
Supplementary Material: On the Faster Alternating Least-Squares for CCA

1 Theoretical Analysis

Theorem 3.1 Given data matrices $(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{d_x \times n} \times \mathbb{R}^{d_y \times n}$, if $\sigma_k^2 > 2\sqrt{\beta} = \sigma_{k+1}^2$, Algorithm 1 then computes Φ_T and Ψ_T which are estimates of top- k canonical subspaces (\mathbf{U}, \mathbf{V}) such that $\sin \theta_T \leq \epsilon$ and $\Phi_T^\top \mathbf{C}_{xx} \Phi_T = \Psi_T^\top \mathbf{C}_{yy} \Psi_T = \mathbf{I}$, in $T = O\left(\sqrt{\frac{\sigma_k^2}{\sigma_k^2 - \sigma_{k+1}^2}} \log \frac{1}{\epsilon(\sigma_k^2 - \sigma_{k+1}^2) \cos \theta_0}\right)$ iterations. If accelerated gradient descent is used as the least-squares solver, the overall running time is at most

$$O\left((dk^2 + k \text{nnz}(\mathbf{X}, \mathbf{Y})) \kappa^{\frac{1}{2}}(\mathbf{X}, \mathbf{Y}) \log \frac{c_1 c_2}{(\sigma_k^2 - \sigma_{k+1}^2) \cos \theta_0}\right) \sqrt{\frac{\sigma_k^2}{\sigma_k^2 - \sigma_{k+1}^2}} \log \frac{1}{\epsilon(\sigma_k^2 - \sigma_{k+1}^2) \cos \theta_0},$$

where $d = \max\{d_x, d_y\}$, $\text{nnz}(\mathbf{X}, \mathbf{Y}) = \text{nnz}(\mathbf{X}) + \text{nnz}(\mathbf{Y})$, $\kappa(\mathbf{X}, \mathbf{Y}) = \max\{\kappa(\mathbf{C}_{xx}), \kappa(\mathbf{C}_{yy})\}$,

$$c_1 = \max_t \frac{\tan \max\{\theta_t, \theta_{\max}(\hat{\Phi}_t, \mathbf{U}), \theta_{\max}(\hat{\Psi}_t, \mathbf{V})\}}{\tan \min\{\theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}}), \theta_{\max}(\mathbf{Q}_t, \tilde{\mathbf{V}})\}}, \text{ and } c_2 = \max_t \left\{ \max_t \frac{\theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{\theta_{\min}(\mathbf{P}_t, \tilde{\mathbf{U}})}, \max_t \frac{\theta_{\max}(\mathbf{Q}_t, \tilde{\mathbf{V}})}{\theta_{\min}(\mathbf{Q}_t, \tilde{\mathbf{V}})} \right\}.$$

Proof We now follow the proof sketch given in the main text to conduct the analysis. Note that $\theta_t \triangleq \max\{\theta_{\max}(\Phi_t, \mathbf{U}), \theta_{\max}(\Psi_t, \mathbf{V})\}$ and we need to show $\sin \theta_{\max}(\Phi_t, \mathbf{U}) \leq \epsilon$ and $\sin \theta_{\max}(\Psi_t, \mathbf{V}) \leq \epsilon$ hold simultaneously. Recall that our update is

$$\begin{cases} \Phi_{t+1} \mathbf{R}_{t+1} = \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} (\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top \Phi_t + \xi_t) - \beta \Phi_{t-1} \mathbf{R}_t^{-1} + \hat{\xi}_t \\ \Psi_{t+1} \mathbf{S}_{t+1} = \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top (\mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_t + \eta_t) - \beta \Psi_{t-1} \mathbf{S}_t^{-1} + \hat{\eta}_t \end{cases}.$$

It can be equivalently written as

$$\begin{cases} \Phi_{t+1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t = \left(\mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} (\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top \Phi_t + \xi_t) - \beta \Phi_{t-1} \mathbf{R}_t^{-1} + \hat{\xi}_t \right) \tilde{\mathbf{R}}_t \\ \Psi_{t+1} \mathbf{S}_{t+1} \tilde{\mathbf{S}}_t = \left(\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top (\mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_t + \eta_t) - \beta \Psi_{t-1} \mathbf{S}_t^{-1} + \hat{\eta}_t \right) \tilde{\mathbf{S}}_t \end{cases},$$

where

$$\tilde{\mathbf{R}}_t = \begin{cases} (\mathbf{I} + \mathbf{R}_t^{-\top} \mathbf{R}_t^{-1})^{-\frac{1}{2}}, & t > 0 \\ \mathbf{I}, & t = 0 \end{cases}, \quad \tilde{\mathbf{S}}_t = \begin{cases} (\mathbf{I} + \mathbf{S}_t^{-\top} \mathbf{S}_t^{-1})^{-\frac{1}{2}}, & t > 0 \\ \mathbf{I}, & t = 0 \end{cases}.$$

We now focus on the first equation. Together with $\Phi_t \tilde{\mathbf{R}}_t = \Phi_t \tilde{\mathbf{R}}_t$, it leads to the following augmented system:

$$\mathbf{P}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t = \mathbf{B}_\phi^{-1} \mathbf{A}_\phi \mathbf{P}_t + \delta_t,$$

where

$$\mathbf{A}_\phi = \begin{pmatrix} \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top & -\beta \mathbf{C}_{xx} \\ \mathbf{C}_{xx} & \mathbf{0} \end{pmatrix}, \quad \mathbf{B}_\phi = \begin{pmatrix} \mathbf{C}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{xx} \end{pmatrix},$$

$$\mathbf{P}_t = \begin{pmatrix} \Phi_t \\ \Phi_{t-1} \mathbf{R}_t^{-1} \end{pmatrix} \tilde{\mathbf{R}}_t, \quad \delta_t = \begin{pmatrix} \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \xi_t + \hat{\xi}_t \\ \mathbf{0} \end{pmatrix} \tilde{\mathbf{R}}_t.$$

Since Φ_t is \mathbf{C}_{xx} -orthonormal, it is easy to see that \mathbf{P}_t is \mathbf{B}_ϕ -orthonormal. To continue, we can write the SVD of \mathbf{C}_{xy} as follows:

$$\mathbf{C}_{xy} = \mathbf{C}_{xx} (\mathbf{U} \Sigma \mathbf{V}^\top + \mathbf{U}_\perp \Sigma_\perp \mathbf{V}_\perp^\top) \mathbf{C}_{yy},$$

where $(\mathbf{U}_\perp, \boldsymbol{\Sigma}_\perp, \mathbf{V}_\perp)$ consists of the $(\text{rank}(\mathbf{C}_{xy}) - k)$ remaining triples of the left singular vector in metric \mathbf{C}_{xx} , singular value, and right singular vector in metric \mathbf{C}_{yy} , other than those in $(\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{V})$. It thus holds that

$$\mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top = \mathbf{C}_{xx} (\mathbf{U} \boldsymbol{\Sigma}^2 \mathbf{U}^\top + \mathbf{U}_\perp \boldsymbol{\Sigma}_\perp^2 \mathbf{U}_\perp^\top) \mathbf{C}_{xx},$$

and accordingly, by Lemma I in Section 2, \mathbf{A}_ϕ has the following Schur decomposition in metric \mathbf{B}_ϕ :

$$\mathbf{A}_\phi = \mathbf{B}_\phi \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \tilde{\boldsymbol{\Sigma}} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi,$$

where notations can be found in Lemma I. Further, since $\tilde{\boldsymbol{\Sigma}}$ and $\tilde{\boldsymbol{\Sigma}}_\perp$ (see Lemma I in Section 2) don't share eigenvalues by assumption that $\sigma_k > \sigma_{k+1}$, there exists (Golub and Van Loan, 2013) a matrix $\boldsymbol{\Xi}$ such that $\tilde{\boldsymbol{\Sigma}} \boldsymbol{\Xi} - \boldsymbol{\Xi} \tilde{\boldsymbol{\Sigma}}_\perp = -\tilde{\boldsymbol{\Sigma}}$ and thus

$$\begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \tilde{\boldsymbol{\Sigma}} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \mathbf{0} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix} \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1}.$$

Plugging the above equation and the Schur decomposition of \mathbf{A}_ϕ into the augmented system and then pre-multiply both sides by $\begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi$, results in the following equation:

$$\begin{aligned} & \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi \mathbf{P}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t = \\ & \begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \mathbf{0} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix} \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi \mathbf{P}_t + \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t. \end{aligned}$$

