

SUPPLEMENTAL MATERIAL FOR “PROCRUSTES ANALYSIS FOR HIGH-DIMENSIONAL  
DATA”

**A. Proof of Theorems and Lemmas**

*Theorem 1.* (Theobald & Wuttke (2006)) Consider the perturbation model described in Definition 2, and the singular value decomposition  $\mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} = \mathbf{U}_i \mathbf{D}_i \mathbf{V}_i^\top$ . So, the maximum likelihood estimators equal  $\hat{\mathbf{R}}_i = \mathbf{U}_i \mathbf{V}_i^\top$ , and  $\hat{\alpha}_i \hat{\mathbf{R}}_i = \|\boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2}\|^2 / \text{tr}(\mathbf{D}_i)$ .

*Proof.* The proof comes directly from Theobald & Wuttke (2006), we report here the final part that we will use for other proofs. We can write Equation (1) as:

$$\frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M} = \mathbf{E}_i \sim \mathcal{MN}(0, \boldsymbol{\Sigma}_n, \boldsymbol{\Sigma}_m).$$

The log-likelihood for  $\mathbf{R}_i$  equals:

$$\ell(\mathbf{R}_i) = -\frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \boldsymbol{\Sigma}_m^{-1} \left( \frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M} \right)^\top \boldsymbol{\Sigma}_n^{-1} \left( \frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M} \right) \right\} + C,$$

where  $(\mathbf{M}, \boldsymbol{\Sigma}_n, \boldsymbol{\Sigma}_m)$  are known nuisance parameter and  $C$  is a constant value. So, we have

$$\ell(\mathbf{R}_i) = \frac{1}{\alpha_i} \text{tr} \{ \mathbf{R}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} \} + C^*,$$

where  $C^*$  is a constant. The maximization of the log-likelihood function  $\ell(\mathbf{R}_i)$  leads to

$$\begin{aligned} \hat{\mathbf{R}}_i &= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{R}_i, \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} \rangle \\ &= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{R}_i, \mathbf{U}_i \mathbf{D}_i \mathbf{V}_i \rangle = \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{D}_i, \mathbf{U}_i^\top \mathbf{R}_i \mathbf{V}_i \rangle \\ &= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{D}_i, \mathbf{R}_i^o \rangle = \mathbf{U}_i \mathbf{V}_i^\top, \end{aligned} \tag{1}$$

where  $\mathbf{R}_i^o = \mathbf{U}_i^\top \mathbf{R}_i \mathbf{V}_i \in \mathcal{O}(m)$ . The step (1) is proved by Gower & Dijkstra (2004), i.e.,  $\langle \mathbf{D}_i, \mathbf{R}_i^o \rangle$  is maximum when  $\mathbf{R}_i^o = \mathbf{I}_m$ , giving  $\mathbf{I}_m = \mathbf{U}_i^\top \hat{\mathbf{R}}_i \mathbf{V}_i$ . We can note that the maximum likelihood estimator  $\hat{\mathbf{R}}_i$  does not depend on  $\alpha_i$ .

Consider the profile log-likelihood for  $\alpha_i$ :

$$\begin{aligned} \ell_p(\alpha_i) &= -\frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \boldsymbol{\Sigma}_m^{-1/2} \left( \frac{1}{\alpha_i} \mathbf{X}_i \hat{\mathbf{R}}_i - \mathbf{M} \right)^\top \boldsymbol{\Sigma}_n^{-1/2} \boldsymbol{\Sigma}_n^{-1/2} \left( \frac{1}{\alpha_i} \mathbf{X}_i \hat{\mathbf{R}}_i - \mathbf{M} \right) \boldsymbol{\Sigma}_m^{-1/2} \right\} + S \\ &= -\frac{1}{2} \left\{ \sum_{i=1}^N \frac{1}{\alpha_i^2} \text{tr} \left( \boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2} \boldsymbol{\Sigma}_n^{-1/2} \mathbf{X}_i \hat{\mathbf{R}}_i \boldsymbol{\Sigma}_m^{-1/2} \right) \right. \\ &\quad + \text{tr} \left( \boldsymbol{\Sigma}_m^{-1/2} \mathbf{M}^\top \boldsymbol{\Sigma}_n^{-1/2} \boldsymbol{\Sigma}_n^{-1/2} \mathbf{M} \boldsymbol{\Sigma}_m^{-1/2} \right) \\ &\quad \left. - 2 \frac{1}{\alpha_i} \text{tr} \left( \boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2} \boldsymbol{\Sigma}_n^{-1/2} \mathbf{M} \boldsymbol{\Sigma}_m^{-1/2} \right) \right\} + S \\ &= \sum_{i=1}^N -\frac{1}{2\alpha_i^2} \left\| \boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2} \right\|^2 + \frac{1}{\alpha_i} \text{tr}(\mathbf{D}_i) + S^*, \end{aligned}$$

