligned} \hat{\mathbf{D}} &= \hat{\mathbf{K}}^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \hat{\mathbf{K}} \\ &= (\Gamma_0 - \hat{\mathbf{H}} \Gamma^T) \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} (\Gamma_0 - \hat{\mathbf{H}} \Gamma^T)^T \\ &= ((\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1} \Gamma_0)^{-1} \Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1}) \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} ((\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1} \Gamma_0)^{-1} \Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1})^T \\ &= (\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1} \Gamma)^{-1}. \end{aligned}$$

Finally by solving $\frac{\partial L}{\partial \Gamma^T \boldsymbol{\Sigma}_k \Gamma} = 0$, we get $\Gamma^T \boldsymbol{\Sigma}_k \Gamma = \Gamma^T \hat{\boldsymbol{\Sigma}}_k \Gamma$. After plugging back $\Gamma^T \boldsymbol{\Sigma}_k \Gamma = \Gamma^T \hat{\boldsymbol{\Sigma}}_k \Gamma$ into $L(\Gamma)$, we prove the asserted result. \square

S.2.8 Proof of Proposition 7

Proof. Let $\Gamma = \bigotimes_{m=M}^1 \Gamma_m$ be the basis matrix of $\mathcal{E}_{Y|X}$, and (Γ, Γ_0) be orthonormal matrix, $\Gamma^T \Gamma_0 = 0$. Under FELAD model assumption, we have $\boldsymbol{\Sigma}_k = (\bigotimes_{m=M}^1 \Gamma_m) \boldsymbol{\Omega}_k (\bigotimes_{m=M}^1 \Gamma_m^T) + \Gamma_0 \boldsymbol{\Omega}_0 \Gamma_0^T$, by which we can show that $|\boldsymbol{\Sigma}_k| = |\boldsymbol{\Omega}_k| |\boldsymbol{\Omega}_0|$ and $\boldsymbol{\Sigma}_k^{-1} = (\bigotimes_{m=M}^1 \Gamma_m) \boldsymbol{\Omega}_k^{-1} (\bigotimes_{m=M}^1 \Gamma_m^T) + \Gamma_0 \boldsymbol{\Omega}_0^{-1} \Gamma_0^T$. By Proposition 1, $\text{vec}(\boldsymbol{\mu}_k - \boldsymbol{\mu}) = \Gamma \boldsymbol{\Sigma} \boldsymbol{\nu}_k$ for some vector $\boldsymbol{\nu}_k$. For the FELAD model, $\boldsymbol{\Sigma} \Gamma \boldsymbol{\nu}_k =$

$\Gamma(\sum_{k=1}^K \pi_k \mathbf{\Omega}_k) \boldsymbol{\nu}_k$. Let $\boldsymbol{\alpha}_k = (\sum_{k=1}^K \pi_k \mathbf{\Omega}_k) \boldsymbol{\nu}_k$. Then

$$\begin{aligned} L(\Gamma) &= -\frac{np}{2} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K n_k \log |\mathbf{\Omega}_k| - \frac{n}{2} \log |\mathbf{\Omega}_0| \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_k} (\Gamma^T \text{vec}(\mathbf{X}_{ki} - \bar{\boldsymbol{\mu}}) - \boldsymbol{\alpha}_k)^T \mathbf{\Omega}_k^{-1} (\Gamma^T \text{vec}(\mathbf{X}_{ki} - \bar{\boldsymbol{\mu}}) - \boldsymbol{\alpha}_k) \\ &\quad - \frac{n}{2} \sum_{i=1}^n (\Gamma_0^T \text{vec}(\mathbf{X}_i - \bar{\boldsymbol{\mu}}))^T \mathbf{\Omega}_0^{-1} (\Gamma_0^T \text{vec}(\mathbf{X}_i - \bar{\boldsymbol{\mu}})). \end{aligned}$$

By solving the equation $\frac{\partial L}{\partial \boldsymbol{\alpha}_k} = 0$, we have

$$\mathbf{\Omega}_k^{-1} (\Gamma^T \text{vec}(\mathbf{X}_{ki} - \bar{\boldsymbol{\mu}}) - \boldsymbol{\alpha}_k) = 0,$$

$$\boldsymbol{\alpha}_k = \Gamma^T \text{vec}(\bar{\mathbf{X}}_k - \bar{\boldsymbol{\mu}}).$$

