Probabilistic Low-Rank Subspace Clustering Supplementary Material

S. Derin Babacan University of Illinois at Urbana-Champaign Urbana, IL 61801, USA dbabacan@gmail.com Shinichi Nakajima Nikon Corporation Tokyo, 140-8601, Japan nakajima.s@nikon.co.jp

Minh N. Do University of Illinois at Urbana-Champaign Urbana, IL 61801, USA minhdo@illinois.edu

In this supplementary material, we provide the derivation of the global solution of the expectationmaximization method in Sec. 2.2 and the required statistics in the variational Bayesian methods in Secs. 3 and 4. Equation numbers are denoted with preceding "S-", and the ones without "S-" refer to the main text.

1 Global Solution of the EM method

The log-likelihood is given by

$$\mathcal{L} = \sum_{i=1}^{N} \log p(\mathbf{d}_i, \mathbf{y}_i | \mathbf{D}, \mathbf{A})$$

$$= -\frac{N}{2} \left(M \log(\sigma_y^2) + \log |\mathbf{K}| - \frac{1}{N} \operatorname{tr} \left(\mathbf{K}^{-1} \mathbf{D} \mathbf{D}^T \right) \right) - \frac{1}{2\sigma_y^2} \operatorname{tr} \left((\mathbf{Y} - \mathbf{D})^T (\mathbf{Y} - \mathbf{D}) \right) + \operatorname{const} \mathbf{A}$$
(S-1)

with $\mathbf{K} = \sigma_d^2 \mathbf{I} + \mathbf{D} \mathbf{A} \mathbf{A}^T \mathbf{D}^T$. To maximize the log-likelihood w.r.t. \mathbf{A} , we take its gradient w.r.t. \mathbf{A} using matrix differentiation identities [2] and set it equal to zero, which yields

$$\mathbf{D}^{T}\mathbf{K}^{-1}\mathbf{D}\mathbf{A} = \frac{1}{N}\mathbf{D}^{T}\mathbf{K}^{-1}\mathbf{D}\mathbf{D}^{T}\mathbf{K}^{-1}\mathbf{D}\mathbf{A}.$$
 (S-2)

This has three possible solutions: (i) $\mathbf{DA} = \mathbf{0}$, (ii) $\mathbf{K} = \frac{1}{N}\mathbf{DD}^T$, and (iii) $\mathbf{DA} \neq \mathbf{0}$ and $\mathbf{K} \neq \frac{1}{N}\mathbf{DD}^T$. We consider the latter two cases, as the first one is not interesting for subspace clustering. In the last case, assuming $\sigma_d^2 > 0$ and thus \mathbf{K}^{-1} exists, we have

$$\mathbf{DA} = \frac{1}{N} \mathbf{D} \mathbf{D}^T \mathbf{K}^{-1} \mathbf{D} \mathbf{A} \,. \tag{S-3}$$

We first solve this system w.r.t. **DA**. Let the SVDs of **D** and **DA** be¹ **D** = **U** Λ **V**^T and **DA** = $\hat{\mathbf{U}}\hat{\Lambda}\hat{\mathbf{V}}^{T}$, respectively, such that we have

$$\mathbf{K}^{-1}\mathbf{D}\mathbf{A} = \left(\sigma_d^2 \mathbf{I} + \mathbf{D}\mathbf{A}\mathbf{A}^T \mathbf{D}^T\right)^{-1} \mathbf{D}\mathbf{A}, \qquad (S-4)$$

$$= \mathbf{D}\mathbf{A} \left(\sigma_d^2 \mathbf{I} + \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \right)^{-1} , \qquad (S-5)$$

$$= \hat{\mathbf{U}}\hat{\mathbf{\Lambda}} \left(\sigma_d^2 \mathbf{I} + \hat{\mathbf{\Lambda}}^2\right)^{-1} \hat{\mathbf{V}}^T \,. \tag{S-6}$$

¹At this point, we do not know if the singular vectors of **DA** and **D** are related.