Letting

$$\begin{aligned} \begin{pmatrix} \tilde{\mathbf{X}}_t \\ \tilde{\mathbf{Y}}_t \end{pmatrix} &= \begin{pmatrix} \mathbf{I} & \boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi \mathbf{P}_t \\ &= \begin{pmatrix} \mathbf{I} & -\boldsymbol{\Xi} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^\text{H} \mathbf{B}_\phi \mathbf{P}_t = \begin{pmatrix} \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \mathbf{P}_t \\ \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \mathbf{P}_t \end{pmatrix}, \end{aligned}$$

where $\mathbf{W}_\text{U} = \tilde{\mathbf{U}} - \tilde{\mathbf{U}}_\perp \boldsymbol{\Xi}^\text{H}$, the above equation then can be split into the following two equations:

$$\begin{cases} \tilde{\mathbf{X}}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t = \tilde{\boldsymbol{\Sigma}} \tilde{\mathbf{X}}_t + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \\ \tilde{\mathbf{Y}}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t = \tilde{\boldsymbol{\Sigma}}_\perp \tilde{\mathbf{Y}}_t + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \end{cases}.$$

Hence, we have that

$$\begin{aligned} \tilde{\mathbf{Z}}_{t+1} &= \tilde{\mathbf{Y}}_{t+1} \tilde{\mathbf{X}}_{t+1}^{-1} = \left(\tilde{\mathbf{Y}}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t \right) \left(\tilde{\mathbf{X}}_{t+1} \tilde{\mathbf{R}}_{t+1}^{-1} \mathbf{R}_{t+1} \tilde{\mathbf{R}}_t \right)^{-1} \\ &= \left(\tilde{\boldsymbol{\Sigma}}_\perp \tilde{\mathbf{Y}}_t + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \right) \left(\tilde{\boldsymbol{\Sigma}} \tilde{\mathbf{X}}_t + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \right)^{-1} \\ &= \left(\tilde{\boldsymbol{\Sigma}}_\perp \tilde{\mathbf{Y}}_t \tilde{\mathbf{X}}_t^{-1} + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \right) \left(\tilde{\boldsymbol{\Sigma}} + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \right)^{-1} \\ &= \left(\tilde{\boldsymbol{\Sigma}}_\perp \tilde{\mathbf{Z}}_t + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \left(\tilde{\mathbf{Z}}_t^\text{H} \tilde{\mathbf{Z}}_t \right)^{-1} \tilde{\mathbf{Z}}_t^\text{H} \tilde{\mathbf{Z}}_t \right) \left(\tilde{\boldsymbol{\Sigma}} + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \right)^{-1} \\ &= \left(\tilde{\boldsymbol{\Sigma}}_\perp + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \left(\tilde{\mathbf{Z}}_t^\text{H} \tilde{\mathbf{Z}}_t \right)^{-1} \tilde{\mathbf{Z}}_t^\text{H} \right) \tilde{\mathbf{Z}}_t \left(\tilde{\boldsymbol{\Sigma}} + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \right)^{-1}, \end{aligned}$$

and thus that

$$\tilde{\mathbf{Z}}_T = \prod_{t=T-1}^0 \left(\tilde{\boldsymbol{\Sigma}}_\perp + \tilde{\mathbf{U}}_\perp^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_t \tilde{\mathbf{X}}_t^{-1} \left(\tilde{\mathbf{Z}}_t^\text{H} \tilde{\mathbf{Z}}_t \right)^{-1} \tilde{\mathbf{Z}}_t^\text{H} \right) \tilde{\mathbf{Z}}_0 \prod_{t'=0}^{T-1} \left(\tilde{\boldsymbol{\Sigma}} + \mathbf{W}_\text{U}^\text{H} \mathbf{B}_\phi \boldsymbol{\delta}_{t'} \tilde{\mathbf{X}}_{t'}^{-1} \right)^{-1}.$$

By Lemma 12 in Ge et al. (2016), $\sin \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}}) = \|\tilde{\mathbf{U}}_{\perp}^{\top} \mathbf{B}_{\phi} \mathbf{P}_t\|_2$. We then can write that

$$\sin \theta_{\max}(\mathbf{P}_T, \tilde{\mathbf{U}}) = \|\tilde{\mathbf{Y}}_T\|_2 \leq \|\tilde{\mathbf{Z}}_T\|_2 \|\tilde{\mathbf{X}}_T\|_2,$$

where

$$\begin{aligned} \|\tilde{\mathbf{X}}_T\|_2 &= \left\| \begin{pmatrix} \mathbf{I} & -\mathbf{\Xi} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix}^{\text{H}} \mathbf{B}_{\phi} \mathbf{P}_t \right\|_2 \\ &\leq \left\| \begin{pmatrix} \mathbf{I} & -\mathbf{\Xi} \end{pmatrix} \right\|_2 \left\| \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} \right\|_{\mathbf{B}_{\phi,2}} \|\mathbf{P}_t\|_{\mathbf{B}_{\phi,2}} = 1 + \|\mathbf{\Xi}\|_2, \end{aligned}$$

where the last equality with $\|\mathbf{P}_t\|_{\mathbf{B}_{\phi}} = \|\mathbf{B}_{\phi}^{\frac{1}{2}} \mathbf{P}_t\|_2 = 1$ is by the \mathbf{B}_{ϕ} -orthonormality of both $\begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix}$ and \mathbf{P}_t . What's more, we have permutation matrix $\mathbf{\Pi}$, constant $1 + \gamma = \frac{8(1+\beta)}{\sigma_k^{\pm} - \sqrt{\beta}}$, and diagonal matrix $\mathbf{\Gamma} = \text{diag}(\text{diag}(1, 1 + \gamma), \text{diag}(1, 1 + \gamma), \mathbf{I})$ such that

$$\mathbf{\Pi} \tilde{\mathbf{\Sigma}}_{\perp} \mathbf{\Pi}^{\top} = \text{diag}(\mathbf{\Sigma}_{k+1}, \dots, \mathbf{\Sigma}_r, \mathbf{I}), \quad \mathbf{\Gamma} \mathbf{\Pi} \tilde{\mathbf{\Sigma}}_{\perp} \mathbf{\Pi}^{\top} \mathbf{\Gamma}^{-1} = \text{diag}(\mathbf{\Sigma}_{k+1}^{(\gamma)}, \dots, \mathbf{\Sigma}_r^{(\gamma)}, \mathbf{I}),$$

where σ_j^{\pm} is defined in Lemma I of Section 2, and

$$\mathbf{\Sigma}_j = \begin{pmatrix} \sigma_j^+ & -(1 + (\sigma_j^+)^2) \\ 0 & \sigma_j^- \end{pmatrix}, \quad \mathbf{\Sigma}_j^{(\gamma)} = \begin{pmatrix} \sigma_j^+ & -\frac{1 + (\sigma_j^+)^2}{1 + \gamma} \\ 0 & \sigma_j^- \end{pmatrix}.$$

Thus, we can write that

$$\begin{aligned} \|\tilde{\mathbf{Z}}_T\|_2 &= \left\| \mathbf{\Pi}^{\top} \mathbf{\Gamma}^{-1} \prod_{t=T-1}^0 \left(\mathbf{\Gamma} \mathbf{\Pi} \tilde{\mathbf{\Sigma}}_{\perp} \mathbf{\Pi}^{\top} \mathbf{\Gamma}^{-1} + \mathbf{\Gamma} \mathbf{\Pi} \tilde{\mathbf{U}}_{\perp}^{\text{H}} \mathbf{B}_{\phi} \delta_t \right. \right. \\ &\quad \left. \left. \tilde{\mathbf{X}}_t^{-1} (\tilde{\mathbf{Z}}_t^{\text{H}} \tilde{\mathbf{Z}}_t)^{-1} \tilde{\mathbf{Z}}_t^{\text{H}} \mathbf{\Pi}^{\top} \mathbf{\Gamma}^{-1} \right) \mathbf{\Gamma} \mathbf{\Pi} \tilde{\mathbf{Z}}_0 \prod_{t'=0}^{T-1} \left(\tilde{\mathbf{\Sigma}} + \mathbf{W}_{\mathbf{U}}^{\text{H}} \mathbf{B}_{\phi} \delta_{t'} \tilde{\mathbf{X}}_{t'}^{-1} \right)^{-1} \right\| \\ &\leq \|\mathbf{\Gamma}\|_2 \|\mathbf{\Gamma}^{-1}\|_2 \|\tilde{\mathbf{Z}}_0\|_2 \prod_{t=0}^{T-1} \left\| \left(\tilde{\mathbf{\Sigma}} + \mathbf{W}_{\mathbf{U}}^{\text{H}} \mathbf{B}_{\phi} \delta_t \tilde{\mathbf{X}}_t^{-1} \right)^{-1} \right\|_2 \left(\|\mathbf{\Gamma} \mathbf{\Pi} \tilde{\mathbf{\Sigma}}_{\perp} \mathbf{\Pi}^{\top} \mathbf{\Gamma}^{-1}\|_2 \right. \\ &\quad \left. + \|\mathbf{\Gamma}\|_2 \|\mathbf{\Gamma}^{-1}\|_2 \|\tilde{\mathbf{U}}_{\perp}\|_{\mathbf{B}_{\phi,2}} \|\delta_t\|_{\mathbf{B}_{\phi,2}} \|\tilde{\mathbf{X}}_t^{-1} (\tilde{\mathbf{Z}}_t^{\text{H}} \tilde{\mathbf{Z}}_t)^{-1} \tilde{\mathbf{Z}}_t^{\text{H}}\|_2 \right), \end{aligned}$$