Hence

$$L(\Gamma) = -\frac{np}{2} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K n_k \log |\mathbf{\Omega}_k| - \frac{n}{2} \log |\mathbf{\Omega}_0| - \frac{1}{2} \sum_{k=1}^K n_k \text{tr}(\Gamma^T \tilde{\boldsymbol{\Sigma}}_k \Gamma^T \mathbf{\Omega}_k^{-1}) - \frac{n}{2} \text{tr}\{\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma_0 \mathbf{\Omega}_0^{-1}\}.$$

By solving the equation $\frac{\partial L}{\partial \mathbf{\Omega}_0} = 0$ and applying Lemma 4, we have

$$\mathbf{\Omega}_0 = \Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma_0.$$

By solving $\frac{\partial L}{\partial \mathbf{\Omega}_k} = 0$ and using Lemma 4 again, we get

$$\mathbf{\Omega}_k = \Gamma^T \tilde{\boldsymbol{\Sigma}}_k \Gamma.$$

Plugging the above relation into $L(\Gamma)$, we have

$$L(\Gamma) = -\frac{np}{2} \log(2\pi) - \frac{np}{2} - \frac{n}{2} \log |\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma_0| - \frac{1}{2} \sum_{k=1}^K n_k \log |\Gamma^T \tilde{\boldsymbol{\Sigma}}_k \Gamma|.$$

By Lemma 2, we have $\log |\Gamma_0^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}} \Gamma_0| = \log |\Gamma^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1} \Gamma| + \log |\tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}|$. Finally, we have

$$L(\Gamma) = -\frac{np}{2} \log(2\pi + 1) - \frac{n}{2} \log |\tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}| - \frac{n}{2} \log |\Gamma^T \tilde{\boldsymbol{\Sigma}}_{\mathbf{X}}^{-1} \Gamma| - \frac{1}{2} \sum_{k=1}^K n_k \log |\Gamma^T \tilde{\boldsymbol{\Sigma}}_k \Gamma|.$$

□

S.2.9 Proof of Proposition 8

Proof. For the FLAD model, the free parameter ϕ^T can be written as $(\text{vec}^T(\boldsymbol{\mu}), \text{vec}^T(\boldsymbol{\alpha}_1), \dots, \text{vec}^T(\boldsymbol{\alpha}_{K-1}), \text{vec}^T(\boldsymbol{\Gamma}_1), \dots, \text{vec}^T(\boldsymbol{\Gamma}_M), \text{vech}^T(\boldsymbol{\Sigma}), \text{vech}^T(\mathbf{M}_1), \dots, \text{vech}^T(\mathbf{M}_K))^T = (\phi_1^T, \dots, \phi_{2K+M}^T)^T$ with parameter space Θ being the Cartesian product of the parameter space for the individual components. Because the elements of \mathbf{h} are analytic, they are twice continuously differentiable over Θ and every point in Θ is regular (Shapiro 1986, Definition 2.1), except on a set of Lebesgue measure 0.

Let $\widehat{\boldsymbol{\Sigma}}^T = (\text{vec}^T(\widehat{\boldsymbol{\mu}}_1), \dots, \text{vec}^T(\widehat{\boldsymbol{\mu}}_K), \text{vech}^T(\widehat{\boldsymbol{\Sigma}}_1), \dots, \text{vech}^T(\widehat{\boldsymbol{\Sigma}}_K))^T$ be the MLE under the full model. The FLAD discrepancy function is defined as $F_{FLAD}(\widehat{\boldsymbol{\Sigma}}, \mathbf{h}) = L_p(\widehat{\boldsymbol{\Sigma}} | \widehat{\boldsymbol{\Sigma}}) - L_d(h | \widehat{\boldsymbol{\Sigma}})$, where

$$L_d(\mathbf{h} | \widehat{\boldsymbol{\Sigma}}) = - \sum_{k=1}^K \frac{n_k}{2} \{ \log |\boldsymbol{\Sigma}_k| + \text{tr}(\widehat{\boldsymbol{\Sigma}}_k \boldsymbol{\Sigma}_k^{-1}) + (\bar{\mathbf{X}}_k - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\bar{\mathbf{X}}_k - \boldsymbol{\mu}_k) \}.$$

F_{FLAD} is an analytic function of $\widehat{\boldsymbol{\Sigma}}$ and \mathbf{h} . Also note that $F_{FLAD} \geq 0$ for all $\widehat{\boldsymbol{\Sigma}}$ and \mathbf{h} , and $F_{FLAD} = 0$ if and only if $\widehat{\boldsymbol{\Sigma}} = \mathbf{h}$. So F_{FLAD} satisfies the necessary condition for a discrepancy function (Shapiro 1986).