Plugging this in (S-2), we have at the stationary points

$$\hat{\mathbf{U}}\hat{\mathbf{\Lambda}}\hat{\mathbf{V}}^{T} = \frac{1}{N}\mathbf{D}\mathbf{D}^{T}\hat{\mathbf{U}}\hat{\mathbf{\Lambda}}\left(\sigma_{d}^{2}\mathbf{I} + \hat{\mathbf{\Lambda}}^{2}\right)^{-1}\hat{\mathbf{V}}^{T},$$
(S-7)

$$\hat{\mathbf{U}}\left(\sigma_d^2 \mathbf{I} + \hat{\mathbf{\Lambda}}^2\right) \hat{\mathbf{\Lambda}} = \frac{1}{N} \mathbf{D} \mathbf{D}^T \hat{\mathbf{U}} \hat{\mathbf{\Lambda}}, \qquad (S-8)$$

from which it can be observed that $\hat{\mathbf{U}}$ contains the eigenvectors of $\mathbf{D}\mathbf{D}^T$ and hence the left singular vectors of \mathbf{D} , such that $\hat{\mathbf{U}} = \mathbf{U}$. Moreover, $\sigma_d^2 \mathbf{I} + \hat{\mathbf{\Lambda}}^2$ contains the eigenvalues of $\frac{1}{N}\mathbf{D}\mathbf{D}^T$. Therefore, similarly to [3], we have the solution

$$\mathbf{DA} = \mathbf{U}_q \left(\frac{1}{N}\mathbf{\Lambda}_q^2 - \sigma_d^2 \mathbf{I}\right)^{1/2} \mathbf{R}, \qquad (S-9)$$

where **R** is an arbitrary orthogonal rotation matrix, and \mathbf{U}_q is a $M \times q$ matrix consisting of q left singular vectors of **D** with corresponding singular values that are larger than $\sqrt{N}\sigma_d$. Therefore, the singular values of **DA** satisfy $l_i = (\frac{\lambda_i^2}{N} - \sigma_d^2)^{1/2}$.

In the case (ii), we have the same solution (S-9) where the last M - q smallest singular values of **D** are equal to $\sqrt{N\sigma_d}$. This is an unrealistic case and is analyzed also in PPCA [3].

Using the solution (S-9), we can solve for the optimal B using (9) as

$$\langle \mathbf{B} \rangle = \boldsymbol{\Sigma}_{\mathbf{B}} \frac{1}{\sigma_d^2} \mathbf{A}^T \mathbf{D}^T \mathbf{D}, \qquad (S-10)$$

$$= \left(\sigma_d^2 \mathbf{I} + \mathbf{R}^T (\frac{1}{N} \mathbf{\Lambda}_q^2 - \sigma_d^2 \mathbf{I}) \mathbf{R}\right)^{-1} \mathbf{R}^T (\frac{1}{N} \mathbf{\Lambda}_q^2 - \sigma_d^2 \mathbf{I})^{1/2} \mathbf{U}_q^T \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T,$$
(S-11)

$$= \mathbf{R}^{T} \sigma_{d}^{-2} \left(\frac{1}{N} \mathbf{\Lambda}_{q}^{2} - \sigma_{d}^{2} \mathbf{I}\right)^{1/2} \left(\mathbf{I} + \sigma_{d}^{-2} \left(\frac{1}{N} \mathbf{\Lambda}_{q}^{2} - \sigma_{d}^{2} \mathbf{I}\right) \mathbf{R} \mathbf{R}^{T}\right)^{-1} \mathbf{\Lambda}_{q} \mathbf{V}_{q}^{T}, \qquad (S-12)$$

$$= \mathbf{R}^{T} (\frac{1}{N} \mathbf{\Lambda}_{q}^{2} - \sigma_{d}^{2} \mathbf{I})^{1/2} \mathbf{\Lambda}_{q}^{-1} N \mathbf{V}_{q}^{T}.$$
(S-13)

Now we have an expression for **DA** and $\langle \mathbf{B} \rangle$. Combining,

$$\mathbf{DA}\langle \mathbf{B} \rangle = \mathbf{U}_q (\mathbf{\Lambda}_q^2 - N\sigma_d^2 \mathbf{I}) \mathbf{\Lambda}_q^{-1} \mathbf{V}_q^T \,. \tag{S-14}$$

Plugging $\mathbf{D} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T$ in (S-14) yields the final solution

$$\mathbf{A} \langle \mathbf{B} \rangle = \mathbf{V}_q (\mathbf{\Lambda}_q^2 - N\sigma_d^2 \mathbf{I}) \mathbf{\Lambda}_q^{-2} \mathbf{V}_q^T = \mathbf{V}_q \tilde{\mathbf{\Lambda}}_q \mathbf{V}_q^T, \qquad (S-15)$$

with $\tilde{\Lambda}_q$ is a diagonal matrix with $1 - \frac{N\sigma_d^2}{\lambda_j^2}$ on the diagonal. The optimal solution for **A** can easily be extracted from this expression.

