

Supplementary Materials for “MM Algorithms for Variance Components Models”

Liuyi Hu Hua Zhou Jin Zhou Kenneth Lange

S.1 Proof for Lemma 2

Proof. Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a random normal vector with mean $\mathbf{0}$ and positive definite covariance matrix \mathbf{A} . Let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ be a random normal vector independent of \mathbf{X} with mean $\mathbf{0}$ and positive semidefinite covariance matrix \mathbf{B} having positive diagonal entries. Then $\mathbf{Z} = \mathbf{X} \odot \mathbf{Y}$ has covariances $E(Z_i Z_j) = E(X_i Y_i X_j Y_j) = E(X_i X_j) E(Y_i Y_j) = a_{ij} b_{ij}$. It follows that $\text{Cov}(\mathbf{Z}) = \mathbf{A} \odot \mathbf{B}$. To show $\mathbf{A} \odot \mathbf{B}$ is positive definite, suppose on the contrary that $\mathbf{v}^T (\mathbf{A} \odot \mathbf{B}) \mathbf{v} = \text{Var}(\mathbf{v}^T \mathbf{Z}) = 0$ for some $\mathbf{v} \neq \mathbf{0}$. Then

$$0 = \text{Var}(\mathbf{v}^T \mathbf{Z}) = E\left(\sum_i v_i X_i Y_i\right)^2 = E\left[\left(\sum_i v_i X_i Y_i\right)^2 \mid \mathbf{Y}\right] = E[(\mathbf{v} \odot \mathbf{Y})^T \mathbf{A} (\mathbf{v} \odot \mathbf{Y})]$$

implies $\mathbf{v} \odot \mathbf{Y} = \mathbf{0}$ with probability 1. Since $\mathbf{v} \neq \mathbf{0}$, $Y_i = 0$ with probability 1 for some i . This contradicts the assumption $b_{ii} = \text{Var}(Y_i) > 0$ for all i . \square

S.2 Objective Values in ANOVA Simulation Example

Table 1 summarizes the converged objective values for the two-way ANOVA example in Section 3. Reported in table are average and standard error based on 50 simulation replicates.

S.3 Proof of Theorem 1

We need three technical Lemmas to show the global convergence result in Theorem 1.

Lemma S.1. *Under Assumption 1 or 2, the log-likelihood function (1) is coercive in the sense that the super-level set $S_c = \{\boldsymbol{\sigma}^2 \geq \mathbf{0} : L(\boldsymbol{\sigma}^2) \geq c\}$ is compact for every c .*

Proof. Let us first prove the assertion when all of the covariance matrices \mathbf{V}_i are positive definite. If we set $r = \|\boldsymbol{\sigma}^2\|_1$ and $\alpha_i = r^{-1} \sigma_i^2$ for each i , then the log-likelihood satisfies

$$L(\boldsymbol{\sigma}^2) = -\frac{n}{2} \ln r - \frac{1}{2} \ln \det \left(\sum_{i=1}^m \alpha_i \mathbf{V}_i \right) - \frac{1}{2r} \mathbf{y}^T \left(\sum_{i=1}^m \alpha_i \mathbf{V}_i \right)^{-1} \mathbf{y}.$$

The functions $\ln \det \left(\sum_{i=1}^m \alpha_i \mathbf{V}_i \right)$ and $\mathbf{y}^T \left(\sum_{i=1}^m \alpha_i \mathbf{V}_i \right)^{-1} \mathbf{y}$ of $\boldsymbol{\alpha}$ are defined and continuous on the unit simplex and hence bounded there. The dominant term $-\frac{n}{2} \ln r$ of the loglikelihood tends to $-\infty$ as r tends to ∞ .

Table 1: MM, EM, Fisher scoring, and lme4 converge to similar objective values. Shown below are average objective values for fitting a two-way ANOVA model with $a = b = 5$ levels of both factors. Standard errors are given in parentheses.