where  $S$ , and  $S^*$  are constant values. Taking the first derivative, we have:

$$\begin{aligned} \frac{\partial \ell_p(\alpha_i)}{\partial \alpha_i} &= \alpha_i^{-1} \left\| \boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2} \right\|^2 - \text{tr}(\mathbf{D}_i) \\ \hat{\alpha}_i \hat{\mathbf{R}}_i &= \frac{\left\| \boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2} \right\|^2}{\text{tr}(\mathbf{D}_i)}, \end{aligned}$$

having  $\hat{\mathbf{R}}_i = \mathbf{U}_i \mathbf{V}_i^\top$  and  $\alpha_i \in \mathbb{R}^+$ .

*Lemma 2.* Consider the perturbation model of Definition 2, with  $\mathbf{R}_i$  distributed accordingly to (4), then the posterior distribution  $f(\mathbf{R}_i | k, \mathbf{F}, \mathbf{X}_i)$  is conjugate distribution to the von Mises-Fisher prior distribution with location posterior parameter equalling the following:

$$\mathbf{F}^* = \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} + k \mathbf{F}.$$

*Proof.* The joint posterior distribution is defined as

$$\begin{aligned}
\prod_{i=1}^N f(\mathbf{R}_i | \mathbf{X}_i, \mathbf{M}, \boldsymbol{\Sigma}_m, \boldsymbol{\Sigma}_n, k, \mathbf{F}) &= \prod_{i=1}^N \exp\left[-\frac{1}{2} \text{tr}\left\{\boldsymbol{\Sigma}_m^{-1/2} \left(\frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M}\right)^\top \boldsymbol{\Sigma}_n^{-1/2} \boldsymbol{\Sigma}_n^{-1/2} \left(\frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M}\right) \boldsymbol{\Sigma}_m^{-1/2}\right\}\right] \\
&\quad \cdot \exp\left\{k \text{tr}(\mathbf{F}^\top \mathbf{R}_i)\right\} \cdot C \\
&= \exp\left(-\sum_{i=1}^N \frac{1}{2} \Psi_i\right) \\
&\quad \cdot \exp\left(\sum_{i=1}^N \langle \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} + k \mathbf{F}, \mathbf{R}_i \rangle\right), \tag{2}
\end{aligned}$$

where  $\Psi_i = f(X_i)$  and  $C$  is a constant value. The quantity (2) is a kernel of a matrix von Mises-Fisher distribution with location parameter equals

$$\mathbf{F}^* = \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} + k \mathbf{F}.$$

*Theorem 2.* The ProMises model is defined as the perturbation model specified in Definition 2 imposing the prior distribution (8) for  $\alpha_i \mathbf{R}_i$ . Let the singular value decomposition of  $\mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1} + k \mathbf{F}$  be  $\mathbf{U}_i \mathbf{D}_i \mathbf{V}_i^\top$ . Then, the maximum a posteriori estimators equal  $\hat{\mathbf{R}}'_i = \mathbf{U}_i \mathbf{V}_i^\top$ , and  $\hat{\alpha}'_{\hat{\mathbf{R}}'_i} = \|\boldsymbol{\Sigma}_m^{-1/2} \hat{\mathbf{R}}'_i \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1/2}\|^2 / \text{tr}(\mathbf{D}_i)$ .

*Proof.* Consider the same assumption of Theorem 1, the log-posterior distribution for  $\mathbf{R}_i$  and  $\alpha_i$  equals:

$$\begin{aligned}
\log f(\alpha_i, \mathbf{R}_i | \mathbf{X}_i, \mathbf{M}, \boldsymbol{\Sigma}_m, \boldsymbol{\Sigma}_n, k, \mathbf{F}) &= -\frac{1}{2} \sum_{i=1}^N \text{tr}\left\{\boldsymbol{\Sigma}_m^{-1} \left(\frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M}\right)^\top \boldsymbol{\Sigma}_n^{-1} \left(\frac{1}{\alpha_i} \mathbf{X}_i \mathbf{R}_i - \mathbf{M}\right) \right. \\
&\quad \left. - 2 \frac{k}{\alpha_i} \mathbf{F}^\top \mathbf{R}_i\right\} - \log(\alpha_i) + K,
\end{aligned}$$

where  $K$  is a constant value. Following the same steps of Theorem 1's proof, the maximum a

