# 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 expectation-maximization 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

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

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

with $\mathbf{K} = \sigma_d^2 \mathbf{I} + \mathbf{DAA}^T\mathbf{D}^T$. 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

$$\mathbf{D}^T\mathbf{K}^{-1}\mathbf{DA} = \frac{1}{N}\mathbf{D}^T\mathbf{K}^{-1}\mathbf{DD}^T\mathbf{K}^{-1}\mathbf{DA}. \tag{S-2}$$

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

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

We first solve this system w.r.t. $\mathbf{DA}$. Let the SVDs of $\mathbf{D}$ and $\mathbf{DA}$ be[1] $\mathbf{D} = \mathbf{U\Lambda V}^T$ and $\mathbf{DA} = \hat{\mathbf{U}}\hat{\mathbf{\Lambda}}\hat{\mathbf{V}}^T$, respectively, such that we have

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

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

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

---

[1] At this point, we do not know if the singular vectors of $\mathbf{DA}$ and $\mathbf{D}$ are related.

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

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

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

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

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

where $\mathbf{R}$ is an arbitrary orthogonal rotation matrix, and $\mathbf{U}_q$ is a $M \times q$ matrix consisting of $q$ left singular vectors of $\mathbf{D}$ with corresponding singular values that are larger than $\sqrt{N}\sigma_d$. Therefore, the singular values of $\mathbf{D}\mathbf{A}$ 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 $\mathbf{D}$ 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 $\mathbf{B}$ using (9) as

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

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

$$= \mathbf{R}^T\sigma_d^{-2}(\frac{1}{N}\mathbf{\Lambda}_q^2 - \sigma_d^2\mathbf{I})^{1/2}\left(\mathbf{I}+\sigma_d^{-2}(\frac{1}{N}\mathbf{\Lambda}_q^2 - \sigma_d^2\mathbf{I})\mathbf{R}\mathbf{R}^T\right)^{-1}\mathbf{\Lambda}_q\mathbf{V}_q^T,\tag{S-12}$$

$$= \mathbf{R}^T(\frac{1}{N}\mathbf{\Lambda}_q^2 - \sigma_d^2\mathbf{I})^{1/2}\mathbf{\Lambda}_q^{-1}N\mathbf{V}_q^T.\tag{S-13}$$

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

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

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

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

with $\tilde{\mathbf{\Lambda}}_q$ is a diagonal matrix with $1 - \frac{N\sigma_d^2}{\lambda_j^2}$ on the diagonal. The optimal solution for $\mathbf{A}$ 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

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

$$= \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T\left[\mathbf{I}+N\sigma_y^2\mathbf{V}_q\mathbf{\Lambda}_q^{-2}\mathbf{V}_q^T\right],\tag{S-17}$$

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

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

Therefore, the eigenvectors of $\mathbf{D}$ and $\mathbf{Y}$ are the same, but the eigenvalues are related via

$$\xi_j = \begin{cases}\lambda_j + N\sigma_y^2\,\lambda_j^{-1}, & \text{if } \lambda_j > \sqrt{N}\sigma_d \\ \lambda_j\frac{\sigma_y^2+\sigma_d^2}{\sigma_d^2}, & \text{if } \lambda_j \leq \sqrt{N}\sigma_d\end{cases}\tag{S-19}$$
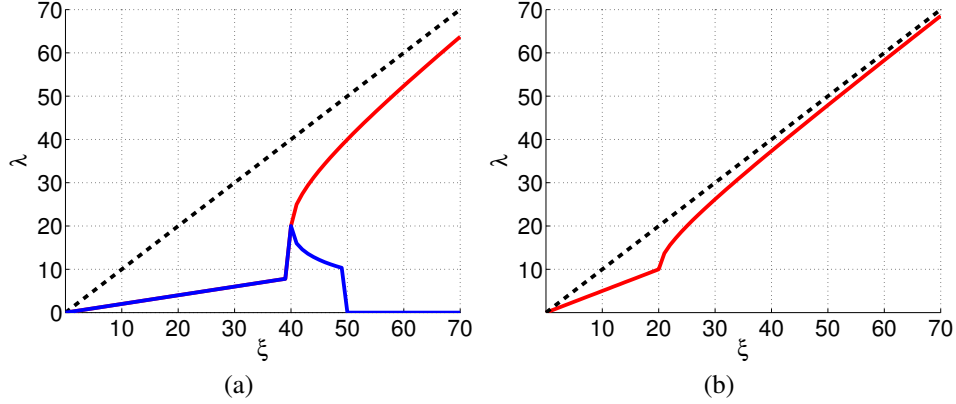
Figure 1: Estimates of singular values $\lambda$ of $\mathbf{D}$ given singular values $\xi$ of $\mathbf{Y}$ (N = 100). The dashed line is $\lambda = \xi$. In (a), $\sigma_d = 1$, $\sigma_y = 2$, in (b), $\sigma_d = \sigma_y = 1$.