where $\|\mathbf{\Gamma}\|_2 \|\mathbf{\Gamma}^{-1}\|_2 = 1 + \gamma$ and $\|\tilde{\mathbf{U}}_{\perp}\|_{\mathbf{B}_{\phi,2}} = 1$. The remaining factors above can be derived as follows. First, note that $\mathbf{W}_{\mathbf{U}}$ spans the top- k generalized eigenspace of the matrix pair $(\mathbf{A}_{\phi}^{\text{H}}, \mathbf{B}_{\phi})$ which has \mathbf{B}_{ϕ} -orthonormal basis $\tilde{\mathbf{U}}$ (see Lemma I in Section 2). In fact, by the Schur decomposition of \mathbf{A}_{ϕ} , it holds that

$$\begin{aligned} \mathbf{A}_{\phi}^{\text{H}} \mathbf{W}_{\mathbf{U}} &= \mathbf{B}_{\phi} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{\Sigma}} & \mathbf{0}^{\text{H}} \\ \tilde{\mathbf{\Sigma}}^{\text{H}} & \tilde{\mathbf{\Sigma}}_{\perp}^{\text{H}} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{\Xi} \end{pmatrix}^{\text{H}} = \mathbf{B}_{\phi} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{\Sigma}} \\ (\tilde{\mathbf{\Sigma}} - \mathbf{\Xi} \tilde{\mathbf{\Sigma}}_{\perp})^{\text{H}} \end{pmatrix} \\ &= \mathbf{B}_{\phi} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{\Sigma}} \\ (-\tilde{\mathbf{\Sigma}} \mathbf{\Xi})^{\text{H}} \end{pmatrix} = \mathbf{B}_{\phi} \mathbf{W}_{\mathbf{U}} \tilde{\mathbf{\Sigma}}. \end{aligned}$$

Letting $\mathbf{G} = (\mathbf{I} + \mathbf{\Xi} \mathbf{\Xi}^{\text{H}})^{\frac{1}{2}}$, we have that $\|\mathbf{G}^{-\text{H}}\|_2 \leq 1$ due to $\mathbf{G} \succcurlyeq \mathbf{I}$, and $\mathbf{W}_{\mathbf{U}} \mathbf{G}^{-1} = \tilde{\mathbf{U}} \mathbf{O}$ for a certain unitary matrix $\mathbf{O} \in \mathbb{C}^{k \times k}$. Consequently,

$$\begin{aligned} \|\tilde{\mathbf{Z}}_0\|_2 &= \|\tilde{\mathbf{Y}}_0 \tilde{\mathbf{X}}_0^{-1}\|_2 \leq \|\tilde{\mathbf{X}}_0^{-1}\|_2 = \|((\mathbf{W}_{\mathbf{U}} \mathbf{G}^{-1})^{\text{H}} \mathbf{B}_{\phi} \mathbf{P}_0)^{-1} \mathbf{G}^{-\text{H}}\|_2 \\ &\leq \|(\tilde{\mathbf{U}}^{\text{H}} \mathbf{B}_{\phi} \mathbf{P}_0)^{-1}\|_2 = \|(\mathbf{D}(\beta) \mathbf{U}^{\text{H}} \mathbf{C}_{xx} \mathbf{\Phi}_0)^{-1}\|_2 = \frac{1}{\sigma_{\min}(\mathbf{D}(\beta) \mathbf{U}^{\text{H}} \mathbf{C}_{xx} \mathbf{\Phi}_0)} \\ &\leq \frac{1}{\sigma_{\min}(\mathbf{D}(\beta)) \sigma_{\min}(\mathbf{U}^{\text{H}} \mathbf{C}_{xx} \mathbf{\Phi}_0)} \leq \frac{\sqrt{\beta^2 + (\sigma_k^+)^2}}{\sigma_k^+ \sigma_{\min}(\mathbf{U}^{\text{H}} \mathbf{C}_{xx} \mathbf{\Phi}_0)}, \end{aligned}$$

where $\sigma_{\min}(\cdot)$ represents the minimum singular value of a matrix. Similarly,

$$\begin{aligned} \|\tilde{\mathbf{X}}_t^{-1}\|_2 &\leq \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1}\|_2 \\ &\leq \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1}(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)\|_2 \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1}\|_2 = \frac{a_{\phi,t}}{\cos \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}, \end{aligned}$$

where $a_{\phi,t} = \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1}(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)\|_2$ and Lemma 12 in Ge et al. (2016) is used in the last equality. Noting that $\|\mathbf{C}\|_2 \leq \sigma_1 \leq 1$ and $\|\tilde{\mathbf{R}}_t\|_2 \leq 1$ as $\mathbf{0} \prec \tilde{\mathbf{R}}_t \prec \mathbf{I}$, it holds that

$$\begin{aligned} \|\delta_t\|_{\mathbf{B}_\phi,2} &= \left\| \mathbf{B}_\phi^{\frac{1}{2}} \begin{pmatrix} \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \hat{\boldsymbol{\xi}}_t + \hat{\boldsymbol{\xi}}_t \\ \mathbf{0} \end{pmatrix} \tilde{\mathbf{R}}_t \right\|_2 \leq (\|\mathbf{C}\|_2 \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{yy},2} + \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},2}) \|\tilde{\mathbf{R}}_t\|_2 \\ &\leq \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{yy},2} + \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},2}, \end{aligned}$$

$$\begin{aligned} \|(\tilde{\boldsymbol{\Sigma}} + \mathbf{W}_U^H \mathbf{B}_\phi \delta_t \tilde{\mathbf{X}}_t^{-1})^{-1}\|_2 &\leq \frac{1}{\sigma_{\min}(\tilde{\boldsymbol{\Sigma}}) - \|\mathbf{W}_U\|_{\mathbf{B}_\phi,2} \|\delta\|_{\mathbf{B}_\phi,2} \|\tilde{\mathbf{X}}_t^{-1}\|_2} \\ &\leq \frac{1}{\sigma_k^+ - a_{\phi,t}(1 + \|\boldsymbol{\Xi}\|_2)(\|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{yy},2} + \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},2}) / \cos \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}. \end{aligned}$$

Also, we have that

$$\begin{aligned} \|\boldsymbol{\Gamma} \boldsymbol{\Pi} \tilde{\boldsymbol{\Sigma}}_{\perp} \boldsymbol{\Pi}^T \boldsymbol{\Gamma}^{-1}\|_2 &= \max_{k+1 \leq j \leq r} \|\boldsymbol{\Sigma}_j^{(\gamma)}\|_2 \leq \max_{k+1 \leq j \leq r} (|\sigma_j^+| + \frac{1+\beta}{1+\gamma}) \\ &= \sqrt{\beta} + \frac{1+\beta}{1+\gamma} = \sqrt{\beta} + \frac{\sigma_k^+ - \sqrt{\beta}}{8}, \end{aligned}$$

and that

$$\begin{aligned} \|\tilde{\mathbf{X}}_t^{-1}(\tilde{\mathbf{Z}}_t^H \tilde{\mathbf{Z}}_t)^{-1} \tilde{\mathbf{Z}}_t^H\|_2 &= \|\tilde{\mathbf{X}}_t^{-1}((\tilde{\mathbf{Y}}_t \tilde{\mathbf{X}}_t^{-1})^H (\tilde{\mathbf{Y}}_t \tilde{\mathbf{X}}_t^{-1}))^{-1} (\tilde{\mathbf{Y}}_t \tilde{\mathbf{X}}_t^{-1})^H\|_2 = \|(\tilde{\mathbf{Y}}_t^H \tilde{\mathbf{Y}}_t)^{-1} \tilde{\mathbf{Y}}_t^H\|_2 \\ &= \|(\tilde{\mathbf{Y}}_t^H \tilde{\mathbf{Y}}_t)^{-\frac{1}{2}}\|_2 \|(\tilde{\mathbf{Y}}_t^H \tilde{\mathbf{Y}}_t)^{-\frac{1}{2}} \tilde{\mathbf{Y}}_t^H\|_2 = \|(\tilde{\mathbf{Y}}_t^H \tilde{\mathbf{Y}}_t)^{-\frac{1}{2}}\|_2 \\ &= \|((\tilde{\mathbf{U}}_{\perp}^H \mathbf{B}_\phi \mathbf{P}_t)^H (\tilde{\mathbf{U}}_{\perp}^H \mathbf{B}_\phi \mathbf{P}_t))^{-\frac{1}{2}}\|_2 = \|(\mathbf{I} - \mathbf{P}_t^H \mathbf{B}_\phi \tilde{\mathbf{U}} \tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-\frac{1}{2}}\|_2 \\ &= \frac{1}{\sqrt{1 - \sigma_{\max}^2(\mathbf{P}_t^H \mathbf{B}_\phi \tilde{\mathbf{U}})}} = \frac{1}{\sin \theta_{\min}(\mathbf{P}_t, \tilde{\mathbf{U}})} = \frac{c_{\phi,t}}{c_{\phi,t} \sin \theta_{\min}(\mathbf{P}_t, \tilde{\mathbf{U}})} \\ &\leq \frac{c_{\phi,t}}{\sin(c_{\phi,t} \theta_{\min}(\mathbf{P}_t, \tilde{\mathbf{U}}))} = \frac{c_{\phi,t}}{\sin \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}, \end{aligned}$$