Let $\mathbf{W} = \frac{1}{2} \frac{\partial^2 F_{FLAD}}{\partial \mathbf{h} \partial \mathbf{h}^T}$ evaluated at $(\mathbf{h}_0, \mathbf{h}_0)$, where \mathbf{h}_0 is the true value of \mathbf{h} . It is easy to check $\mathbf{W} = \mathbf{J}$, which is also a sufficient condition for FLAD to give asymptotic efficient estimators (Shapiro 1986)(eqn.5.1). Then, using Proposition 4.1 in (Shapiro 1986), we have $\sqrt{n}(\widehat{\mathbf{h}} - \mathbf{h}) \xrightarrow{D} N(0, \mathbf{V})$, where $\mathbf{V} = \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J} \mathbf{J}^{-1} \mathbf{J} \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T = \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T$. By calculation, we have $\mathbf{V}_0^{-1/2}(\mathbf{V}_0 - \mathbf{V})\mathbf{V}_0^{-1/2} = \mathbf{I} - \mathbf{J}^{1/2} \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J}^{1/2} = \mathbf{Q}_{\mathbf{J}^{1/2} \mathbf{H}} \geq 0$.

The proof of the situation of FELAD is similar to that of FLAD, so we omit the details here.

Let $\mathbf{P}_0 = \mathbf{V}_0^{-1/2} \mathbf{V} \mathbf{V}_0^{-1/2}$, and $\mathbf{P}_1 = \mathbf{V}_0^{-1/2} \mathbf{V}_1 \mathbf{V}_0^{-1/2}$. We have $\mathbf{P}_0 = \mathbf{J}^{1/2} \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J}^{1/2}$, $\mathbf{P}_1 = \mathbf{J}^{1/2} \mathbf{H} \mathbf{G}_1 (\mathbf{G}_1^T \mathbf{H}^T \mathbf{J} \mathbf{H} \mathbf{G}_1)^\dagger \mathbf{G}_1^T \mathbf{H}^T \mathbf{J}^{1/2}$, and $\mathbf{P}_0 \mathbf{P}_1 = \mathbf{J}^{1/2} \{ \mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J} \mathbf{H} \} \mathbf{G}_1 (\mathbf{G}_1^T \mathbf{H}^T \mathbf{J} \mathbf{H} \mathbf{G}_1)^\dagger \mathbf{G}_1^T \mathbf{H}^T \mathbf{J}^{1/2}$.

Using matrix identity (Rao & Mitra 1971)(Theo 2.4(c))

$$\mathbf{H}(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J} \mathbf{H} = \mathbf{H},$$

and

$$(\mathbf{H}^T \mathbf{J} \mathbf{H})(\mathbf{H}^T \mathbf{J} \mathbf{H})^\dagger \mathbf{H}^T = \mathbf{H}^T,$$

we get $\mathbf{P}_0\mathbf{P}_1 = \mathbf{P}_1$. Similarly, we also have $\mathbf{P}_1\mathbf{P}_0 = \mathbf{P}_1$. Then $\mathbf{P}_0 - \mathbf{P}_1 = \mathbf{P}_0 - \mathbf{P}_0\mathbf{P}_1 = \mathbf{P}_0(\mathbf{I} - \mathbf{P}_1) = \mathbf{P}_{\mathbf{J}^{1/2}\mathbf{H}}\mathbf{Q}_{\mathbf{J}^{1/2}\mathbf{H}\mathbf{G}_1}$. Also $\mathbf{P}_0 - \mathbf{P}_1 = \mathbf{P}_0 - \mathbf{P}_1\mathbf{P}_0 = (\mathbf{I} - \mathbf{P}_1)\mathbf{P}_0 = \mathbf{Q}_{\mathbf{J}^{1/2}\mathbf{H}\mathbf{G}_1}\mathbf{P}_{\mathbf{J}^{1/2}\mathbf{H}}$.