The explicit solutions for $\lambda_j$ are given by

$$
\lambda_j = \begin{cases} \xi_j \frac{\sigma_d^2}{\sigma_y^2 + \sigma_d^2}, & \xi_j < 2\sqrt{N}\sigma_y \\ \frac{\xi_j}{2} + \frac{1}{2}\sqrt{\xi_j^2 - 4N\sigma_y^2} & \xi_j \geq 2\sqrt{N}\sigma_y, \xi_j \geq \min\left(2\sqrt{N}\sigma_d, \frac{\sqrt{N}}{\sigma_d}(\sigma_d^2 + \sigma_y^2)\right) \\ \frac{\xi_j}{2} - \frac{1}{2}\sqrt{\xi_j^2 - 4N\sigma_y^2} & \sigma_y \geq \sigma_d, 2\sqrt{N}\sigma_y \leq \xi_j \leq \frac{\sqrt{N}}{\sigma_d}(\sigma_d^2 + \sigma_y^2) \end{cases} \tag{S-20}
$$

The solution for $\lambda_j$ is unique except when $\sigma_y \geq \sigma_d$ and $2\sqrt{N}\sigma_y \leq \xi_j \leq \frac{\sqrt{N}}{\sigma_d}(\sigma_d^2 + \sigma_y^2)$, 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 down-scaling 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: $\xi_j$ 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

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

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

3

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

$$
\begin{aligned}
-\log \mathrm{q}(\mathbf{a}) &= \text{tr}\left(\langle\sigma_d^{-2}\rangle\,\langle\|\mathbf{D} - \mathbf{DAB}\|_F\rangle^2 + \mathbf{AC_A^{-1}A}^T\right) + \frac{N}{2}\log|\mathbf{C_A}| + \text{const}\\
&= \langle\sigma_d^{-2}\rangle\,\langle\|\mathbf{d} - (\mathbf{B}^T\otimes\mathbf{D})\mathbf{a}\|_2^2\rangle + \mathbf{a}^T(\mathbf{C_A^{-1}}\otimes\mathbf{I})\mathbf{a} + \frac{N}{2}\log|\mathbf{C_A}| + \text{const}\\
&= \langle\sigma_d^{-2}\rangle\,\langle\left(\mathbf{d}^T\mathbf{d} + \mathbf{a}^T(\mathbf{B}^T\otimes\mathbf{D})^T(\mathbf{B}^T\otimes\mathbf{D})\mathbf{a} - 2\mathbf{a}^T(\mathbf{B}^T\otimes\mathbf{D})^T\mathbf{d}\right)\rangle + \mathbf{a}^T(\mathbf{C_A^{-1}}\otimes\mathbf{I})\mathbf{a} + \frac{N}{2}\log|\mathbf{C_A}| + \text{const}\\
&= \mathbf{a}^T\left[\langle\sigma_d^{-2}\rangle\langle(\mathbf{B}^T\otimes\mathbf{D})^T(\mathbf{B}^T\otimes\mathbf{D})\rangle + \mathbf{C_A^{-1}}\otimes\mathbf{I}\right]\mathbf{a} - 2\mathbf{a}^T(\langle\mathbf{B}\rangle^T\otimes\langle\mathbf{D}\rangle)^T\langle\mathbf{d}\rangle + \frac{N}{2}\log|\mathbf{C_A}| + \text{const}\\
&= \mathbf{a}^T\left[\langle\sigma_d^{-2}\rangle(\langle\mathbf{B}^T\mathbf{B}\rangle\otimes\langle\mathbf{D}^T\mathbf{D}\rangle) + \mathbf{C_A^{-1}}\otimes\mathbf{I}\right]\mathbf{a} - 2\mathbf{a}^T(\langle\mathbf{B}\rangle^T\otimes\langle\mathbf{D}\rangle)^T\langle\mathbf{d}\rangle + \frac{N}{2}\log|\mathbf{C_A}| + \text{const}
\end{aligned}
\tag{S-22}
$$