where $c_{\phi,t} = \frac{\theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{\theta_{\min}(\mathbf{P}_t, \tilde{\mathbf{U}})}$, the inequality is by the fact that $|\sin(nx)| \leq n \sin(x)$ for any natural number and real x , and the fifth equality is due to that $\begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix}$ is unitary in non-Euclidean metric \mathbf{B}_ϕ , i.e.,

$$\begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix}^H \mathbf{B}_\phi \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} = \mathbf{B}_\phi^{\frac{1}{2}} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_{\perp} \end{pmatrix}^H \mathbf{B}_\phi^{\frac{1}{2}} = \mathbf{I},$$

and thus

$$\mathbf{B}_\phi^{\frac{1}{2}} \tilde{\mathbf{U}}_{\perp} \tilde{\mathbf{U}}_{\perp}^H \mathbf{B}_\phi^{\frac{1}{2}} = \mathbf{I} - \mathbf{B}_\phi^{\frac{1}{2}} \tilde{\mathbf{U}} \tilde{\mathbf{U}}^H \mathbf{B}_\phi^{\frac{1}{2}}.$$

If

$$\max \{ \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{yy},F}, \|\hat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},F} \} \leq \frac{\sigma_k^+ - \sqrt{\beta}}{16 \max \{1 + \gamma, 1 + \|\boldsymbol{\Xi}\|_2\}} \min \left\{ \frac{\cos \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{a_{\phi,t}}, \frac{\sin \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{c_{\phi,t}} \right\},$$

we then get from the results derived above that

$$\begin{aligned}
\sin \theta_{\max}(\mathbf{P}_T, \tilde{\mathbf{U}}) &\leq (1+\gamma)(1+\|\Xi\|_2) \frac{\sqrt{\beta^2 + (\sigma_k^+)^2}/\sigma_k^+}{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0)} \times \\
&\quad \prod_{t=0}^{T-1} \frac{\sqrt{\beta} + \frac{\sigma_k^+ - \sqrt{\beta}}{4} + \frac{2c_{\phi,t}(1+\gamma) \max\{\|\xi_t\|_{\mathbf{C}_{yy,2}}, \|\hat{\xi}_t\|_{\mathbf{C}_{xx,2}}\}}{\sin \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}}{\sigma_k^+ - \frac{2a_{\phi,t}(1+\|\Xi\|_2) \max\{\|\xi_t\|_{\mathbf{C}_{yy,2}}, \|\hat{\xi}_t\|_{\mathbf{C}_{xx,2}}\}}{\cos \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}}} \\
&\leq (1+\gamma)(1+\|\Xi\|_2) \frac{\sqrt{\beta^2 + (\sigma_k^+)^2}/\sigma_k^+}{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0)} \left(\frac{\sqrt{\beta} + \frac{\sigma_k^+ - \sqrt{\beta}}{4}}{\sigma_k^+ - \frac{\sigma_k^+ - \sqrt{\beta}}{4}} \right)^T,
\end{aligned}$$

where

$$\begin{aligned}
\left(\frac{\sqrt{\beta} + \frac{\sigma_k^+ - \sqrt{\beta}}{4}}{\sigma_k^+ - \frac{\sigma_k^+ - \sqrt{\beta}}{4}} \right)^T &= \left(\frac{\sigma_k^+ + 3\sqrt{\beta}}{3\sigma_k^+ + \sqrt{\beta}} \right)^T \leq \left(1 - \frac{2(\sigma_k^+ - \sqrt{\beta})}{3\sigma_k^+ + \sqrt{\beta}} \right)^T \\
&\leq \left(1 - \frac{\sigma_k^+ - \sqrt{\beta}}{2\sigma_k^+} \right)^T \leq \exp\left\{ -\frac{\sigma_k^+ - \sqrt{\beta}}{\sigma_k^+} \frac{T}{2} \right\},
\end{aligned}$$

and

$$\begin{aligned}
\frac{\sigma_k^+ - \sqrt{\beta}}{\sigma_k^+} &= \frac{\sigma_k^2 - 2\sqrt{\beta} + \sqrt{\sigma_k^4 - 4\beta}}{\sigma_k^2 + \sqrt{\sigma_k^4 - 4\beta}} = \frac{\sqrt{\sigma_k^2 - 2\sqrt{\beta}}(\sqrt{\sigma_k^2 - 2\sqrt{\beta}} + \sqrt{\sigma_k^2 + 2\sqrt{\beta}})}{\sigma_k^2 + \sqrt{\sigma_k^4 - 4\beta}} \\
&\geq \frac{\sqrt{\sigma_k^2 - 2\sqrt{\beta}} \cdot \sigma_k}{2\sigma_k^2} = \frac{1}{2} \sqrt{\frac{\sigma_k^2 - 2\sqrt{\beta}}{\sigma_k^2}} = \frac{1}{2} \sqrt{\frac{\sigma_k^2 - \sigma_{k+1}^2}{\sigma_k^2}}.
\end{aligned}$$

Thus, we have that

$$\begin{aligned}
\sin \theta_{\max}(\mathbf{P}_T, \tilde{\mathbf{U}}) &\leq \frac{(1+\gamma)(1+\|\Xi\|_2) \sqrt{\beta^2 + (\sigma_k^+)^2}/\sigma_k^+}{\min\{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0), \sigma_{\min}(\mathbf{V}^H \mathbf{C}_{yy} \Psi_0)\}} \times \\
&\quad \exp\left\{ -\sqrt{\frac{\sigma_k^2 - \sigma_{k+1}^2}{\sigma_k^2}} \cdot \frac{T}{4} \right\} \triangleq \frac{\sigma_k^+}{k\sqrt{1 + (\sigma_k^+)^2}} \epsilon,
\end{aligned}$$

and thus $\sin \theta_{\max}(\Phi_T, \mathbf{U}) < \epsilon$ by Lemma 3.4. Solving the last equation for T yields that

$$T = 4 \sqrt{\frac{\sigma_k^2}{\sigma_k^2 - \sigma_{k+1}^2}} \log \frac{k(1+\gamma)(1+\|\Xi\|_2) \sqrt{(1 + (\sigma_k^+)^2)(\beta^2 + (\sigma_k^+)^2)}}{(\sigma_k^+)^2 \epsilon \cos \theta_0},$$

where $\theta_0 = \max\{\theta_{\max}(\Phi_0, \mathbf{U}), \theta_{\max}(\Psi_0, \mathbf{V})\}$. By $\tilde{\Sigma}\Xi - \Xi\tilde{\Sigma}_{\perp} = -\tilde{\Sigma}$, we have that (Golub and Van Loan, 2013)

$$\|\tilde{\Sigma}\|_F = \|\tilde{\Sigma}\Xi - \Xi\tilde{\Sigma}_{\perp}\|_F = \|\Xi\|_F \frac{\|\tilde{\Sigma}\Xi - \Xi\tilde{\Sigma}_{\perp}\|_F}{\|\Xi\|_F} \geq \|\Xi\|_F \text{sep}_F(\tilde{\Sigma}, \tilde{\Sigma}_{\perp}),$$

and that

$$\begin{aligned}
\text{sep}_F(\tilde{\Sigma}, \tilde{\Sigma}_{\perp}) &\leq \min|\lambda(\tilde{\Sigma}) - \lambda(\tilde{\Sigma}_{\perp})| \leq \sigma_k^+ - \sigma_k^- = \sqrt{\sigma_k^4 - 4\beta} \\
&= \sqrt{\sigma_k^4 - \sigma_{k+1}^4} = \sqrt{(\sigma_k^2 + \sigma_{k+1}^2)(\sigma_k^2 - \sigma_{k+1}^2)},
\end{aligned}$$

where $\text{sep}_F(\cdot, \cdot)$ represents the separation between two matrices in Frobenius norm and $\lambda(\cdot)$ represent an arbitrary eigenvalue of a matrix. Thus, we have that $\|\Xi\|_2 = O(\frac{1+\beta}{\sqrt{\sigma_k^2 - \sigma_{k+1}^2}})$. Besides, $1 + \gamma = \frac{1+\beta}{8(\sigma_k^+ - \sqrt{\beta})} \leq \frac{1+\beta}{4\sigma_k \sqrt{\sigma_k^2 - \sigma_{k+1}^2}}$. Thus, we can write that

$$T = O\left(\sqrt{\frac{\sigma_k^2}{\sigma_k^2 - \sigma_{k+1}^2}} \log \frac{1}{(\sigma_k^2 - \sigma_{k+1}^2)\epsilon \cos \theta_0}\right).$$