Finally, using this expression for $\mathbf{A} \langle \mathbf{B} \rangle$ in (10), we solve for \mathbf{D} as

$$\mathbf{Y} = \mathbf{D} \left[\mathbf{I} + \frac{\sigma_y^2}{\sigma_d^2} \langle \left(\mathbf{I} - \mathbf{A} \mathbf{B} \right) \left(\mathbf{I} - \mathbf{A} \mathbf{B} \right)^T \rangle \right], \qquad (S-16)$$

$$= \mathbf{U} \mathbf{\Lambda} \mathbf{V}^{T} \left[\mathbf{I} + N \sigma_{y}^{2} \mathbf{V}_{q} \mathbf{\Lambda}_{q}^{-2} \mathbf{V}_{q}^{T} \right] , \qquad (S-17)$$

Using the partitioning $\mathbf{D} = [\mathbf{U}_q, \mathbf{U}_{N-q}] \operatorname{diag}(\mathbf{\Lambda}_q, \mathbf{\Lambda}_{N-q}) [\mathbf{V}_q, \mathbf{V}_{N-q}]^T$, we have the final solution

$$\mathbf{Y} = [\mathbf{U}_q, \mathbf{U}_{N-q}] \begin{bmatrix} \mathbf{\Lambda}_q + N \sigma_y^2 \mathbf{\Lambda}_q^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{\sigma_y^2 + \sigma_d^2}{\sigma_d^2} \mathbf{\Lambda}_{N-q} \end{bmatrix} [\mathbf{V}_q \mathbf{V}_{N-q}]^T.$$
(S-18)

Therefore, the eigenvectors of D and Y are the same, but the eigenvalues are related via

$$\xi_{j} = \begin{cases} \lambda_{j} + N\sigma_{y}^{2} \lambda_{j}^{-1}, & \text{if } \lambda_{j} > \sqrt{N}\sigma_{d} \\ \lambda_{j} \frac{\sigma_{y}^{2} + \sigma_{d}^{2}}{\sigma_{d}^{2}}, & \text{if } \lambda_{j} \le \sqrt{N}\sigma_{d} \end{cases}$$
(S-19)

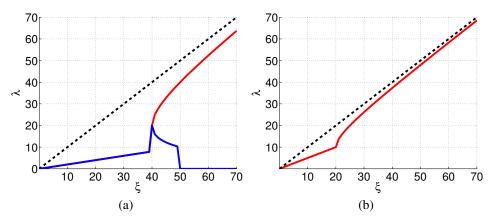


Figure 1: Estimates of singular values λ of **D** given singular values ξ of **Y** (N = 100). The dashed line is $\lambda = \xi$. In (a), $\sigma_d = 1$, $\sigma_y = 2$, in (b), $\sigma_d = \sigma_y = 1$.

The explicit solutions for λ_j are given by

$$\lambda_{j} = \begin{cases} \xi_{j} \frac{\sigma_{d}^{2}}{\sigma_{y}^{2} + \sigma_{d}^{2}}, & \xi_{j} < 2\sqrt{N}\sigma_{y} \\ \frac{\xi_{j}}{2} + \frac{1}{2}\sqrt{\xi_{j}^{2} - 4N\sigma_{y}^{2}} & \xi_{j} \ge 2\sqrt{N}\sigma_{y}, \xi_{j} \ge \min\left(2\sqrt{N}\sigma_{d}, \frac{\sqrt{N}}{\sigma_{d}}(\sigma_{d}^{2} + \sigma_{y}^{2})\right) \\ \frac{\xi_{j}}{2} - \frac{1}{2}\sqrt{\xi_{j}^{2} - 4N\sigma_{y}^{2}} & \sigma_{y} \ge \sigma_{d}, 2\sqrt{N}\sigma_{y} \le \xi_{j} \le \frac{\sqrt{N}}{\sigma_{d}}(\sigma_{d}^{2} + \sigma_{y}^{2}) \end{cases}$$
(S-20)