posteriori estimate equals

$$\begin{aligned}
\hat{\mathbf{R}}'_i &= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{R}_i, \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{M} \Sigma_m^{-1} \rangle + k \langle \mathbf{R}_i, \mathbf{F} \rangle \\
&= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{R}_i, \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{M} \Sigma_m^{-1} + k \mathbf{F} \rangle \\
&= \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{R}_i, \mathbf{U}_i \mathbf{D}_i \mathbf{V}_i \rangle = \arg \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{D}_i, \mathbf{U}_i^\top \mathbf{R}_i \mathbf{V}_i \rangle \\
&= \max_{\mathbf{R}_i \in \mathcal{O}(m)} \langle \mathbf{D}_i, \mathbf{R}_i^o \rangle = \mathbf{U}_i \mathbf{V}_i^\top,
\end{aligned} \tag{3}$$

where step (3) is proved in the same way as step (1) of Theorem 1' proof.

Then we compute the maximum a posteriori estimate for  $\alpha_i$ :

$$\begin{aligned}
\hat{\alpha}'_{\hat{\mathbf{R}}'_i} &= \arg \max_{\alpha_i \in \mathbb{R}^+} -\frac{1}{2\alpha_i^2} \|\Sigma_m^{-1/2} \hat{\mathbf{R}}'^\top_i \mathbf{X}_i^\top \Sigma_n^{-1/2}\|^2 + \frac{1}{\alpha_i} \langle \hat{\mathbf{R}}'^\top_i, \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{M} \Sigma_m^{-1} \rangle \\
&\quad + \frac{k}{\alpha_i} \text{tr}(\mathbf{F}^\top \hat{\mathbf{R}}'_i) - \log(\alpha_i) + P,
\end{aligned}$$

where  $P$  is a constant value. Compute the first derivative and set it to zero:

$$\alpha_i^{-2} \|\Sigma_m^{-1/2} \hat{\mathbf{R}}'^\top_i \mathbf{X}_i^\top \Sigma_n^{-1/2}\|^2 - \alpha_i^{-1} \langle \hat{\mathbf{R}}'^\top_i, \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{M} \Sigma_m^{-1} + k \mathbf{F} \rangle - 1 = 0,$$

so, applying the Viéte theorem (Viéte, 1646),  $\hat{\alpha}'_{\hat{\mathbf{R}}'_i}$  equals:

$$\hat{\alpha}'_{\hat{\mathbf{R}}'_i} = \frac{\|\Sigma_m^{-1/2} \hat{\mathbf{R}}'^\top_i \mathbf{X}_i^\top \Sigma_n^{-1/2}\|^2}{\text{tr}(\mathbf{D}_i)},$$

under the condition  $\frac{\text{tr}(\mathbf{D}_i)}{\|\Sigma_m^{-1/2} \hat{\mathbf{R}}'^\top_i \mathbf{X}_i^\top \Sigma_n^{-1/2}\|^2} \gg \frac{1}{\text{tr}(\mathbf{D}_i)}$ , and  $\mathbf{D}_i$  coming from the singular value decomposition of  $\mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{M} \Sigma_m^{-1} + k \mathbf{F}$ .

*Lemma 4.* Consider  $\mathbf{X}_i \in \mathbb{R}^{n \times m}$ , if  $n < m$ , then the maximum likelihood estimate for  $\mathbf{R}_i$  defined in Theorem 1 is not unique.

*Proof.* In practice, without loss of generality, the Procrustes problem can be resumed as:

$$\max_{\mathbf{R}_i \in \mathcal{O}(m)} \text{tr}(\mathbf{A}_i^\top \mathbf{R}_i), \tag{4}$$

where  $\mathbf{A}_i = \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1}$ . Trendafilov & Lippert (2002) and Myronenko & Song (2009, Lemma 1) proved that the solution for (4) is unique if and only if the matrix  $\mathbf{A}_i$  has full rank. In Theorem 1 with  $n < m$ ,  $\mathbf{A}_i$  is equal to  $\mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{M} \boldsymbol{\Sigma}_m^{-1}$  having rank lower than  $m$ , so the solution is not unique. Please refer to Trendafilov & Lippert (2002) and Myronenko & Song (2009, Lemma 1) for further details about the complete proof.