Analogously, we also have $\sin \theta_{\max}(\mathbf{Q}_T, \tilde{\mathbf{V}}) < \frac{\sigma_k^+}{k\sqrt{1+(\sigma_k^+)^2}}\epsilon$, where $\mathbf{Q}_t = \begin{pmatrix} \Psi_t \\ \Psi_{t-1} \mathbf{S}_t^{-1} \end{pmatrix} \tilde{\mathbf{S}}_t$, and thus $\sin \theta_{\max}(\Psi_T, \mathbf{V}) < \epsilon$. Therefore, we get that

$$\sin \theta_T = \sin \max\{\theta_{\max}(\Phi_T, \mathbf{U}), \theta_{\max}(\Psi_T, \mathbf{V})\} < \epsilon.$$

On the other hand, in each iteration, we need to solve four least-squares subproblems in order to make the error conditions on, e.g., ξ_t and $\hat{\xi}_t$, satisfied. This part of complexity can be obtained using Lemmas 3.2-3.3 as follows. First, by Lemma 3.2, the complexity of getting ξ_t as small as required previously is

$$O\left(\text{nnz}(\mathbf{X}) + \text{nnz}(\mathbf{Y}) \sqrt{\kappa(\mathbf{C}_{yy})} \log \frac{\epsilon_t(\Psi_t^{(0)})}{\epsilon_t(\hat{\Psi}_t)}\right),$$

and the ratio of the initial to final error can be further written as

$$\begin{aligned} \log \frac{\epsilon_t(\Psi_t^{(0)})}{\epsilon_t(\hat{\Psi}_t)} &= O\left(\log \frac{\tan^2 \max\{\theta_{\max}(\Phi_t, \mathbf{U}), \theta_{\max}(\Psi_t, \mathbf{V})\}}{\|\xi_t\|_{\mathbf{C}_{yy}, F}^2}\right) \\ &= O\left(\log \frac{c_1^2 \tan^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{(\sigma_k^+ - \sqrt{\beta})^2 \min\left\{\frac{\cos^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{a_{\phi, t}^2}, \frac{\sin^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{c_{\phi, t}^2}\right\}}\right) \\ &= O\left(\log \left(\frac{a_{\phi, t} c_{\phi, t} c_1}{\sigma_k^+ - \sigma_{k+1}^2} \max\left\{\frac{\sin^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{\cos^4 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}, \frac{1}{\cos^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}\right\}\right)\right), \end{aligned}$$

where the numerator in the second equation is by c_1 's definition in the theorem, and $O\left(\log \max\left\{\frac{\sin^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}{\cos^4 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}, \frac{1}{\cos^2 \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})}\right\}\right) = O\left(\log \frac{1}{\cos \theta_{\max}(\mathbf{P}_0, \tilde{\mathbf{U}})}\right)$ when $\theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})$ is large, otherwise $O(1)$. Also, we have that

$$\begin{aligned} a_{\phi, t} &= \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1} (\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)\|_2 \leq \|(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_t)^{-1}\|_2 \\ &= \frac{1}{\cos \theta_{\max}(\mathbf{P}_t, \tilde{\mathbf{U}})} \leq \frac{1}{\cos \theta_{\max}(\mathbf{P}_0, \tilde{\mathbf{U}})} = \frac{1}{\sigma_{\min}(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_0)} \\ &\leq \frac{\sqrt{\beta^2 + (\sigma_k^+)^2} / \sigma_k^+}{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0)} \quad (\text{see Page 3 on } \|\tilde{\mathbf{Z}}_0\|_2), \end{aligned}$$

and similarly,

$$\frac{1}{\cos \theta_{\max}(\mathbf{P}_0, \tilde{\mathbf{U}})} = \frac{1}{\sigma_{\min}(\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_0)} \leq \frac{\sqrt{1 + (\sigma_k^+)^2} / \sigma_k^+}{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0)}.$$

Thus, we can write that

$$\begin{aligned} \log \frac{\epsilon_t(\Psi_t^{(0)})}{\epsilon_t(\hat{\Psi}_t)} &= O\left(\log \frac{c_1 c_2}{(\sigma_k^2 - \sigma_{k+1}^2) \sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0)}\right) \\ &= O\left(\log \frac{c_1 c_2}{(\sigma_k^2 - \sigma_{k+1}^2) \min\{\sigma_{\min}(\mathbf{U}^H \mathbf{C}_{xx} \Phi_0), \sigma_{\min}(\mathbf{V}^H \mathbf{C}_{yy} \Psi_0)\}}\right), \end{aligned}$$

where $c_2 = \max\{\max_t\{c_{\phi,t}\}, \max_t\{c_{\psi,t}\}\}$. Similarly, by Lemma 3.3, the complexity of getting $\widehat{\boldsymbol{\xi}}_t$ as small as required previously is

$$O\left(\text{nnz}(\mathbf{Y}) + \text{nnz}(\mathbf{X})\sqrt{\kappa(\mathbf{C}_{xx})} \log \frac{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t^{(0)})}{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t)}\right),$$

and $\log \frac{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t^{(0)})}{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t)}$ can be further written as

$$\begin{aligned} & O\left(\log \frac{(\sigma_1^2 + \|\boldsymbol{\xi}_t\|_{\mathbf{C}_{yy},F}^2) \tan^2 \max\{\theta_{\max}(\widehat{\boldsymbol{\Phi}}_t, \mathbf{U}), \theta_{\max}(\widehat{\boldsymbol{\Psi}}_t, \mathbf{V})\}}{\|\widehat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},F}^2}\right) \\ &= O\left(\log \frac{\tan^2 \max\{\theta_{\max}(\widehat{\boldsymbol{\Phi}}_t, \mathbf{U}), \theta_{\max}(\widehat{\boldsymbol{\Psi}}_t, \mathbf{V})\}}{\|\widehat{\boldsymbol{\xi}}_t\|_{\mathbf{C}_{xx},F}^2}\right), \end{aligned}$$

where we have used that $\sigma_1 \leq 1$ and $\|\boldsymbol{\xi}_t\|_{\mathbf{C}_{yy},F} \leq 1$. The cases for $\boldsymbol{\eta}_t$ and $\widehat{\boldsymbol{\eta}}_t$ are similar as well. We thus have the following overall complexity:

$$O\left(dk^2T + kT\text{nnz}(\mathbf{X}, \mathbf{Y})\kappa^{\frac{1}{2}}(\mathbf{X}, \mathbf{Y}) \log \frac{c_1c_2}{(\sigma_k^2 - \sigma_{k+1}^2) \cos \theta_0}\right),$$

where $d = \max\{d_x, d_y\}$, $\text{nnz}(\mathbf{X}, \mathbf{Y}) = \text{nnz}(\mathbf{X}) + \text{nnz}(\mathbf{Y})$, and $\kappa(\mathbf{X}, \mathbf{Y}) = \max\{\kappa(\mathbf{C}_{xx}), \kappa(\mathbf{C}_{yy})\}$. \square

Lemma 3.3 Consider the least-squares subproblem $\min_{\boldsymbol{\Phi}} \widehat{l}_t(\boldsymbol{\Phi})$, for which the minimizer and the objective sub-optimality gap can be expressed as $\widehat{\boldsymbol{\Phi}}_t^* = \mathbf{C}_{xx}^{-1}\mathbf{C}_{xy}\widehat{\boldsymbol{\Psi}}_t$ and $\widehat{\epsilon}_t(\boldsymbol{\Phi}) = \widehat{l}_t(\boldsymbol{\Phi}) - \widehat{l}_t(\widehat{\boldsymbol{\Phi}}_t^*) = \frac{1}{2}\|\boldsymbol{\Phi} - \widehat{\boldsymbol{\Phi}}_t^*\|_{\mathbf{C}_{xx},F}^2$. We have that

$$\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t^{(0)}) \leq 8k(\sigma_1^2 + \|\boldsymbol{\xi}_t\|_{\mathbf{C}_{yy},2}^2) \tan^2 \widehat{\theta}_t,$$

for the initial sub-optimality, and accelerated gradient descent takes $O(\text{nnz}(\mathbf{Y}) + \text{nnz}(\mathbf{X})\kappa^{\frac{1}{2}}(\mathbf{C}_{xx}) \log \frac{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t^{(0)})}{\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t)})$ complexity to get the final sub-optimality $\widehat{\epsilon}_t(\widehat{\boldsymbol{\Phi}}_t)$, where $\widehat{\boldsymbol{\Phi}}_t^{(0)} = \widehat{\boldsymbol{\Phi}}_t(\widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xx} \widehat{\boldsymbol{\Phi}}_t)^{-1}(\widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t)$, $\|\mathbf{A}\|_{\mathbf{B},2} = \|\mathbf{B}^{\frac{1}{2}}\mathbf{A}\|_2$, and $\widehat{\theta}_t = \max\{\theta_{\max}(\widehat{\boldsymbol{\Phi}}_t, \mathbf{U}), \theta_{\max}(\widehat{\boldsymbol{\Psi}}_t, \mathbf{V})\}$. Parallel results hold for $\min_{\boldsymbol{\Psi}} \widehat{h}_t(\boldsymbol{\Psi})$ as well.