σ_1^2/σ_e^2	Method	$c = \#$ observations per combination			
		5	10	20	50
0.00	MM	-176.67(7.94)	-353.59(10.19)	-713.42(14.90)	-1776.40(25.02)
	EM	-176.68(7.94)	-353.60(10.18)	-713.43(14.90)	-1776.41(25.02)
	FS	-176.67(7.94)	-353.59(10.19)	-713.42(14.90)	-1776.40(25.02)
	lme4	-176.67(7.94)	-353.59(10.19)	-713.42(14.90)	-1776.40(25.02)
0.05	MM	-181.06(7.24)	-365.39(10.92)	-722.16(15.36)	-1794.26(25.96)
	EM	-181.06(7.24)	-365.39(10.92)	-722.16(15.36)	-1794.26(25.96)
	FS	-181.06(7.24)	-365.39(10.92)	-722.16(15.36)	-1794.26(25.96)
	lme4	-181.06(7.24)	-365.39(10.92)	-722.16(15.36)	-1794.26(25.96)
0.10	MM	-185.10(6.77)	-368.83(10.28)	-726.21(13.35)	-1813.88(20.24)
	EM	-185.10(6.77)	-368.83(10.28)	-726.21(13.35)	-1813.88(20.24)
	FS	-185.10(6.77)	-368.83(10.28)	-726.21(13.35)	-1813.88(20.24)
	lme4	-185.10(6.77)	-368.83(10.28)	-726.21(13.35)	-1813.88(20.24)
1.00	MM	-204.05(8.18)	-392.35(10.94)	-754.95(13.93)	-1831.11(22.63)
	EM	-204.05(8.18)	-392.35(10.94)	-754.95(13.93)	-1831.12(22.63)
	FS	-204.05(8.18)	-392.35(10.94)	-754.95(13.93)	-1831.11(22.63)
	lme4	-204.05(8.18)	-392.35(10.94)	-754.95(13.93)	-1831.11(22.63)
10.00	MM	-233.65(7.90)	-416.56(11.79)	-777.85(15.65)	-1862.05(28.29)
	EM	-233.65(7.90)	-416.56(11.79)	-777.85(15.65)	-1862.05(28.29)
	FS	-233.65(7.90)	-416.56(11.79)	-777.85(15.65)	-1862.05(28.29)
	lme4	-233.65(7.90)	-416.56(11.79)	-777.85(15.65)	-1862.05(28.29)
20.00	MM	-242.21(8.20)	-424.68(10.11)	-795.63(15.56)	-1864.07(24.56)
	EM	-242.21(8.20)	-424.68(10.11)	-795.63(15.56)	-1864.07(24.56)
	FS	-242.21(8.20)	-424.68(10.11)	-795.63(15.56)	-1864.07(24.56)
	lme4	-242.21(8.20)	-424.68(10.11)	-795.63(15.56)	-1864.07(24.56)

To prove the assertion under Assumption 2, consider first the case $\mathbf{V}_1 = \mathbf{I}_n$. Setting $\alpha_i = \sigma_i^2/\sigma_1^2$ for $i = 2, \dots, m$ reduces the loglikelihood to

$$L(\sigma_1^2, \boldsymbol{\alpha}) = -\frac{n}{2} \ln \sigma_1^2 - \frac{1}{2} \ln \det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) - \frac{1}{2\sigma_1^2} \mathbf{y}^T \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right)^{-1} \mathbf{y}. \quad (1)$$

The middle term on the right satisfies

$$-\frac{1}{2} \ln \det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) \leq 0$$

because $\det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) \geq \det \mathbf{I}_n = 1$. Now let $\mathbf{U} = (\mathbf{U}_q, \mathbf{U}_{n-q})$ be an $n \times n$ orthogonal matrix whose left columns \mathbf{U}_q span \mathcal{H} and whose right columns \mathbf{U}_{n-q} span \mathcal{H}^\perp . The identity

$$\mathbf{U}^T \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) \mathbf{U} = \begin{pmatrix} \mathbf{I}_q + \sum_{i=2}^m \alpha_i \mathbf{U}_q^T \mathbf{V}_i \mathbf{U}_q & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-q} \end{pmatrix}$$

follows from the orthogonality relations $\mathbf{U}_{n-q}^T \mathbf{V}_i = \mathbf{U}_{n-q}^T \mathbf{U}_q = \mathbf{0}_{(n-q) \times n}$. This in turn implies

$$\begin{aligned} \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right)^{-1} &= \mathbf{U} \begin{pmatrix} (\mathbf{I}_q + \sum_{i=2}^m \alpha_i \mathbf{U}_q^T \mathbf{V}_i \mathbf{U}_q)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-q} \end{pmatrix} \mathbf{U}^T \\ &\succeq \mathbf{U} \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-q} \end{pmatrix} \mathbf{U}^T \\ &= \mathbf{U}_{n-q} \mathbf{U}_{n-q}^T. \end{aligned}$$

Therefore the quadratic term in equation (1) is bounded below by the positive constant

$$\mathbf{y}^T \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right)^{-1} \mathbf{y} \geq \mathbf{y}^T \mathbf{U}_{n-q} \mathbf{U}_{n-q}^T \mathbf{y} = \|\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y}\|^2 > 0.$$

Here the assumption $\mathbf{y} \notin \mathcal{H}$ guarantees the projection property $\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y} \neq \mathbf{0}$.