The solution for λ_j is unique except when $\sigma_y \ge \sigma_d$ and $2\sqrt{N}\sigma_y \le \xi_j \le \frac{\sqrt{N}}{\sigma_d}(\sigma_d^2 + \sigma_y^2)$, where we have the latter two cases as solutions. As shown in Fig. 1(a), the last solution is only valid in a comparably small region. To achieve continuity in the solutions, we always choose the first two solutions (S-20).

As can be observed from Fig. 1, the solution (S-20) is a combination of two operations: a downscaling when $\xi_j < 2\sqrt{N}\sigma_y$ and a polynomial thresholding operation for larger singular values. The polynomial thresholding preserves the larger singular values as the shrinkage amount gets smaller: ξ_j gets larger compared to $2N\sigma_y$, and for very large values $\lambda_j \approx \xi_j$. On the other hand, small singular values get shrunk via down-scaling. Obviously, when $\sigma_d = 0$, no shrinkage is applied and $\mathbf{D} = \mathbf{Y}$.

2 Derivation of the Variational Bayesian Methods

The explicit form of the variational free energy in (17) is given by

$$\begin{split} \mathcal{F} &= \langle \log \mathbf{q}(\mathbf{D}, \mathbf{A}, \mathbf{B}, \sigma_d^2, \sigma_y^2) - \log \mathbf{p}(\mathbf{D}, \mathbf{A}, \mathbf{B}, \sigma_d^2, \sigma_y^2) \rangle_{q(\mathbf{D}, \mathbf{A}, \mathbf{B}, \sigma_d^2, \sigma_y^2)} \\ &= \langle \log \mathbf{q}(\mathbf{D}) \mathbf{q}(\mathbf{A}) \mathbf{q}(\mathbf{B}) \mathbf{q}(\sigma_d^2) \mathbf{q}(\sigma_y^2) \rangle \\ &+ \frac{MN}{2} \langle \log \sigma_d^2 \rangle + \frac{MN}{2} \langle \log \sigma_y^2 \rangle + \frac{1}{2} \operatorname{tr}(\langle \mathbf{A} \mathbf{C}_{\mathbf{A}}^{-1} \mathbf{A}^T \rangle) + \frac{1}{2} \operatorname{tr}(\langle \mathbf{C}_{\mathbf{B}}^{-1} \mathbf{B} \mathbf{B}^T \rangle) + \frac{1}{2} \operatorname{tr}(\langle \mathbf{D} \mathbf{D}^T \rangle) \\ &+ \left(\frac{1}{2 \langle \sigma_y^2 \rangle} + \frac{1}{2 \langle \sigma_d^2 \rangle} \right) \operatorname{tr}(\langle \mathbf{D} \mathbf{D}^T \rangle) + \frac{1}{2 \langle \sigma_y^2 \rangle} \| \mathbf{Y} \|_{\mathrm{F}}^2 - \frac{1}{\langle \sigma_y^2 \rangle} \operatorname{tr}(\langle \mathbf{D} \rangle^T \mathbf{Y}) \\ &- \frac{1}{\langle \sigma_d^2 \rangle} \operatorname{tr}(\langle \mathbf{B}^T \mathbf{A}^T \mathbf{D}^T \mathbf{D} \rangle) + \frac{1}{2 \langle \sigma_d^2 \rangle} \operatorname{tr}(\langle \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \mathbf{B} \mathbf{B}^T \rangle) \tag{S-21} \\ &+ \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \frac{N}{2} \log |\mathbf{C}_{\mathbf{B}}| + \operatorname{const}. \end{split}$$