Proof Noting that $\widehat{\boldsymbol{\Phi}}_t^* = \mathbf{C}_{xx}^{-1}\mathbf{C}_{xy}\widehat{\boldsymbol{\Psi}}_t$ and

$$\widehat{l}_t(\widehat{\boldsymbol{\Phi}}_t^*) = -\frac{1}{2}\text{tr}\left(\widehat{\boldsymbol{\Psi}}_t^\top \mathbf{C}_{xy}^\top \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t\right) + \frac{1}{2n}\|\mathbf{Y}^\top \widehat{\boldsymbol{\Psi}}_t\|_F^2.$$

we have that

$$\begin{aligned} \frac{1}{2}\|\boldsymbol{\Phi} - \widehat{\boldsymbol{\Phi}}_t^*\|_{\mathbf{C}_{xx},F}^2 &= \frac{1}{2}\text{tr}\left((\boldsymbol{\Phi} - \widehat{\boldsymbol{\Phi}}_t^*)^\top \mathbf{C}_{xx} (\boldsymbol{\Phi} - \widehat{\boldsymbol{\Phi}}_t^*)\right) \\ &= \text{tr}\left(\frac{1}{2}\boldsymbol{\Phi}^\top \mathbf{C}_{xx} \boldsymbol{\Phi} - \boldsymbol{\Phi}^\top \mathbf{C}_{xx} \widehat{\boldsymbol{\Phi}}_t^* + \frac{1}{2}(\widehat{\boldsymbol{\Phi}}_t^*)^\top \mathbf{C}_{xx} \widehat{\boldsymbol{\Phi}}_t^*\right) \\ &= \text{tr}\left(\frac{1}{2}\boldsymbol{\Phi}^\top \mathbf{C}_{xx} \boldsymbol{\Phi} - \boldsymbol{\Phi}^\top \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t + \frac{1}{2}\widehat{\boldsymbol{\Psi}}_t^\top \mathbf{C}_{xy}^\top \mathbf{C}_{xx} \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t\right) \\ &= \widehat{l}_t(\boldsymbol{\Phi}) - \widehat{l}_t(\widehat{\boldsymbol{\Phi}}_t^*) = \widehat{\epsilon}_t(\boldsymbol{\Phi}). \end{aligned}$$

Setting $\frac{\partial}{\partial \boldsymbol{\Gamma}} \widehat{l}_t(\widehat{\boldsymbol{\Phi}}_t \boldsymbol{\Gamma}) = \widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xx} \widehat{\boldsymbol{\Phi}}_t \boldsymbol{\Gamma} - \widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t = 0$ yields the optimal

$$\boldsymbol{\Gamma}^* = (\widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xx} \widehat{\boldsymbol{\Phi}}_t)^{-1} \widehat{\boldsymbol{\Phi}}_t^\top \mathbf{C}_{xy} \widehat{\boldsymbol{\Psi}}_t.$$

Thus, $\widehat{\Phi}_t^{(0)} = \widehat{\Phi}_t \Gamma^*$. Noting that $\mathbf{C}_{xy} = \mathbf{C}_{xx}(\mathbf{U}\Sigma\mathbf{V}^\top + \mathbf{U}_\perp \Sigma_\perp \mathbf{V}_\perp^\top) \mathbf{C}_{yy}$, it holds that

$$\begin{aligned}
 2\widehat{\epsilon}_t(\widehat{\Phi}_t^{(0)}) &\leq 2\widehat{\epsilon}_t(\widehat{\Phi}_t \widetilde{\Gamma}) = \|\widehat{\Phi}_t \widetilde{\Gamma} - \widehat{\Phi}_t^*\|_{\mathbf{C}_{xx}, F}^2 \\
 &= \|\mathbf{U}^\top \mathbf{C}_{xx}(\widehat{\Phi}_t \widetilde{\Gamma} - \widehat{\Phi}_t^*)\|_F^2 + \|\mathbf{U}_\perp^\top \mathbf{C}_{xx}(\widehat{\Phi}_t \widetilde{\Gamma} - \widehat{\Phi}_t^*)\|_F^2 \\
 &= \|\mathbf{U}_\perp^\top \mathbf{C}_{xx} \widehat{\Phi}_t \widetilde{\Gamma} - \mathbf{U}_\perp^\top \mathbf{C}_{xy} \widehat{\Psi}_t\|_F^2 \quad \left(\text{let } \widetilde{\Gamma} = (\mathbf{U}^\top \mathbf{C}_{xx} \widehat{\Phi}_t)^{-1} \mathbf{U}^\top \mathbf{C}_{xx} \widehat{\Phi}_t^*\right) \\
 &= \|\mathbf{U}_\perp^\top \mathbf{C}_{xx} \widehat{\Phi}_t (\mathbf{U}^\top \mathbf{C}_{xx} \widehat{\Phi}_t)^{-1} \Sigma \mathbf{V}^\top \mathbf{C}_{yy} \widehat{\Psi}_t - \Sigma_\perp^\top \mathbf{V}_\perp^\top \mathbf{C}_{yy} \widehat{\Psi}_t\|_F^2 \\
 &\leq 2k \left(\|\mathbf{U}_\perp^\top \mathbf{C}_{xx} \widehat{\Phi}_t (\mathbf{U}^\top \mathbf{C}_{xx} \widehat{\Phi}_t)^{-1}\|_2^2 \|\Sigma\|_2^2 \|\mathbf{V}^\top\|_{\mathbf{C}_{yy}, 2}^2 \|\widehat{\Psi}_t\|_{\mathbf{C}_{yy}, 2}^2 \right. \\
 &\quad \left. + \|\Sigma_\perp\|_2^2 \|\mathbf{V}_\perp^\top \mathbf{C}_{yy} \widehat{\Psi}_t (\widehat{\Psi}_t^\top \mathbf{C}_{yy} \widehat{\Psi}_t)^{-\frac{1}{2}}\|_2^2 \|(\widehat{\Psi}_t^\top \mathbf{C}_{yy} \widehat{\Psi}_t)^{\frac{1}{2}}\|_2^2 \right) \\
 &= 2k \left(\sigma_1^2 \tan^2 \theta_{\max}(\widehat{\Phi}_t, \mathbf{U}) + \sigma_{k+1}^2 \sin^2 \theta_{\max}(\widehat{\Psi}_t, \mathbf{V}) \right) \|\widehat{\Psi}_t\|_{\mathbf{C}_{yy}, 2}^2 \\
 &\leq 4k \|\widehat{\Psi}_t\|_{\mathbf{C}_{yy}, 2}^2 \tan^2 \max \{ \theta_{\max}(\widehat{\Phi}_t, \mathbf{U}), \theta_{\max}(\widehat{\Psi}_t, \mathbf{V}) \},
 \end{aligned}$$

where we have used that $\sigma_1 \leq 1$ and additionally,

$$\begin{aligned}
 \|\widehat{\Psi}_t\|_{\mathbf{C}_{yy}, 2}^2 &= \|\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top \Phi_t + \xi_t\|_{\mathbf{C}_{yy}, 2}^2 \\
 &\leq 2 \left(\|\mathbf{C}\|_2^2 \|\Phi_t\|_{\mathbf{C}_{xx}, 2}^2 + \|\xi_t\|_{\mathbf{C}_{yy}, 2}^2 \right) = 2 \left(\sigma_1^2 + \|\xi_t\|_{\mathbf{C}_{yy}, 2}^2 \right).
 \end{aligned}$$

Thus, we can write that

$$\widehat{\epsilon}_t(\widehat{\Phi}_t^{(0)}) \leq 8k \left(\sigma_1^2 + \|\xi_t\|_{\mathbf{C}_{yy}, 2}^2 \right) \tan^2 \max \{ \theta_{\max}(\widehat{\Phi}_t, \mathbf{U}), \theta_{\max}(\widehat{\Psi}_t, \mathbf{V}) \}.$$

The proof completes by noting that $\widehat{l}_t(\Phi)$ is $\lambda_{\max}(\mathbf{C}_{xx})$ -smooth and $\lambda_{\min}(\mathbf{C}_{xx})$ -strongly convex and using the complexity of Nesterov's accelerated gradient descent (Nesterov, 2014). The case for the least-squares subproblem $\min_{\Psi} \widehat{h}_t(\Psi)$ is analogous. \square

Lemma 3.4 If $\sin \max \{ \theta_{\max}(\mathbf{P}_T, \widetilde{\mathbf{U}}), \theta_{\max}(\mathbf{Q}_T, \widetilde{\mathbf{V}}) \} < \frac{\sigma_k^+ \epsilon}{k \sqrt{1 + (\sigma_k^+)^2}}$ where $\sigma_k^+ = \frac{\sigma_k^2 + \sqrt{\sigma_k^4 - 4\beta}}{2}$, it holds that $\sin \max \{ \theta_{\max}(\Phi_T, \mathbf{U}), \theta_{\max}(\Psi_T, \mathbf{V}) \} < \epsilon$.