Next we show that the loglikelihood tends to $-\infty$ when σ_1^2 tends to 0 or ∞ or when $\|\boldsymbol{\alpha}\|_2$ tends to ∞ . The second of the two inequalities

$$\begin{aligned} L(\sigma_0^2, \boldsymbol{\alpha}) &\leq -\frac{n}{2} \ln \sigma_1^2 - \frac{1}{2} \ln \det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) - \frac{1}{2\sigma_1^2} \|\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y}\|^2 \\ &\leq -\frac{n}{2} \ln \sigma_1^2 - \frac{1}{2\sigma_1^2} \|\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y}\|^2 \end{aligned}$$

renders the claim about σ_1^2 obvious. To prove the claim about $\boldsymbol{\alpha}$, we make the worst case choice $\sigma_i^2 = \|\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y}\|^2$ in the first inequality. It follows that

$$L(\sigma_0^2, \boldsymbol{\alpha}) \leq -\frac{1}{2} \ln \det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) - \frac{n}{2} \ln \|\mathbb{P}_{\mathcal{H}^\perp} \mathbf{y}\|^2 - \frac{n}{2}.$$

If α_j tends to ∞ , then the inequality

$$-\frac{1}{2} \ln \det \left(\mathbf{I}_n + \sum_{i=2}^m \alpha_i \mathbf{V}_i \right) \leq -\frac{1}{2} \ln \det \left(\mathbf{I}_n + \alpha_j \mathbf{V}_j \right) = -\frac{1}{2} \sum_{k=1}^n \ln(1 + \alpha_j \lambda_{jk})$$

holds, where the λ_{jk} are the eigenvalues of \mathbf{V}_j . At least one of these eigenvalues is positive because \mathbf{V}_j is nontrivial. It follows that $L(\sigma_0^2, \boldsymbol{\alpha})$ tends to $-\infty$ in this case as well.

For the general case where \mathbf{V}_1 is non-singular but not necessarily \mathbf{I}_n , let $\mathbf{V}_1^{1/2}$ be the symmetric square root of \mathbf{V}_1 and write

$$\mathbf{V}_1 + \sum_{i=2}^m \sigma_i^2 \mathbf{V}_i = \mathbf{V}_1^{1/2} \left(\mathbf{I} + \sum_{i=2}^m \sigma_i^2 \mathbf{V}_1^{-1/2} \mathbf{V}_i \mathbf{V}_1^{-1/2} \right) \mathbf{V}_1^{1/2}.$$

The above arguments still apply since each $\mathbf{V}_1^{-1/2} \mathbf{V}_i \mathbf{V}_1^{-1/2}$ is nontrivial and \mathbf{y} belongs to the $\text{span}\{\mathbf{V}_2, \dots, \mathbf{V}_m\} = \mathcal{S}$ if and only if $\mathbf{V}_1^{-1/2} \mathbf{y}$ belongs to $\mathbf{V}_1^{-1/2} \mathcal{S} \mathbf{V}_1^{-1/2}$. \square

Lemma S.2. *The iterates possess the ascent property $L(M(\boldsymbol{\sigma}^{2(t)})) \geq L(\boldsymbol{\sigma}^{2(t)})$. Furthermore, when $L(M(\boldsymbol{\sigma}_*^2)) = L(\boldsymbol{\sigma}_*^2)$, $\boldsymbol{\sigma}_*^2$ fulfills the fixed point condition $M(\boldsymbol{\sigma}_*^2) = \boldsymbol{\sigma}_*^2$, and each component satisfies either (i) $\sigma_{*i}^2 = 0$ or (ii) $\sigma_{*i}^2 > 0$ and $\frac{\partial}{\partial \sigma_i^2} L(\boldsymbol{\sigma}_*^2) = 0$.*

Proof. The ascent property is built into any MM algorithm. Suppose $L(M(\boldsymbol{\sigma}_*^2)) = L(\boldsymbol{\sigma}_*^2)$ at a point $\boldsymbol{\sigma}_*^2 \in \mathbb{R}_+^m$. Then equality must hold in the string of inequalities (3). It follows that

$$g(M(\boldsymbol{\sigma}_*^2) \mid \boldsymbol{\sigma}_*^2) = g(\boldsymbol{\sigma}_*^2 \mid \boldsymbol{\sigma}_*^2).$$

$g(\cdot \mid \boldsymbol{\sigma}_*^2)$ has a unique maximum since its Hessian is diagonal with strictly negative entries, hence $M(\boldsymbol{\sigma}_*^2) = \boldsymbol{\sigma}_*^2$. If $\sigma_{*i}^2 > 0$, the stationarity condition

$$\frac{\partial}{\partial \sigma_i^2} L(\boldsymbol{\sigma}_*^2) = \frac{\partial}{\partial \sigma_i^2} g(\boldsymbol{\sigma}_*^2 \mid \boldsymbol{\sigma}_*^2) = 0$$

applies. The equivalence of the two displayed partial derivatives is a consequence of the fact that the difference $f(\boldsymbol{\sigma}^2) - g(\boldsymbol{\sigma}^2 \mid \boldsymbol{\sigma}_*^2)$ achieves its minimum of 0 at $\boldsymbol{\sigma}^2 = \boldsymbol{\sigma}_*^2$. \square

Lemma S.3. *The distance between successive iterates $\|\boldsymbol{\sigma}^{2(t+1)} - \boldsymbol{\sigma}^{2(t)}\|_2$ converges to 0.*