□

S.2.10 Proof of Proposition 9

The proof of this proposition is parallel to the proof of Proposition 8. However, the Fisher information \mathbf{J} is changed and may not equal to $\mathbf{W} = \frac{1}{2} \frac{\partial^2 \mathbb{E}_{FLAD}}{\partial \mathbf{h} \partial \mathbf{h}^T}$ at $(\mathbf{h}_0, \mathbf{h}_0)$. In this case, by Proposition 4.1 in Shapiro (1986), the estimator $\hat{\mathbf{h}}$ is still \sqrt{n} -consistent, but may not be the most efficient one.

S.2.11 Proof of Proposition 10

Proof. The population objective function of FLAD is

$$\begin{aligned} L(\mathcal{S}) &= \frac{1}{2} \log |\mathbf{\Gamma}_0^T \mathbf{\Sigma}_{\mathbf{X}}^{-1} \mathbf{\Gamma}_0| - \frac{1}{2} \sum_{k=1}^K \pi_k \log |\mathbf{\Gamma}_0^T \mathbf{\Sigma}_k^{-1} \mathbf{\Gamma}_0| - \frac{1}{2} \log |\mathbf{\Sigma}_{\mathbf{X}}| + \frac{1}{2} \sum_{k=1}^K \pi_k \log |\mathbf{\Sigma}_k| \\ &\leq \frac{1}{2} \log |\mathbf{\Gamma}_0^T (\sum_{k=1}^K \pi_k \mathbf{\Sigma}_k)^{-1} \mathbf{\Gamma}_0| - \frac{1}{2} \sum_{k=1}^K \pi_k \log |\mathbf{\Gamma}_0^T \mathbf{\Sigma}_k^{-1} \mathbf{\Gamma}_0| - \frac{1}{2} \log |\mathbf{\Sigma}_{\mathbf{X}}| + \frac{1}{2} \sum_{k=1}^K \pi_k \log |\mathbf{\Sigma}_k| \\ &\leq -\frac{1}{2} \log |\mathbf{\Sigma}_{\mathbf{X}}| + \frac{1}{2} \sum_{k=1}^K \pi_k \log |\mathbf{\Sigma}_k|. \end{aligned}$$

where $\mathbf{\Gamma}_m \in \mathbb{R}^{p_m \times d_m}$, $\mathbf{\Gamma}_0$ is the orthogonal complement of $\mathbf{\Gamma} = \bigotimes_{m=M}^1 \mathbf{\Gamma}_m$, and $\mathcal{S} = \text{span}(\mathbf{\Gamma})$.

The first inequality follows since $\mathbf{\Sigma}_{\mathbf{X}}^{-1} \leq (\sum_{k=1}^K \pi_k \mathbf{\Sigma}_k)^{-1}$ and the second inequality follows since the function $\log |\mathbf{\Gamma}_0^T \mathbf{\Delta}^{-1} \mathbf{\Gamma}_0|$ is convex in $\mathbf{\Delta}$ on the space formed by symmetric and positive definite matrices (See Section A.6 of Cook & Forzani (2009)). Define $\mathcal{S}_{FLAD} = \text{argmax} L(\mathcal{S})$. If we find an orthogonal matrix $(\boldsymbol{\eta}, \boldsymbol{\eta}_0)$, where $\boldsymbol{\eta}_m \in \mathbb{R}^{p_m \times d_m}$ and $\boldsymbol{\eta} = \bigotimes_{m=M}^1 \boldsymbol{\eta}_m$, such that $\frac{1}{2} \log |\boldsymbol{\eta}_0^T \mathbf{\Sigma}_{\mathbf{X}}^{-1} \boldsymbol{\eta}_0| = \frac{1}{2} \sum_{k=1}^K \pi_k \log |\boldsymbol{\eta}_0^T \mathbf{\Sigma}_k^{-1} \boldsymbol{\eta}_0|$, then $\mathcal{S}_{FLAD} = \text{span}(\boldsymbol{\eta})$. Next, we show how to get the desired $(\boldsymbol{\eta}, \boldsymbol{\eta}_0)$ such that $\text{span}(\boldsymbol{\eta}) \subseteq \mathcal{T}_{Y|\mathbf{X}}$ and $\frac{1}{2} \log |\boldsymbol{\eta}_0^T \mathbf{\Sigma}_{\mathbf{X}}^{-1} \boldsymbol{\eta}_0| = \frac{1}{2} \sum_{k=1}^K \pi_k \log |\boldsymbol{\eta}_0^T \mathbf{\Sigma}_k^{-1} \boldsymbol{\eta}_0|$.