The optimal forms of $q(\mathbf{D})$ and $q(\mathbf{B})$ can be found as matrix-variate normal distributions by inspection. The optimal $q(\mathbf{A})$ does not have a matrix-variate normal form. The optimal distribution is

found in terms of $\mathbf{a} = \text{vec}(\mathbf{A})$, by rewriting the terms involving \mathbf{A} in (S-23) as

$$\begin{aligned} -\log \mathbf{q}(\mathbf{a}) &= \operatorname{tr} \left(\langle \sigma_d^{-2} \rangle \langle \| \mathbf{D} - \mathbf{D} \mathbf{A} \mathbf{B} \|_F \rangle^2 + \mathbf{A} \mathbf{C}_{\mathbf{A}}^{-1} \mathbf{A}^T \right) + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &= \langle \sigma_d^{-2} \rangle \langle \| \mathbf{d} - (\mathbf{B}^T \otimes \mathbf{D}) \mathbf{a} \|_2^2 \rangle + \mathbf{a}^T (\mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I}) \mathbf{a} + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &= \langle \sigma_d^{-2} \rangle \langle (\mathbf{d}^T \mathbf{d} + \mathbf{a}^T (\mathbf{B}^T \otimes \mathbf{D})^T (\mathbf{B}^T \otimes \mathbf{D}) \mathbf{a} - 2\mathbf{a}^T (\mathbf{B}^T \otimes \mathbf{D})^T \mathbf{d} \rangle \rangle + \mathbf{a}^T (\mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I}) \mathbf{a} + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &= \mathbf{a}^T \left[\langle \sigma_d^{-2} \rangle \langle (\mathbf{B}^T \otimes \mathbf{D})^T (\mathbf{B}^T \otimes \mathbf{D}) \rangle + \mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I} \right] \mathbf{a} - 2\mathbf{a}^T (\langle \mathbf{B} \rangle^T \otimes \langle \mathbf{D} \rangle)^T \langle \mathbf{d} \rangle + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &= \mathbf{a}^T \left[\langle \sigma_d^{-2} \rangle (\langle \mathbf{B}^T \mathbf{B} \rangle \otimes \langle \mathbf{D}^T \mathbf{D} \rangle) + \mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I} \right] \mathbf{a} - 2\mathbf{a}^T (\langle \mathbf{B} \rangle^T \otimes \langle \mathbf{D} \rangle)^T \langle \mathbf{d} \rangle + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &= \mathbf{a}^T \left[\langle \sigma_d^{-2} \rangle (\langle \mathbf{B}^T \mathbf{B} \rangle \otimes \langle \mathbf{D}^T \mathbf{D} \rangle) + \mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I} \right] \mathbf{a} - 2\mathbf{a}^T (\langle \mathbf{B} \rangle^T \otimes \langle \mathbf{D} \rangle)^T \langle \mathbf{d} \rangle + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \operatorname{const} \\ &\quad (\mathbf{S}\text{-}\mathbf{2}) \end{aligned}$$

where we used $\operatorname{vec}(\mathbf{DAB}) = (\mathbf{B}^T \otimes \mathbf{D}) \operatorname{vec}(\mathbf{A})$, and $\mathbf{d} = \operatorname{vec}(\mathbf{D})$, $\mathbf{b} = \operatorname{vec}(\mathbf{B})$. It can be derived from here that $q(\mathbf{a})$ has a multivariate normal distribution with mean $\Sigma_{\mathbf{a}} (\langle \mathbf{B} \rangle^T \otimes \langle \mathbf{D} \rangle)^T \langle \mathbf{d} \rangle$ and covariance $\Sigma_{\mathbf{a}} = [\langle \sigma_d^{-2} \rangle (\langle \mathbf{B}^T \mathbf{B} \rangle \otimes \langle \mathbf{D}^T \mathbf{D} \rangle) + \mathbf{C}_{\mathbf{A}}^{-1} \otimes \mathbf{I}]^{-1}$. However, computing \mathbf{A} in this manner can be very inefficient, as $\Sigma_{\mathbf{A}}$ might get extremely big $(MN \times MN \text{ for } \mathbf{A} \text{ of size } N \times N)$ and \mathbf{D} of size $M \times N$).