Proof. Suppose on the contrary that $\|\boldsymbol{\sigma}^{2(t+1)} - \boldsymbol{\sigma}^{2(t)}\|_2$ does not converge to 0. Then one can extract a subsequence $\{t_k\}_{k \geq 1}$ such that

$$\|\boldsymbol{\sigma}^{2(t_k+1)} - \boldsymbol{\sigma}^{2(t_k)}\|_2 \geq \epsilon > 0 \quad (2)$$

for all k . Let C_0 be the compact super-level set $\{\boldsymbol{\sigma}^2 : L(\boldsymbol{\sigma}^2) \geq L(\boldsymbol{\sigma}^{2(0)})\}$. Since the sequence $\{\boldsymbol{\sigma}^{2(t_k)}\}_{k \geq 1}$ is confined to C_0 , one can pass to a subsequence if necessary and assume that $\boldsymbol{\sigma}^{2(t_k)}$ converges to a limit $\boldsymbol{\sigma}_*^2$ and that $\boldsymbol{\sigma}^{2(t_k+1)}$ converges to a limit $\boldsymbol{\sigma}_{**}^2$. Taking limits in the relation $\boldsymbol{\sigma}^{2(t_k+1)} = M(\boldsymbol{\sigma}^{2(t_k)})$ and invoking the continuity $M(\boldsymbol{\sigma}^2)$ imply that $\boldsymbol{\sigma}_{**}^2 = M(\boldsymbol{\sigma}_*^2)$. Because the sequence $L(\boldsymbol{\sigma}^{2(t_k)})$ is monotonically increasing in k and bounded above on C_0 , it converges to a limit L_* . Hence, the continuity of $L(\boldsymbol{\sigma}^2)$ implies

$$L(\boldsymbol{\sigma}_*^2) = \lim_k L(\boldsymbol{\sigma}^{2(t_k)}) = L_* = \lim_k L(\boldsymbol{\sigma}^{2(t_k+1)}) = L(\boldsymbol{\sigma}_{**}^2) = L(M(\boldsymbol{\sigma}_*^2)).$$

Lemma S.2 therefore gives $\boldsymbol{\sigma}_{**}^2 = M(\boldsymbol{\sigma}_*^2) = \boldsymbol{\sigma}_*^2$, contradicting the bound $\|\boldsymbol{\sigma}_*^2 - \boldsymbol{\sigma}_{**}^2\|_2 \geq \epsilon$ entailed by inequality (2). \square

With Lemmas S.1-S.3, we are ready to prove Theorem 1.

Proof. The sequence $\{\boldsymbol{\sigma}^{2(t)}\}_{t \geq 0}$ is contained in the super-level compact set C_0 defined in Lemma S.3 and therefore admits a convergent subsequence $\boldsymbol{\sigma}^{2(t_k)}$ with limit $\boldsymbol{\sigma}^{2(\infty)}$. As argued in Lemma S.3, $L(\boldsymbol{\sigma}^{2(\infty)}) = L(M(\boldsymbol{\sigma}^{2(\infty)}))$. Lemma S.2 now implies that $\boldsymbol{\sigma}^{2(\infty)}$ is a fixed point of the algorithm map $M(\boldsymbol{\sigma}^2)$.

According to Ostrowski's theorem (Lange, 2010, Proposition 8.2.1), the set of limit points of a bounded sequence $\{\boldsymbol{\sigma}^{2(t)}\}_{t \geq 0}$ is connected and compact provided $\|\boldsymbol{\sigma}^{2(t+1)} - \boldsymbol{\sigma}^{2(t)}\|_2 \rightarrow 0$. If the set of fixed points is discrete, then the connected subset of limit points reduces to a single point. Hence, the bounded sequence $\boldsymbol{\sigma}^{2(t)}$ converges to this point. When the limit exists, one can check that $\boldsymbol{\sigma}^{2(\infty)}$ satisfies the KKT conditions by proving that each zero component of $\boldsymbol{\sigma}^{2(\infty)}$ has a non-positive partial derivative. Suppose on the contrary $\sigma_i^{2(\infty)} = 0$ and $\frac{\partial}{\partial \sigma_i^2} L(\boldsymbol{\sigma}^{2(\infty)}) > 0$. By continuity $\frac{\partial}{\partial \sigma_i^2} L(\boldsymbol{\sigma}^{2(t)}) > 0$ for all large t . Therefore, $\sigma_i^{2(t+1)} > \sigma_i^{2(t)}$ for all large t by the observation made after equation (9). This behavior is inconsistent with the assumption that $\sigma_i^{2(t)} \rightarrow 0$. \square

S.4 Proof of Theorem 2

MM algorithm: The minorizing function for the MM algorithm is

$$\begin{aligned}
& g_{\text{MM}}(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) \\
&= -\frac{1}{2} \text{tr}(\boldsymbol{\Omega}^{-(t)} \boldsymbol{\Omega}) - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \left(\sum_{i=1}^m \frac{\sigma_i^{4(t)}}{\sigma_i^2} \mathbf{V}_i \right) \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}) + c^{(t)} \\
&= \sum_{i=1}^m -\frac{\sigma_i^2}{2} \text{tr}(\boldsymbol{\Omega}^{-(t)} \mathbf{V}_i) - \frac{\sigma_i^{4(t)}}{2\sigma_i^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \mathbf{V}_i \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}) + c^{(t)},
\end{aligned}$$

where

$$c^{(t)} = -\frac{n}{2} \ln 2\pi - \frac{1}{2} \ln \det \boldsymbol{\Omega}^{(t)} + \frac{n}{2}.$$