We first introduce the definition of ‘‘Kronecker envelope’’ proposed by Li et al. (2010).

Definition 1. Let $\mathbf{U} \in \mathbb{R}^{\prod_{m=1}^M p_M \times l}$ be a random matrix. There are subspaces $\mathcal{S}_1 \subseteq \mathbb{P}^{p_m}$, $m = 1, \dots, M$ such that (i) $\text{span}(\mathbf{U}) \subseteq \bigotimes_{m=1}^M \mathcal{S}_m$ almost surely; and (ii) If there exists other subspaces $\mathcal{S}'_1 \subseteq \mathbb{P}^{p_m}$, $m = 1, \dots, M$, that satisfies Condition (i), then $\bigotimes_{m=1}^M \mathcal{S}_m \subseteq \bigotimes_{m=1}^M \mathcal{S}'_m$. Then the subspace $\bigotimes_{m=1}^M \mathcal{S}_m$ is called the Kronecker envelope of \mathbf{U} , denoted as $\mathcal{E}^\otimes(\mathbf{U})$.

We choose the random matrix \mathbf{U} to be $\Sigma_{\mathbf{X}}^{-1}(\Sigma_{\mathbf{X}} - \text{cov}(\mathbf{X} | Y))$, where Y take values in $\{1, \dots, K\}$. Under the condition that $E(\text{vec}(\mathbf{X}) | \beta^T \text{vec}(\mathbf{X}))$ is linear in $\beta^T \text{vec}(\mathbf{X})$, and $\text{var}(\text{vec}(\mathbf{X}) | \beta^T \text{vec}(\mathbf{X}))$ is nonrandom, Cook & Weisberg (1991) showed that the subspace spanned by $\mathbf{U} = \Sigma_{\mathbf{X}}^{-1}(\Sigma_{\mathbf{X}} - \text{cov}(\mathbf{X} | Y))$ is subspace of $\mathcal{S}_{Y|\text{vec}(\mathbf{X})}$. By Theorem 5 of Li et al. (2010), we have $\mathcal{E}^\otimes(\mathbf{U}) \subseteq \mathcal{T}_{Y|\mathbf{X}}$. Let $\boldsymbol{\eta} = \bigotimes_{m=1}^M \boldsymbol{\eta}_m$ be a basis matrix for $\mathcal{E}^\otimes(\mathbf{U})$. Then $\Sigma_{\mathbf{X}}^{-1}(\Sigma_{\mathbf{X}} - \Sigma_k) = \boldsymbol{\eta} \mathbf{w}_k$ for some matrix \mathbf{w}_k , $k = 1, \dots, K$. Consequently, $\Sigma_{\mathbf{X}}^{-1}(\Sigma_{\mathbf{X}} - \Sigma_k) = \mathbf{P}_{\boldsymbol{\eta}(\Sigma_{\mathbf{X}})} \Sigma_{\mathbf{X}}^{-1}(\Sigma_{\mathbf{X}} - \Sigma_k)$. Thus, $\Sigma_{\mathbf{X}} = \Sigma_k + \mathbf{P}_{\boldsymbol{\eta}(\Sigma_{\mathbf{X}})}^T (\Sigma_{\mathbf{X}} - \Sigma_k) = \Sigma_k + \mathbf{P}_{\boldsymbol{\eta}(\Sigma_{\mathbf{X}})}^T (\Sigma_{\mathbf{X}} - \Sigma_k) \mathbf{P}_{\boldsymbol{\eta}(\Sigma_{\mathbf{X}})}$. Then, by direct multiplication, we have $\Sigma_k^{-1} = \Sigma_{\mathbf{X}}^{-1} + \boldsymbol{\eta} \{(\boldsymbol{\eta}^T \Sigma_k \boldsymbol{\eta})^{-1} - (\boldsymbol{\eta}^T \Sigma_{\mathbf{X}} \boldsymbol{\eta})^{-1}\} \boldsymbol{\eta}^T$, which implies $\frac{1}{2} \log |\boldsymbol{\eta}_0^T \Sigma_{\mathbf{X}}^{-1} \boldsymbol{\eta}_0| = \frac{1}{2} \sum_{k=1}^K \pi_k \log |\boldsymbol{\eta}_0^T \Sigma_k^{-1} \boldsymbol{\eta}_0|$. As such, $\mathcal{S}_{FLAD} = \text{span}(\boldsymbol{\eta}) \subseteq \mathcal{T}_{Y|\mathbf{X}}$. The \sqrt{n} -consistent property of \mathcal{S}_{FLAD} can be obtained from Proposition 4.1 of Shapiro (1986), the proof is similar to that of Proposition 8 and 9, we omit the details here. \square