Therefore, we force $q(\mathbf{A})$ to have a matrix-variate form $\mathcal{N}(\langle \mathbf{A} \rangle, \Sigma_{\mathbf{A}}, \Omega_{\mathbf{A}})$, which leads to an efficient algorithm. Under this constraint, the variational free energy can be rewritten as (treating all terms not involving \mathbf{A} as constant)

$$\mathcal{F} = \frac{1}{2} \operatorname{tr}(\langle \mathbf{A} \mathbf{C}_{\mathbf{A}}^{-1} \mathbf{A}^{T} \rangle) - \frac{1}{\langle \sigma_{d}^{2} \rangle} \operatorname{tr}(\langle \mathbf{B}^{T} \mathbf{A}^{T} \mathbf{D}^{T} \mathbf{D} \rangle) + \frac{1}{2 \langle \sigma_{d}^{2} \rangle} \operatorname{tr}(\langle \mathbf{A}^{T} \mathbf{D}^{T} \mathbf{D} \mathbf{A} \mathbf{B} \mathbf{B}^{T} \rangle) \quad (S-23)$$
$$- \frac{N}{2} \log |\mathbf{\Sigma}_{\mathbf{A}}| - \frac{N}{2} \log |\mathbf{\Omega}_{\mathbf{A}}| + \frac{N}{2} \log |\mathbf{C}_{\mathbf{A}}| + \frac{N}{2} \log |\mathbf{C}_{\mathbf{B}}| + \operatorname{const.}$$

Evaluating the expectations using the matrix-variate normal form for $q(\mathbf{A})$ (see the next section), we minimize \mathcal{F} with respect to $\Sigma_{\mathbf{A}}$, resulting in

$$\boldsymbol{\Sigma}_{\mathbf{A}}^{-1} = \frac{1}{N} \operatorname{tr}(\mathbf{C}_{\mathbf{A}}^{-1} \boldsymbol{\Omega}_{\mathbf{A}}) \mathbf{I} + \frac{1}{N \sigma_d^2} \operatorname{tr}(\boldsymbol{\Omega}_{\mathbf{A}} \langle \mathbf{B} \mathbf{B}^T \rangle) \langle \mathbf{D}^T \mathbf{D} \rangle$$
(S-24)

Similarly, minimization with respect to Ω_A yields

$$\mathbf{\Omega}_{\mathbf{A}}^{-1} = \frac{1}{N} \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{A}}) \mathbf{C}_{\mathbf{A}}^{-1} + \frac{1}{N \sigma_d^2} \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{A}} \langle \mathbf{D}^T \mathbf{D} \rangle) \langle \mathbf{B} \mathbf{B}^T \rangle .$$
(S-25)

Finally, the update of $\langle \mathbf{A} \rangle$ is given by

$$\langle \mathbf{A} \rangle \mathbf{C}_{\mathbf{A}}^{-1} + \frac{1}{\sigma_d^2} \langle \mathbf{D}^T \mathbf{D} \rangle \langle \mathbf{A} \rangle \langle \mathbf{B} \mathbf{B}^T \rangle = \frac{1}{\sigma_d^2} \langle \mathbf{D}^T \mathbf{D} \rangle \langle \mathbf{B} \rangle^T$$
(S-26)

The closed form solution for $\langle \mathbf{A} \rangle$ cannot be found, but it can be solved using a fixed-point iteration starting from an initial estimate.

2.1 Required Statistics for the Variational Bayesian Methods

For a general matrix-variate Gaussian distribution $p(\mathbf{X}|\mathbf{M}, \Omega, \Sigma) = \mathcal{N}(\mathbf{X}|\mathbf{M}, \Sigma, \Omega)$, we have [1]

$$\langle \mathbf{X}^T \mathbf{K} \mathbf{X} \rangle = \operatorname{tr}(\mathbf{\Sigma} \mathbf{K}^T) \mathbf{\Omega} + \mathbf{M}^T \mathbf{K} \mathbf{M},$$
 (S-27)

$$\langle \mathbf{X}\mathbf{K}\mathbf{X}^T \rangle = \operatorname{tr}(\mathbf{K}^T \mathbf{\Omega}) \mathbf{\Sigma} + \mathbf{M}\mathbf{K}\mathbf{M}^T.$$
 (S-28)