Taking derivatives, we have

$$\begin{aligned}
\frac{\partial}{\partial \sigma_i^2} g_{\text{MM}}(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) &= -\frac{1}{2} \text{tr}(\boldsymbol{\Omega}^{-(t)} \mathbf{V}_i) + \frac{\sigma_i^{4(t)}}{2\sigma_i^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \mathbf{V}_i \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}), \\
\frac{\partial^2}{\partial \sigma_i^2 \partial \sigma_j^2} g_{\text{MM}}(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) &= \begin{cases} -\frac{\sigma_i^{4(t)}}{\sigma_i^6} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \mathbf{V}_i \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}) & i = j \\ 0 & i \neq j. \end{cases}
\end{aligned}$$

EM algorithm: Assume $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \sum_{i=1}^m \mathbf{Z}_i$, where $\mathbf{Z}_i \sim N(\mathbf{0}, \sigma_i^2 \mathbf{V}_i)$ are independent. Then the complete data is $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_m)$. From the information inequality, we have

$$L(\mathbf{y} | \boldsymbol{\sigma}^2) \geq Q(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) - Q(\boldsymbol{\sigma}^{2(t)} | \boldsymbol{\sigma}^{2(t)}) + L(\mathbf{y} | \boldsymbol{\sigma}^{2(t)}),$$

where

$$\begin{aligned}
L(\mathbf{y} | \boldsymbol{\sigma}^2) &= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln \det \boldsymbol{\Omega} - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \\
Q(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) &= -\frac{1}{2} \sum_{i=1}^m \left[\text{rank}(\mathbf{V}_i) \ln \sigma_i^2 + \frac{\sigma_i^{2(t)}}{\sigma_i^2} \text{rank}(\mathbf{V}_i) - \frac{\sigma_i^{4(t)}}{\sigma_i^2} \text{tr}(\boldsymbol{\Omega}^{-(t)} \mathbf{V}_i) \right] \\
&\quad - \frac{1}{2} \sum_{i=1}^m \left[\frac{\sigma_i^{4(t)}}{\sigma_i^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \mathbf{V}_i \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}) \right].
\end{aligned}$$

We derive this expression in Section S.6. The minorizing function

$$\begin{aligned}
& g_{\text{EM}}(\boldsymbol{\sigma} | \boldsymbol{\sigma}^{2(t)}) \\
&= Q(\boldsymbol{\sigma}^2 | \boldsymbol{\sigma}^{2(t)}) - Q(\boldsymbol{\sigma}^{2(t)} | \boldsymbol{\sigma}^{2(t)}) + L(\mathbf{y} | \boldsymbol{\sigma}^{2(t)}) \\
&= -\frac{1}{2} \sum_{i=1}^m \left[\text{rank}(\mathbf{V}_i) \ln \sigma_i^2 + \frac{\sigma_i^{2(t)}}{\sigma_i^2} \text{rank}(\mathbf{V}_i) - \frac{\sigma_i^{4(t)}}{\sigma_i^2} \text{tr}(\boldsymbol{\Omega}^{-(t)} \mathbf{V}_i) \right] \\
&\quad - \frac{1}{2} \sum_{i=1}^m \left[\frac{\sigma_i^{4(t)}}{\sigma_i^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)})^T \boldsymbol{\Omega}^{-(t)} \mathbf{V}_i \boldsymbol{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^{(t)}) \right] \\
&\quad - \frac{1}{2} \sum_{i=1}^m \left[-\text{rank}(\mathbf{V}_i) \ln \sigma_i^{2(t)} - \text{rank}(\mathbf{V}_i) \right] - \frac{1}{2} \left[n \ln(2\pi) + \ln \det \boldsymbol{\Omega}^{(t)} + n \right]
\end{aligned}$$

of the EM algorithms depends on σ^2 only through $Q(\sigma^2|\sigma^{2(t)})$. Taking derivatives, we have