Proof We only show that if $\sin \theta_{\max}(\mathbf{P}_T, \widetilde{\mathbf{U}}) < \frac{\sigma_k^+ \epsilon}{k \sqrt{1 + (\sigma_k^+)^2}}$ then it holds that $\sin \theta_{\max}(\Phi_T, \mathbf{U}) < \epsilon$, because the case of $\theta_{\max}(\mathbf{Q}_T, \widetilde{\mathbf{V}})$ is analogous and then it is easy to see the lemma holds.

Note that

$$\begin{aligned}
 \|\widetilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_T\|_F^2 &= k - \|\widetilde{\mathbf{U}}_\perp^H \mathbf{B}_\phi \mathbf{P}_T\|_F^2 \geq k - k \|\widetilde{\mathbf{U}}_\perp^H \mathbf{B}_\phi \mathbf{P}_T\|_2^2 \\
 &= k - k \sin^2 \theta_{\max}(\mathbf{P}_T, \widetilde{\mathbf{U}}) \geq k - \frac{(\sigma_k^+)^2 \epsilon^2}{k(1 + (\sigma_k^+)^2)},
 \end{aligned}$$

and

$$\begin{aligned}
 \|\widetilde{\mathbf{u}}_j^H \mathbf{B}_\phi \mathbf{P}_T\|_2^2 &= \left\| \begin{pmatrix} \mu_j(1) \mathbf{u}_j \\ \nu_j(1) \mathbf{u}_j \end{pmatrix}^\top \mathbf{B}_\phi \begin{pmatrix} \Phi_T \\ \Phi_{T-1} \mathbf{R}_T^{-1} \end{pmatrix} \widetilde{\mathbf{R}}_T \right\|_2^2 \\
 &= \left\| \begin{pmatrix} \mu_j(1) \Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_j \\ \nu_j(1) \Phi_{T-1}^\top \mathbf{C}_{xx} \mathbf{u}_j \end{pmatrix}^\top \begin{pmatrix} \mathbf{I} \\ \mathbf{R}_T^{-1} \end{pmatrix} \widetilde{\mathbf{R}}_T \right\|_2^2 \\
 &\leq \left\| \begin{pmatrix} \mu_j(1) \Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_j \\ \nu_j(1) \Phi_{T-1}^\top \mathbf{C}_{xx} \mathbf{u}_j \end{pmatrix}^\top \right\|_2^2 = \mu_j(1)^2 \|\Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_j\|_2^2 + \nu_j(1)^2 \|\Phi_{T-1}^\top \mathbf{C}_{xx} \mathbf{u}_j\|_2^2,
 \end{aligned}$$

for any $j = 1, \dots, k$, where $\tilde{\mathbf{u}}_j$ is the j -th column of $\tilde{\mathbf{U}}$, $\mu_j(\alpha)$ and $\nu_j(\alpha)$ are given in Lemma I of Section 2. If $\|\Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_{j'}\|_2^2 < 1 - \frac{\epsilon^2}{k}$ for some j' , there must be

$$\begin{aligned} \|\tilde{\mathbf{u}}_{j'}^H \mathbf{B}_\phi \mathbf{P}_T\|_2^2 &\leq \mu_{j'}(1)^2 \|\Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_{j'}\|_2^2 + \nu_{j'}(1)^2 \|\Phi_{T-1}^\top \mathbf{C}_{xx} \mathbf{u}_{j'}\|_2^2 \\ &< \mu_{j'}(1)^2 \left(1 - \frac{\epsilon^2}{k}\right) + \nu_{j'}(1)^2 = \mu_{j'}(1)^2 + \nu_{j'}(1)^2 - \mu_{j'}(1)^2 \frac{\epsilon^2}{k} \\ &< 1 - \mu_k(1)^2 \frac{\epsilon^2}{k} = 1 - \frac{(\sigma_k^+)^2 \epsilon^2}{k(1 + (\sigma_k^+)^2)}, \end{aligned}$$

and then

$$\|\tilde{\mathbf{U}}^H \mathbf{B}_\phi \mathbf{P}_T\|_F^2 = \sum_{j=1}^k \|\tilde{\mathbf{u}}_j^H \mathbf{B}_\phi \mathbf{P}_T\|_2^2 < 1 - \frac{(\sigma_k^+)^2 \epsilon^2}{k(1 + (\sigma_k^+)^2)} + k - 1 = k - \frac{(\sigma_k^+)^2 \epsilon^2}{k(1 + (\sigma_k^+)^2)},$$

contradiction. We thus have $\|\Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_j\|_2^2 > 1 - \frac{\epsilon^2}{k}$ for all $j = 1, \dots, k$ and then

$$\|\Phi_T^\top \mathbf{C}_{xx} \mathbf{U}\|_F^2 = \sum_{j=1}^k \|\Phi_T^\top \mathbf{C}_{xx} \mathbf{u}_j\|_2^2 > k \left(1 - \frac{\epsilon^2}{k}\right) = k - \epsilon^2.$$

We thus get that

$$\sin \theta_{\max}(\Phi_T, \mathbf{U}) < \|\sin \theta(\Phi_T, \mathbf{U})\|_2 = \|\mathbf{U}_\perp^\top \mathbf{C}_{xx} \Phi_T\|_F \leq \sqrt{k - \|\mathbf{U}^\top \mathbf{C}_{xx} \Phi_T\|_F^2} = \epsilon.$$

□

2 Auxiliary Lemma

Lemma I. \mathbf{A}_ϕ and \mathbf{A}_ϕ^H have the following Schur decompositions in non-Euclidean metric \mathbf{B}_ϕ :

$$\begin{aligned} \mathbf{A}_\phi &= \mathbf{B}_\phi \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\Sigma} & \tilde{\Sigma} \\ \mathbf{0} & \tilde{\Sigma}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^H \mathbf{B}_\phi, \\ \mathbf{A}_\phi^H &= \mathbf{B}_\phi \begin{pmatrix} \tilde{\tilde{\mathbf{U}}} & \tilde{\tilde{\mathbf{U}}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\tilde{\Lambda}} & \tilde{\tilde{\Lambda}} \\ \mathbf{0} & \tilde{\tilde{\Lambda}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\tilde{\mathbf{U}}} & \tilde{\tilde{\mathbf{U}}}_\perp \end{pmatrix}^H \mathbf{B}_\phi, \end{aligned}$$

respectively, where both $\begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}$ and $\begin{pmatrix} \tilde{\tilde{\mathbf{U}}} & \tilde{\tilde{\mathbf{U}}}_\perp \end{pmatrix}$ are \mathbf{B}_ϕ -unitary, $\tilde{\mathbf{U}}_\perp$ represents $\tilde{\mathbf{U}}$'s \mathbf{B}_ϕ -orthogonal complement, and $\sigma_j^\pm = \frac{\sigma_j^2 \pm \sqrt{\sigma_j^4 - 4\beta}}{2}$,

$$\tilde{\mathbf{U}} = \begin{pmatrix} \mathbf{U}\mathbf{D}(1) \\ \mathbf{U}\mathbf{J}(1) \end{pmatrix}, \quad \tilde{\tilde{\mathbf{U}}} = \begin{pmatrix} \mathbf{U}\mathbf{D}(\beta) \\ -\mathbf{U}\mathbf{J}(\beta) \end{pmatrix}, \quad \mu_j(\alpha) = \frac{\sigma_j^+}{\sqrt{\alpha^2 + (\sigma_j^+)^2}}, \quad \nu_j(\alpha) = \frac{\alpha}{\sqrt{\alpha^2 + (\sigma_j^+)^2}},$$

$$\mathbf{D}(\alpha) = \text{diag}(\mu_1(\alpha), \dots, \mu_k(\alpha)), \quad \mathbf{J}(\alpha) = \text{diag}(\nu_1(\alpha), \dots, \nu_k(\alpha)),$$

$$\tilde{\Sigma}_\perp = \begin{pmatrix} \Sigma_\perp^+ & \mathbf{L}(-1) & \mathbf{0} \\ \mathbf{0} & \Sigma_\perp^- & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma^- \end{pmatrix}, \quad \tilde{\tilde{\Lambda}}_\perp = \begin{pmatrix} \Sigma_\perp^+ & \mathbf{L}(\beta) & \mathbf{0} \\ \mathbf{0} & \Sigma_\perp^- & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma^- \end{pmatrix},$$

$$\tilde{\Sigma} = \tilde{\tilde{\Lambda}} = \Sigma^+, \quad \tilde{\Sigma} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & -(1 + \beta)\mathbf{I} \end{pmatrix}, \quad \tilde{\tilde{\Lambda}} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & (1 + \beta)\mathbf{I} \end{pmatrix},$$

$$\Sigma^+ = \text{diag}(\sigma_1^+, \dots, \sigma_k^+), \quad \Sigma_\perp^+ = \text{diag}(\sigma_{k+1}^+, \dots, \sigma_{2d_x}^+), \quad \mathbf{L}(\alpha) = \alpha \mathbf{I} + \frac{1}{\alpha} (\Sigma_\perp^-)^2,$$

$$\Sigma^- = \text{diag}(\sigma_1^-, \dots, \sigma_k^-), \quad \Sigma_\perp^- = \text{diag}(\sigma_{k+1}^-, \dots, \sigma_{2d_x}^-).$$

If $\sigma_j = 2\sqrt{\beta}$ then corresponding entries in blocks $\mathbf{L}(-1)$ and $\mathbf{L}(\beta)$ are replaced with $(1 + \beta)$.