Thus, for $q(\mathbf{D}) = \mathcal{N}(\langle \mathbf{D} \rangle, \mathbf{I}, \Omega_{\mathbf{D}}), q(\mathbf{A}) = \mathcal{N}(\langle \mathbf{A} \rangle, \Sigma_{\mathbf{A}}, \Omega_{\mathbf{A}})$, and $q(\mathbf{B}) = \mathcal{N}(\langle \mathbf{B} \rangle, \mathbf{I}, \Sigma_{\mathbf{B}})$, we have

$$\langle \mathbf{D}^T \mathbf{D} \rangle = \operatorname{tr}(\mathbf{I}_M) \mathbf{\Omega}_{\mathbf{D}} + \langle \mathbf{D} \rangle^T \langle \mathbf{D} \rangle$$
 (S-29)

$$= M \mathbf{\Omega}_{\mathbf{D}} + \langle \mathbf{D} \rangle^T \langle \mathbf{D} \rangle \tag{S-30}$$

$$\langle \mathbf{A}\mathbf{A}^T \rangle = \operatorname{tr}(\mathbf{\Omega}_{\mathbf{A}})\mathbf{\Sigma}_{\mathbf{A}} + \langle \mathbf{A} \rangle \langle \mathbf{A} \rangle^T$$
 (S-31)

$$\langle \mathbf{A}^T \mathbf{A} \rangle = \operatorname{tr}(\boldsymbol{\Sigma}_{\mathbf{A}}) \boldsymbol{\Omega}_{\mathbf{A}} + \langle \mathbf{A} \rangle^T \langle \mathbf{A} \rangle$$
 (S-32)

$$\langle \mathbf{B}\mathbf{B}^T \rangle = \operatorname{tr}(\mathbf{\Omega}_{\mathbf{B}})\mathbf{\Sigma}_{\mathbf{B}} + \langle \mathbf{B} \rangle \langle \mathbf{B} \rangle^T$$
 (S-33)

$$= N \Sigma_{\mathbf{B}} + \langle \mathbf{B} \rangle \langle \mathbf{B} \rangle^{T}$$
(S-34)

$$\langle \mathbf{B}^T \mathbf{B} \rangle = \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{B}}) \mathbf{I}_N + \langle \mathbf{B} \rangle^T \langle \mathbf{B} \rangle$$
 (S-35)

Combining, we obtain

$$\langle \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \rangle = \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{A}} \langle \mathbf{D}^T \mathbf{D} \rangle) \mathbf{\Omega}_{\mathbf{A}} + \langle \mathbf{A} \rangle^T \langle \mathbf{D}^T \mathbf{D} \rangle \langle \mathbf{A} \rangle$$
(S-36)

$$\langle \mathbf{B}^{T} \mathbf{A}^{T} \mathbf{A} \mathbf{B} \rangle = \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{B}} \langle \mathbf{A}^{T} \mathbf{A} \rangle) \mathbf{I}_{N} + \langle \mathbf{B} \rangle^{T} \langle \mathbf{A}^{T} \mathbf{A} \rangle \langle \mathbf{B} \rangle$$
(S-37)

$$\langle \mathbf{A}\mathbf{B}\mathbf{B}^{T}\mathbf{A}^{T}\rangle = \operatorname{tr}(\langle \mathbf{B}\mathbf{B}^{T}\rangle\boldsymbol{\Omega}_{\mathbf{A}})\boldsymbol{\Sigma}_{\mathbf{A}} + \langle \mathbf{A}\rangle\langle \mathbf{B}\mathbf{B}^{T}\rangle\langle \mathbf{A}\rangle^{T}$$
(S-38)

$$\langle \mathbf{B}^T \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \mathbf{B} \rangle = \operatorname{tr}(\mathbf{\Sigma}_{\mathbf{B}} \langle \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \rangle) \mathbf{I}_N + \mathbf{B}^T \langle \mathbf{A}^T \mathbf{D}^T \mathbf{D} \mathbf{A} \rangle \mathbf{B}$$
(S-39)

References

- [1] A. K. Gupta and D. K. Nagar. Matrix Variate Distributions. Chapman & Hall/CRC, New York, 2000.
- [2] K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate analysis*. Academic Press; New York, 1979.
- [3] M. E. Tipping and C. M. Bishop. Mixtures of probabilistic principal component analyzers. *Neural Comput.*, 11(2):443–482, February 1999.