$$\begin{aligned}
& \frac{\partial}{\partial \sigma_i^2} g_{\text{EM}}(\sigma^2|\sigma^{2(t)}) \\
&= -\frac{\text{rank}(\mathbf{V}_i)}{2\sigma_i^2} + \frac{\text{rank}(\mathbf{V}_i)\sigma_i^{2(t)} - \sigma_i^{4(t)}\text{tr}(\Omega^{-(t)}\mathbf{V}_i) + \sigma_i^{4(t)}(\mathbf{y} - \mathbf{X}\beta^{(t)})^T \Omega^{-(t)} \mathbf{V}_i \Omega^{-(t)} (\mathbf{y} - \mathbf{X}\beta^{(t)})}{2\sigma_i^4}, \\
& \frac{\partial^2}{\partial \sigma_i^2 \partial \sigma_j^2} g_{\text{EM}}(\sigma^2|\sigma^{2(t)}) \\
&= \begin{cases} \frac{\text{rank}(\mathbf{V}_i)}{2\sigma_i^4} - \frac{\text{rank}(\mathbf{V}_i)\sigma_i^{2(t)} - \sigma_i^{4(t)}\text{tr}(\Omega^{-(t)}\mathbf{V}_i) + \sigma_i^{4(t)}(\mathbf{y} - \mathbf{X}\beta^{(t)})^T \Omega^{-(t)} \mathbf{V}_i \Omega^{-(t)} (\mathbf{y} - \mathbf{X}\beta^{(t)})}{\sigma_i^6} & i = j \\ 0 & i \neq j. \end{cases}
\end{aligned}$$

EM vs MM: Let $\sigma^{2(\infty)}$ be a common limit point of EM and MM. By Lemma S.2, each component of $\sigma^{2(\infty)}$ is either 0 or has vanishing gradient. Therefore

$$\begin{aligned}
\frac{\partial^2}{(\partial \sigma_i^2)^2} g_{\text{EM}}(\sigma^2|\sigma^{2(\infty)}) \big|_{\sigma^2=\sigma^{2(\infty)}} &= -\frac{\text{rank}(\mathbf{V}_i)}{2\sigma_i^{4(\infty)}}, \\
\frac{\partial^2}{(\partial \sigma_i^2)^2} g_{\text{MM}}(\sigma^2|\sigma^{2(\infty)}) \big|_{\sigma^2=\sigma^{2(\infty)}} &= -\frac{\text{tr}(\Omega^{-(\infty)}\mathbf{V}_i)}{\sigma_i^{2(\infty)}}
\end{aligned}$$

and, when all the \mathbf{V}_i all non-singular,

$$\begin{aligned}
& \frac{1}{m} \sum_{i=1}^m \frac{[d^2 g_{\text{MM}}(\sigma^{2(\infty)}|\sigma^{2(\infty)})]_{ii}}{[d^2 g_{\text{EM}}(\sigma^{2(\infty)}|\sigma^{2(\infty)})]_{ii}} \\
&= \frac{2}{m} \sum_{i=1}^m \frac{\sigma_i^{2(\infty)} \text{tr}(\Omega^{-(\infty)}\mathbf{V}_i)}{\text{rank}(\mathbf{V}_i)} \\
&= \frac{2}{m} \leq 1.
\end{aligned}$$

S.5 Derivation of Algorithm 4

When there are $m = 2$ variance components $\Omega = \Gamma_1 \otimes \mathbf{V}_1 + \Gamma_2 \otimes \mathbf{V}_2$, repeated inversion of the $nd \times nd$ covariance matrix Ω can be reduced to one $d \times d$ (generalized) eigen-decomposition per iteration. The generalized eigen-decomposition of the matrix pair $(\mathbf{V}_1, \mathbf{V}_2)$ yields generalized eigenvalues $\mathbf{d} = (d_1, \dots, d_n)^T$ and generalized eigenvectors \mathbf{U} such that $\mathbf{U}^T \mathbf{V}_1 \mathbf{U} = \mathbf{D} = \text{diag}(\mathbf{d})$ and $\mathbf{U}^T \mathbf{V}_2 \mathbf{U} = \mathbf{I}$. Let the generalized eigen-decomposition of $(\Gamma_1^{(t)}, \Gamma_2^{(t)})$ be $(\Lambda^{(t)}, \Phi^{(t)})$ such that $\Phi^{(t)T} \Gamma_1^{(t)} \Phi^{(t)} = \Lambda^{(t)} = \text{diag}(\lambda^{(t)})$ and $\Phi^{(t)T} \Gamma_2 \Phi^{(t)} = \mathbf{I}_d$. It follows that

$$\begin{aligned}
\Omega^{(t)} &= (\Phi^{-(t)} \otimes \mathbf{U}^{-1})^T (\Lambda^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n) (\Phi^{-(t)} \otimes \mathbf{U}^{-1}) \\
\Omega^{-(t)} &= (\Phi^{(t)} \otimes \mathbf{U}) (\Lambda^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} (\Phi^{(t)} \otimes \mathbf{U})^T \\
\det \Omega^{(t)} &= \det(\Lambda^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n) \det(\Phi^{-(t)} \otimes \mathbf{U}^{-1})^T (\Phi^{-(t)} \otimes \mathbf{U}^{-1}) \\
&= \det(\Lambda^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n) \det(\Gamma_2^{(t)} \otimes \mathbf{V}_2) \\
&= \det(\Lambda^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n) \det(\Gamma_2^{(t)})^n \det(\mathbf{V}_2)^d.
\end{aligned}$$