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

Therefore, we force $\mathrm{q}(\mathbf{A})$ to have a matrix-variate form $\mathcal{N}(\langle\mathbf{A}\rangle, \boldsymbol{\Sigma_A}, \boldsymbol{\Omega_A})$, 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)

$$
\begin{aligned}
\mathcal{F} = &\frac{1}{2}\text{tr}(\langle\mathbf{AC_A^{-1}A}^T\rangle) - \frac{1}{\langle\sigma_d^2\rangle}\text{tr}(\langle\mathbf{B}^T\mathbf{A}^T\mathbf{D}^T\mathbf{D}\rangle) + \frac{1}{2\langle\sigma_d^2\rangle}\text{tr}(\langle\mathbf{A}^T\mathbf{D}^T\mathbf{DABB}^T\rangle)\\
&- \frac{N}{2}\log|\boldsymbol{\Sigma_A}| - \frac{N}{2}\log|\boldsymbol{\Omega_A}| + \frac{N}{2}\log|\mathbf{C_A}| + \frac{N}{2}\log|\mathbf{C_B}| + \text{const}\,.
\end{aligned}
\tag{S-23}
$$

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

$$
\boldsymbol{\Sigma_A^{-1}} = \frac{1}{N}\text{tr}(\mathbf{C_A^{-1}}\boldsymbol{\Omega_A})\,\mathbf{I} + \frac{1}{N\sigma_d^2}\text{tr}(\boldsymbol{\Omega_A}\langle\mathbf{BB}^T\rangle)\,\langle\mathbf{D}^T\mathbf{D}\rangle
\tag{S-24}
$$

Similarly, minimization with respect to $\boldsymbol{\Omega_A}$ yields

$$
\boldsymbol{\Omega_A^{-1}} = \frac{1}{N}\text{tr}(\boldsymbol{\Sigma_A})\mathbf{C_A^{-1}} + \frac{1}{N\sigma_d^2}\text{tr}(\boldsymbol{\Sigma_A}\langle\mathbf{D}^T\mathbf{D}\rangle)\,\langle\mathbf{BB}^T\rangle\,.
\tag{S-25}
$$

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

$$
\langle\mathbf{A}\rangle\mathbf{C_A^{-1}} + \frac{1}{\sigma_d^2}\langle\mathbf{D}^T\mathbf{D}\rangle\langle\mathbf{A}\rangle\langle\mathbf{BB}^T\rangle = \frac{1}{\sigma_d^2}\langle\mathbf{D}^T\mathbf{D}\rangle\langle\mathbf{B}\rangle^T
\tag{S-26}
$$

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 $\mathrm{p}(\mathbf{X}|\mathbf{M}, \boldsymbol{\Omega}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{X}|\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Omega})$, we have [1]

$$
\langle\mathbf{X}^T\mathbf{KX}\rangle = \text{tr}(\boldsymbol{\Sigma}\mathbf{K}^T)\boldsymbol{\Omega} + \mathbf{M}^T\mathbf{KM}\,,
\tag{S-27}
$$

$$
\langle\mathbf{XKX}^T\rangle = \text{tr}(\mathbf{K}^T\boldsymbol{\Omega})\boldsymbol{\Sigma} + \mathbf{MKM}^T\,.
\tag{S-28}
$$