S.3 Closed-Form derivatives of Algorithm 1

The objective functions (3.8) in the paper can be solved by standard Stiefel or Grassmann manifold optimization packages, where we can plug in the closed-form derivatives to speed up the computation. Here

we give the closed-form derivatives when $M = 2$.

$$\begin{aligned}
& \frac{\partial(\log |(\mathbf{I}_{d_2} \otimes \mathbf{\Gamma}_1^T)\mathbf{M}(\mathbf{I}_{d_2} \otimes \mathbf{\Gamma}_1)|)}{\partial \text{vec}(\mathbf{\Gamma}_1)} \\
&= 2\text{vec} \left\{ \mathbf{M}(\mathbf{I}_{d_2} \otimes \mathbf{\Gamma}_1) \{ (\mathbf{I}_{d_2} \otimes \mathbf{\Gamma}_1^T) \mathbf{M}(\mathbf{I}_{d_2} \otimes \mathbf{\Gamma}_1) \}^{-1} \right\} (\mathbf{I}_{d_2} \otimes \mathbf{T}_{d_1 d_2} \otimes \mathbf{I}_{p_1}) \{ \text{vec}(\mathbf{I}_{d_2}) \otimes \mathbf{I}_{p_1 d_1} \}, \\
& \frac{\partial(\log |(\mathbf{\Gamma}_2^T \otimes \mathbf{I}_{d_1})\mathbf{M}(\mathbf{\Gamma}_2^T \otimes \mathbf{I}_{d_1})|)}{\partial \text{vec}(\mathbf{\Gamma}_2)} \\
&= 2\text{vec} \left\{ \mathbf{M}(\mathbf{\Gamma}_2 \otimes \mathbf{I}_{d_1}) \{ (\mathbf{\Gamma}_2^T \otimes \mathbf{I}_{d_1}) \mathbf{M}(\mathbf{\Gamma}_2 \otimes \mathbf{I}_{d_1}) \}^{-1} \right\} (\mathbf{I}_{d_2} \otimes \mathbf{T}_{d_1 p_2} \otimes \mathbf{I}_{d_1}) \{ \mathbf{I}_{p_2 d_2} \otimes \text{vec}(\mathbf{I}_{d_1}) \},
\end{aligned}$$

where \mathbf{T}_{mn} is $mn \times mn$ permutation matrix whose ij -th element is 1 if $j = 1 + m(i-1) - (mn-1) \lfloor \frac{i-1}{n} \rfloor$

or 0 otherwise.

References

- Cook, R. D. & Forzani, L. (2009), ‘Likelihood-based sufficient dimension reduction’, *Journal of the American Statistical Association* **104**(485), 197–208.
- Cook, R. D. & Weisberg, S. (1991), ‘Comment: Sliced inverse regression for dimension reduction’, *Journal of the American Statistical Association* **86**(414), 328–332.
- Li, B., Kim, M. K., Altman, N. et al. (2010), ‘On dimension folding of matrix-or array-valued statistical objects’, *The Annals of Statistics* **38**(2), 1094–1121.
- Rao, C. & Mitra, S. K. (1971), ‘Generalized inverse of a matrix and its applications’, *Berkeley Symposium on Mathematical Statistics and Probability* **1**, 601–620.
- Shapiro, A. (1986), ‘Asymototic theory of overparameterized structural models’, *Journal of the American Statistical Association* **81**, 142–149.
- Sheng, W. & Yuan, Q. (2019), ‘Sufficient dimension folding in regression via distance covariance for matrix-valued predictors’, *Statistical Analysis and Data Mining* .
- Wang, W., Zhang, X. & Li, L. (2019), ‘Common reducing subspace model and network alternation analysis’, *Biometrics* **75**(4), 1109–1120.
- Xue, Y. & Yin, X. (2014), ‘Sufficient dimension folding for regression mean fuction’, *Journal of Computioanl and Graphical Statistics* **23**, 1028–1043.