To update the fixed effects \mathbf{B} given $\mathbf{\Gamma}_1^{(t)}$ and $\mathbf{\Gamma}_2^{(t)}$, the general least squares criterion is

$$\begin{aligned}
& \frac{1}{2} [\text{vec}(\mathbf{Y} - \mathbf{XB})]^T \mathbf{\Omega}^{-(t)} [\text{vec}(\mathbf{Y} - \mathbf{XB})] \\
&= \frac{1}{2} [\text{vec}(\mathbf{Y} - \mathbf{XB})]^T (\mathbf{\Phi}^{(t)} \otimes \mathbf{U}) (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} (\mathbf{\Phi}^{(t)} \otimes \mathbf{U})^T [\text{vec}(\mathbf{Y} - \mathbf{XB})] \\
&= \frac{1}{2} \text{vec}[\mathbf{U}^T (\mathbf{Y} - \mathbf{XB}) \mathbf{\Phi}^{(t)}]^T (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} \text{vec}[\mathbf{U}^T (\mathbf{Y} - \mathbf{XB}) \mathbf{\Phi}^{(t)}] \\
&= \frac{1}{2} [\text{vec}(\mathbf{U}^T \mathbf{Y} \mathbf{\Phi}^{(t)}) - (\mathbf{\Phi}^{(t)T} \otimes \mathbf{U}^T \mathbf{X}) \text{vec} \mathbf{B}]^T (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} \\
&\quad \cdot [\text{vec}(\mathbf{U}^T \mathbf{Y} \mathbf{\Phi}^{(t)}) - (\mathbf{\Phi}^{(t)T} \otimes \mathbf{U}^T \mathbf{X}) \text{vec} \mathbf{B}].
\end{aligned}$$

Minimization of this criterion reduces to a weighted least squares problem for the transformed responses $\tilde{\mathbf{Y}} = \mathbf{U}^T \mathbf{Y}$, transformed predictor matrix $\tilde{\mathbf{X}} = \mathbf{U}^T \mathbf{X}$, and observation weights $(\lambda_k^{(t)} d_i + 1)^{-1}$. To update $\mathbf{\Gamma}_1^{(t)}$ and $\mathbf{\Gamma}_2^{(t)}$, we need to evaluate the matrices \mathbf{M}_i and $\mathbf{\Gamma}_i^{(t)} \mathbf{R}^{(t)T} \mathbf{V}_i \mathbf{R}^{(t)} \mathbf{\Gamma}_i^{(t)}$ that appear in the stationarity condition (13).

Evaluation of \mathbf{M}_i : Note the (j, k) -th entry of \mathbf{M}_i is $\text{tr}(\mathbf{\Omega}_{jk}^{-(t)} \mathbf{V}_i)$, where $\mathbf{\Omega}_{jk}^{-(t)}$ is the (j, k) -th block of

$$\mathbf{\Omega}^{-(t)} = (\mathbf{\Phi}^{(t)} \otimes \mathbf{U}) (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} (\mathbf{\Phi}^{(t)} \otimes \mathbf{U})^T,$$

which can be expressed as

$$\mathbf{\Omega}_{jk}^{-(t)} = \sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \mathbf{U} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \mathbf{U}^T.$$

Therefore \mathbf{M}_1 has entries

$$\begin{aligned}
(\mathbf{M}_1)_{jk} &= \text{tr}(\mathbf{V}_1 \mathbf{\Omega}_{ij}^{-(t)}) \\
&= \text{tr} \left[\mathbf{U}^{-T} \mathbf{D} \mathbf{U}^{-1} \sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \mathbf{U} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \mathbf{U}^T \right] \\
&= \text{tr} \left[\sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \mathbf{D} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \right] \\
&= \sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \text{tr} \left[\mathbf{D} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \right],
\end{aligned}$$

and \mathbf{M}_2 has entries

$$\begin{aligned}
(\mathbf{M}_2)_{jk} &= \text{tr}(\mathbf{V}_2 \mathbf{\Omega}_{ij}^{-(t)}) \\
&= \text{tr} \left[\mathbf{U}^{-T} \mathbf{U}^{-1} \sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \mathbf{U} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \mathbf{U}^T \right] \\
&= \text{tr} \left[\sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} (\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \right] \\
&= \sum_{l=1}^d \phi_{jl}^{(t)} \phi_{lk}^{(t)} \text{tr}(\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1}.
\end{aligned}$$

Collectively we have

$$\begin{aligned} \mathbf{M}_1 &= \mathbf{\Phi}^{(t)} \text{diag} \left\{ \text{tr} \left[\mathbf{D}(\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \right] \right\} \mathbf{\Phi}^{(t)T} \\ \mathbf{M}_2 &= \mathbf{\Phi}^{(t)} \text{diag} \left[\text{tr}(\lambda_l^{(t)} \mathbf{D} + \mathbf{I}_n)^{-1} \right] \mathbf{\Phi}^{(t)T}. \end{aligned}$$