*Theorem 3.* Consider the perturbation model in Definition 2 with  $\boldsymbol{\Sigma}_m = \sigma^2 \mathbf{I}_m$ , and the thin singular value decompositions of  $\mathbf{X}_i = \mathbf{L}_i \mathbf{S}_i \mathbf{Q}_i^\top$  for each  $i = 1, \dots, N$ , where  $\mathbf{Q}_i$  has dimensions  $n \times m$ . The following holds

$$\max_{\mathbf{R}_i \in \mathcal{O}(m)} \text{tr}(\mathbf{R}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{X}_j \boldsymbol{\Sigma}_m^{-1}) = \max_{\mathbf{R}_i^* \in \mathcal{O}(n)} \text{tr}(\mathbf{R}_i^{*\top} \mathbf{Q}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{X}_j \boldsymbol{\Sigma}_m^{-1} \mathbf{Q}_j^\top).$$

*Proof.* Without loss of generality we consider  $\boldsymbol{\Sigma}_m = \sigma^2 \mathbf{I}_m$ , and so the following objective function to maximize

$$\text{tr}(\mathbf{R}_i^\top \mathbf{X}_i^\top \boldsymbol{\Sigma}_n^{-1} \mathbf{X}_j \sigma^2 \mathbf{I}_m).$$

We note that it is equivalent to maximize

$$\text{tr}(\boldsymbol{\Sigma}_n^{-1} \mathbf{X}_i \mathbf{R}_i \mathbf{X}_j^\top) \tag{5}$$

thanks to the trace's proprieties.

Let consider the full singular value decomposition  $\mathbf{X}_i = \mathbf{L}_i \mathbf{S}_i \mathbf{C}_i^\top$ , where  $\mathbf{S}_i \in \mathbb{R}^{n \times m}$ . The  $\mathbf{S}_i$  matrix is defined as

$$\mathbf{S}_i = [\mathbf{S}_i^* \ \mathbf{O}],$$

where  $\mathbf{S}_i^* \in \mathbb{R}^{n \times n}$  and  $\mathbf{O}$  is a matrix of zero with  $n \times (m - n)$  dimensions, since  $\text{rank}(\mathbf{X}_i) = n \ \forall i = 1, \dots, N$ . Therefore, Expression (5) equals

$$\text{tr}(\boldsymbol{\Sigma}_n^{-1} \mathbf{X}_i \mathbf{R}_i \mathbf{X}_j^\top) = \text{tr}(\boldsymbol{\Sigma}_n^{-1} \mathbf{L}_i \mathbf{S}_i \mathbf{C}_i^\top \mathbf{R}_i \mathbf{C}_j \mathbf{S}_j^\top \mathbf{L}_j^\top) = \text{tr}(\boldsymbol{\Sigma}_n^{-1} \mathbf{L}_i \mathbf{S}_i \mathbf{R}_i^o \mathbf{S}_j^\top \mathbf{L}_j^\top),$$

where  $\mathbf{R}_i^o = \mathbf{C}_i^\top \mathbf{R}_i \mathbf{C}_j \in \mathcal{O}(m)$  being a product of orthogonal matrices.

Partitioning  $\mathbf{R}_i^o$  in blocks, i.e.,

$$\begin{bmatrix} \mathbf{R}_{11i}^o & \mathbf{R}_{12i}^o \\ \mathbf{R}_{21i}^o & \mathbf{R}_{22i}^o \end{bmatrix} \quad (6)$$

where  $\mathbf{R}_{11i}^o \in \mathbb{R}^{n \times n}$ ,  $\mathbf{R}_{12i}^o \in \mathbb{R}^{n \times m-n}$ ,  $\mathbf{R}_{21i}^o \in \mathbb{R}^{m-n \times n}$ , and  $\mathbf{R}_{22i}^o \in \mathbb{R}^{m-n \times m-n}$ , we have:

$$\begin{aligned} \Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i \mathbf{R}_i^o \mathbf{S}_j^\top \mathbf{L}_j^\top &= \Sigma_n^{-1} \mathbf{L}_i [\mathbf{S}_i^* \ \mathbf{O}] \begin{bmatrix} \mathbf{R}_{11}^o & \mathbf{R}_{12}^o \\ \mathbf{R}_{21}^o & \mathbf{R}_{22}^o \end{bmatrix} \begin{bmatrix} \mathbf{S}_j^{\top*} \\ \mathbf{O}^\top \end{bmatrix} \mathbf{L}_j^\top \\ &= [\Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i^* \ \mathbf{O}] \begin{bmatrix} \mathbf{R}_{11i}^o & \mathbf{R}_{12i}^o \\ \mathbf{R}_{21i}^o & \mathbf{R}_{22i}^o \end{bmatrix} \begin{bmatrix} \mathbf{S}_j^{\top*} \mathbf{L}_j^\top \\ \mathbf{O}^\top \end{bmatrix} \\ &= \Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i^* \mathbf{R}_{11i}^o \mathbf{S}_j^{\top*} \mathbf{L}_j^\top. \end{aligned} \quad (7)$$