Proof We have that

$$\mathbf{A}_\phi = \begin{pmatrix} \mathbf{C}_{xy}\mathbf{C}_{yy}^{-1}\mathbf{C}_{xy}^\top & -\beta\mathbf{C}_{xx} \\ \mathbf{C}_{xx} & \mathbf{0} \end{pmatrix}, \quad \mathbf{B}_\phi = \begin{pmatrix} \mathbf{C}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{xx} \end{pmatrix},$$

where $\mathbf{C}_{xy}\mathbf{C}_{yy}^{-1}\mathbf{C}_{xy}^\top = \mathbf{C}_{xx}(\mathbf{U}\boldsymbol{\Sigma}^2\mathbf{U}^\top + \mathbf{U}_\perp\boldsymbol{\Sigma}_\perp^2\mathbf{U}_\perp^\top)\mathbf{C}_{xx}$. We can assume that $(\mathbf{U} \ \mathbf{U}_\perp) = (\mathbf{u}_1, \dots, \mathbf{u}_{d_x})$ are orthogonal in metric \mathbf{C}_{xx} with $\boldsymbol{\Sigma}_\perp = \text{diag}(\sigma_{k+1}, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{(d_x-k) \times (d_x-k)}$, where $r = \text{rank}(\mathbf{C}_{xy})$. Similar to Proposition 9 in Xu et al. (2018), $(\mathbf{A}_\phi, \mathbf{B}_\phi)$'s generalized eigenpairs can be written as $(\sigma_j^\pm, \mathbf{u}_j^\pm)$ where $\sigma_j^\pm = \frac{\sigma_j^2 \pm \sqrt{\sigma_j^4 - 4\beta}}{2}$ and $\mathbf{u}_j^\pm = \begin{pmatrix} \sigma_j^\pm \mathbf{u}_j \\ \mathbf{u}_j \end{pmatrix}$. Note that when $\sigma_j^2 = 2\sqrt{\beta}$, i.e., σ_j^\pm with algebraic multiplicity of 2, \mathbf{u}_j^\pm supplies only one generalized eigenvector \mathbf{u}_j^+ , i.e., σ_j^+ with geometric multiplicity of only 1. One ‘‘generalized’’ generalized eigenvector of $(\mathbf{A}_\phi, \mathbf{B}_\phi)$ corresponding to $\sigma_j^+ = \sqrt{\beta}$, i.e., the solution to the equation $(\mathbf{A}_\phi - \sqrt{\beta}\mathbf{B}_\phi)\mathbf{z}_j = \mathbf{B}_\phi \begin{pmatrix} \sqrt{\beta}\mathbf{u}_j \\ \mathbf{u}_j \end{pmatrix}$ in \mathbf{z}_j , is needed. It is easy to see $\mathbf{z}_j = \begin{pmatrix} \mathbf{u}_j \\ \mathbf{0} \end{pmatrix}$. For notational convenience, denote $\mathbf{u}_j^- = \mathbf{z}_j$ in this case. $(\mathbf{A}_\phi, \mathbf{B}_\phi)$'s $2d_x$ generalized eigenvectors or ‘‘generalized’’ generalized eigenvectors now can span \mathbb{C}^{2d_x} in metric \mathbf{B}_ϕ and thus we can write that

$$\mathbf{A}_\phi(\mathbf{u}_1^+, \mathbf{u}_1^-, \dots, \mathbf{u}_n^+, \mathbf{u}_n^-) = \mathbf{B}_\phi(\mathbf{u}_1^+, \mathbf{u}_1^-, \dots, \mathbf{u}_n^+, \mathbf{u}_n^-) \text{diag} \left(\begin{pmatrix} \sigma_1^+ & \delta_1 \\ 0 & \sigma_1^- \end{pmatrix}, \dots, \begin{pmatrix} \sigma_n^+ & \delta_n \\ 0 & \sigma_n^- \end{pmatrix} \right),$$

where $\delta_j = 0$ if $\sigma_n^+ \neq \sigma_n^-$ otherwise 1. Letting $(\mathbf{u}_j^+, \mathbf{u}_j^-) = (\hat{\mathbf{u}}_{2j-1}, \hat{\mathbf{u}}_{2j})\hat{\mathbf{R}}_j$ representing \mathbf{B}_ϕ -orthonormalization of $(\mathbf{u}_j^+, \mathbf{u}_j^-)$ in \mathbb{C}^{2n} , we then can write that

$$\begin{aligned} \mathbf{A}_\phi &= \mathbf{B}_\phi(\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_{2d_x}) \text{diag} \left(\hat{\mathbf{R}}_1 \begin{pmatrix} \sigma_1^+ & \delta_1 \\ 0 & \sigma_1^- \end{pmatrix} \hat{\mathbf{R}}_1^{-1}, \dots, \hat{\mathbf{R}}_n \begin{pmatrix} \sigma_n^+ & \delta_n \\ 0 & \sigma_n^- \end{pmatrix} \hat{\mathbf{R}}_n^{-1} \right) (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_{2d_x})^H \mathbf{B}_\phi \\ &\triangleq \mathbf{B}_\phi \hat{\mathbf{U}} \hat{\boldsymbol{\Sigma}} \hat{\mathbf{U}}^H \mathbf{B}_\phi, \end{aligned}$$

where $\hat{\mathbf{U}}$ is \mathbf{B}_ϕ -unitary and thus $\hat{\mathbf{U}}\hat{\mathbf{U}}^H\mathbf{B}_\phi = \mathbf{I}$. After some algebraic manipulations and permutations, we will arrive at $\hat{\mathbf{U}}\hat{\boldsymbol{\Sigma}}\hat{\mathbf{U}}^H = (\tilde{\mathbf{U}}, \tilde{\mathbf{U}}_\perp)\tilde{\boldsymbol{\Sigma}}(\tilde{\mathbf{U}}, \tilde{\mathbf{U}}_\perp)^H$, where $(\tilde{\mathbf{U}}, \tilde{\mathbf{U}}_\perp)$ is \mathbf{B}_ϕ -unitary and

$$\tilde{\mathbf{U}} = \begin{pmatrix} \mathbf{U}\mathbf{D}(1) \\ \mathbf{U}\mathbf{J}(1) \end{pmatrix}, \quad \tilde{\boldsymbol{\Sigma}} = \begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \tilde{\tilde{\boldsymbol{\Sigma}}} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix},$$

as described in the lemma. Thus, we have that

$$\mathbf{A}_\phi = \mathbf{B}_\phi \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\boldsymbol{\Sigma}} & \tilde{\tilde{\boldsymbol{\Sigma}}} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_\perp \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{U}} & \tilde{\mathbf{U}}_\perp \end{pmatrix}^H \mathbf{B}_\phi.$$

It is analogous for $(\mathbf{A}_\phi^H, \mathbf{B}_\phi)$, except for the generalized eigenpair being $(\sigma_j^\pm, \mathbf{u}_j^\pm)$ with $\mathbf{u}_j^\pm = \begin{pmatrix} \sigma_j^\pm \mathbf{u}_j \\ -\beta \mathbf{u}_j \end{pmatrix}$. \square

References

- Rong Ge, Chi Jin, Sham M. Kakade, Praneeth Netrapalli, and Aaron Sidford. Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, pages 2741–2750, New York City, NY, 2016.
- G.H. Golub and C.F. Van Loan. *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, 2013.
- Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. 2014.
- Peng Xu, Bryan D. He, Christopher De Sa, Ioannis Mitliagkas, and Christopher Ré. Accelerated stochastic power iteration. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 58–67, Playa Blanca, Lanzarote, Canary Islands, Spain, 2018.