Thus, for $q(\mathbf{D}) = \mathcal{N}(\langle\mathbf{D}\rangle, \mathbf{I}, \mathbf{\Omega_D})$, $q(\mathbf{A}) = \mathcal{N}(\langle\mathbf{A}\rangle, \mathbf{\Sigma_A}, \mathbf{\Omega_A})$, and $q(\mathbf{B}) = \mathcal{N}(\langle\mathbf{B}\rangle, \mathbf{I}, \mathbf{\Sigma_B})$, we have

$$\langle\mathbf{D}^T\mathbf{D}\rangle = \text{tr}(\mathbf{I}_M)\mathbf{\Omega_D} + \langle\mathbf{D}\rangle^T\langle\mathbf{D}\rangle \tag{S-29}$$

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

$$\langle\mathbf{A}\mathbf{A}^T\rangle = \text{tr}(\mathbf{\Omega_A})\mathbf{\Sigma_A} + \langle\mathbf{A}\rangle\langle\mathbf{A}\rangle^T \tag{S-31}$$

$$\langle\mathbf{A}^T\mathbf{A}\rangle = \text{tr}(\mathbf{\Sigma_A})\mathbf{\Omega_A} + \langle\mathbf{A}\rangle^T\langle\mathbf{A}\rangle \tag{S-32}$$

$$\langle\mathbf{B}\mathbf{B}^T\rangle = \text{tr}(\mathbf{\Omega_B})\mathbf{\Sigma_B} + \langle\mathbf{B}\rangle\langle\mathbf{B}\rangle^T \tag{S-33}$$

$$= N\mathbf{\Sigma_B} + \langle\mathbf{B}\rangle\langle\mathbf{B}\rangle^T \tag{S-34}$$

$$\langle\mathbf{B}^T\mathbf{B}\rangle = \text{tr}(\mathbf{\Sigma_B})\mathbf{I}_N + \langle\mathbf{B}\rangle^T\langle\mathbf{B}\rangle \tag{S-35}$$

Combining, we obtain

$$\langle\mathbf{A}^T\mathbf{D}^T\mathbf{D}\mathbf{A}\rangle = \text{tr}(\mathbf{\Sigma_A}\langle\mathbf{D}^T\mathbf{D}\rangle)\mathbf{\Omega_A} + \langle\mathbf{A}\rangle^T\langle\mathbf{D}^T\mathbf{D}\rangle\langle\mathbf{A}\rangle \tag{S-36}$$

$$\langle\mathbf{B}^T\mathbf{A}^T\mathbf{A}\mathbf{B}\rangle = \text{tr}(\mathbf{\Sigma_B}\langle\mathbf{A}^T\mathbf{A}\rangle)\mathbf{I}_N + \langle\mathbf{B}\rangle^T\langle\mathbf{A}^T\mathbf{A}\rangle\langle\mathbf{B}\rangle \tag{S-37}$$

$$\langle\mathbf{A}\mathbf{B}\mathbf{B}^T\mathbf{A}^T\rangle = \text{tr}(\langle\mathbf{B}\mathbf{B}^T\rangle\mathbf{\Omega_A})\mathbf{\Sigma_A} + \langle\mathbf{A}\rangle\langle\mathbf{B}\mathbf{B}^T\rangle\langle\mathbf{A}\rangle^T \tag{S-38}$$

$$\langle\mathbf{B}^T\mathbf{A}^T\mathbf{D}^T\mathbf{D}\mathbf{A}\mathbf{B}\rangle = \text{tr}(\mathbf{\Sigma_B}\langle\mathbf{A}^T\mathbf{D}^T\mathbf{D}\mathbf{A}\rangle)\mathbf{I}_N + \mathbf{B}^T\langle\mathbf{A}^T\mathbf{D}^T\mathbf{D}\mathbf{A}\rangle\mathbf{B} \tag{S-39}$$

## 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.