Then, we have

$$\begin{aligned} \max_{\mathbf{R}_i \in \mathcal{O}(m)} \text{tr}(\Sigma_n^{-1} \mathbf{X}_i \mathbf{R}_i \mathbf{X}_j^\top) &= \max_{\mathbf{R}_{11i}^o \in \mathcal{O}(n)} \text{tr}(\Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i^* \mathbf{R}_{11i}^o \mathbf{S}_j^{\top*} \mathbf{L}_j^\top) \\ &= \max_{\mathbf{R}_i^* \in \mathcal{O}(n)} \text{tr}(\Sigma_n^{-1} \mathbf{X}_i \mathbf{Q}_i \mathbf{R}_i^* \mathbf{Q}_j^\top \mathbf{X}_j^\top). \end{aligned}$$

The last equality is due to

$$\Sigma_n^{-1} \mathbf{X}_i \mathbf{Q}_i \mathbf{R}_i^* \mathbf{Q}_j^\top \mathbf{X}_j^\top = \Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i^* \mathbf{Q}_i^\top \mathbf{Q}_i \mathbf{R}_i^* \mathbf{Q}_j^\top \mathbf{Q}_j \mathbf{S}_j^* \mathbf{L}_j^\top = \Sigma_n^{-1} \mathbf{L}_i \mathbf{S}_i^* \mathbf{R}_i^* \mathbf{S}_j^{\top*} \mathbf{L}_j^\top.$$

So essentially, only the first  $n$  dimensions are used in maximizing (5), if  $n < m$  in all Procrustes-based problem.

*Lemma 5.* Consider the assumptions of Theorem 3, then

$$\max_{\mathbf{R}_i \in \mathcal{O}(m)} \text{tr}(\mathbf{R}_i^\top \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{X}_j \Sigma_m^{-1} + k\mathbf{F}) = \max_{\mathbf{R}_i^* \in \mathcal{O}(n)} \text{tr}\{\mathbf{R}_i^{*\top} (\mathbf{Q}_i^\top \mathbf{X}_i^\top \Sigma_n^{-1} \mathbf{X}_j \Sigma_m^{-1} \mathbf{Q}_j^\top + k\mathbf{F}^*)\},$$

where  $\mathbf{F} \in \mathbb{R}^{m \times m}$  and  $\mathbf{F}^* \in \mathbb{R}^{n \times n}$ .

*Proof.* Without loss of generality we consider  $\Sigma_m = \sigma^2 \mathbf{I}_m$ , and  $\mathbf{F}^* = \mathbf{Q}_i^\top \mathbf{F} \mathbf{Q}_j$ . So, the following objective function to maximize

$$\max_{\mathbf{R}_i^* \in \mathcal{O}(n)} \text{tr}(\mathbf{R}_i^{\top*} \mathbf{Q}_i^{\top} \mathbf{X}_i^{\top} \Sigma_n^{-1} \mathbf{X}_j \mathbf{Q}_j) + k \text{tr}(\mathbf{R}_i^{\top*} \mathbf{Q}_i^{\top} \mathbf{F} \mathbf{Q}_j). \quad (8)$$

The left part of the maximization (8) equals (7), while  $k \text{tr}(\mathbf{R}_i^{\top*} \mathbf{Q}_i^{\top} \mathbf{F} \mathbf{Q}_j)$  in  $\mathcal{O}(n)$  is equivalent to  $k \text{tr}(\mathbf{R}_i^{\top} \mathbf{F})$  in  $\mathcal{O}(m)$ .

### References

- Gower, J. C. & Dijksterhuis, G. B. (2004). *Procrustes problems*, Volume 30. Oxford University Press on Demand.
- Jupp, P. E. & Mardia, K. V. (1979). Maximum likelihood estimators for the matrix von mises-fisher and bingham distributions. *The Annals of Statistics* **7**(3), 599–606.
- Myronenko, A. & Song, X. (2009). On the closed-form solution of the rotation matrix arising in computer vision problems. *arXiv preprint arXiv:0904.1613* .
- Theobald, D. L. & Wuttke, D. S. (2006). Empirical bayes hierarchical models for regularizing maximum likelihood estimation in the matrix gaussian procrustes problem. *Proceedings of the National Academy of Sciences* **103**(49), 18521–18527.
- Trendafilov, N. T. & Lippert, R. A. (2002). The multimode procrustes problem. *Linear algebra and its applications* **349**(1–3), 245–264.
- Viète, F. (1646). *Opera mathematica*. F. van Schouten.