Evaluation of $\mathbf{\Gamma}_i^{(t)} \mathbf{R}^{(t)T} \mathbf{V}_i \mathbf{R}^{(t)} \mathbf{\Gamma}_i^{(t)}$: Write

$$\begin{aligned} \mathbf{\Gamma}_1^{(t)} \mathbf{R}^{(t)T} \mathbf{V}_1 \mathbf{R}^{(t)} \mathbf{\Gamma}_1^{(t)} &= \mathbf{N}_1^T \mathbf{N}_1 \\ \mathbf{\Gamma}_2^{(t)} \mathbf{R}^{(t)T} \mathbf{V}_2 \mathbf{R}^{(t)} \mathbf{\Gamma}_2^{(t)} &= \mathbf{N}_2^T \mathbf{N}_2, \end{aligned}$$

where

$$\begin{aligned} \mathbf{N}_1 &= \mathbf{D}^{1/2} \mathbf{U}^{-1} \mathbf{R}^{(t)} \mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)} \\ \mathbf{N}_2 &= \mathbf{U}^{-1} \mathbf{R}^{(t)} \mathbf{\Phi}^{-(t)T} \mathbf{\Phi}^{-(t)}. \end{aligned}$$

To further simplify, note

$$\begin{aligned} &\text{vec } \mathbf{N}_1 \\ &= (\mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)} \otimes \mathbf{D}^{1/2} \mathbf{U}^{-1}) \text{vec } \mathbf{R}^{(t)} \\ &= (\mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)} \otimes \mathbf{D}^{1/2} \mathbf{U}^{-1}) \mathbf{\Omega}^{-(t)} \text{vec}(\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) \\ &= (\mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)} \otimes \mathbf{D}^{1/2} \mathbf{U}^{-1}) (\mathbf{\Phi}^{(t)} \otimes \mathbf{U}) (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} (\mathbf{\Phi}^{(t)} \otimes \mathbf{U})^T \text{vec}(\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) \\ &= (\mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \otimes \mathbf{D}^{1/2}) (\mathbf{\Lambda}^{(t)} \otimes \mathbf{D} + \mathbf{I}_d \otimes \mathbf{I}_n)^{-1} \text{vec}(\mathbf{U}^T (\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) \mathbf{\Phi}^{(t)}) \\ &= (\mathbf{\Phi}^{-(t)T} \mathbf{\Lambda}^{(t)} \otimes \mathbf{D}^{1/2}) \text{vec}[\mathbf{U}^T (\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) \mathbf{\Phi}^{(t)} \oslash (\mathbf{d} \boldsymbol{\lambda}^{(t)T} + \mathbf{1}_n \mathbf{1}_d^T)] \\ &= \text{vec} \{ \mathbf{D}^{1/2} [(\mathbf{U}^T (\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) \mathbf{\Phi}^{(t)}) \oslash (\mathbf{d} \boldsymbol{\lambda}^{(t)T} + \mathbf{1}_n \mathbf{1}_d^T)] \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)} \}, \end{aligned}$$

where \oslash denotes a Hadamard quotient. Thus,

$$\mathbf{N}_1 = \mathbf{D}^{1/2} \{[(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{B}^{(t)}) \mathbf{\Phi}^{(t)}] \oslash (\mathbf{d} \boldsymbol{\lambda}^{(t)T} + \mathbf{1}_n \mathbf{1}_d^T)\} \mathbf{\Lambda}^{(t)} \mathbf{\Phi}^{-(t)},$$

and similarly

$$\mathbf{N}_2 = \{[(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{B}^{(t)}) \mathbf{\Phi}^{(t)}] \oslash (\mathbf{d} \boldsymbol{\lambda}^{(t)T} + \mathbf{1}_n \mathbf{1}_d^T)\} \mathbf{\Phi}^{-(t)}.$$

S.6 EM Algorithm for the Multivariate Response Model

In this section we review the derivation of the EM algorithm for the multivariate response model (Glanz and Carvalho, 2013; Reinsel, 1984). If the response matrix \mathbf{Y} can be written as the sum $\mathbf{Y} = \mathbf{X} \mathbf{B} + \mathbf{Z}_1 + \dots + \mathbf{Z}_m$ of independent random matrices with $\text{vec } \mathbf{Z}_i \sim N(\mathbf{0}, \mathbf{\Omega}_i)$, then $\text{vec } \mathbf{Y} \sim N(\text{vec}(\mathbf{X} \boldsymbol{\beta}), \mathbf{\Omega})$, where $\mathbf{\Omega} = \sum_{i=1}^m \mathbf{\Omega}_i$. Under the matrix normal assumption, $\mathbf{\Omega}_i = \mathbf{\Gamma}_i \otimes \mathbf{V}_i$. As in the text, the $p \times d$ coefficient matrix \mathbf{B} collects the fixed effects, the $\mathbf{\Gamma}_i$ are unknown $d \times d$ covariance matrices, and the \mathbf{V}_i are known $n \times n$ covariance matrices. The complete data log-likelihood for the unobserved \mathbf{Z}_i is

$$-\frac{1}{2} \sum_{i=1}^m \ln \det^+ \mathbf{\Omega}_i - \frac{1}{2} \sum_{i=1}^m \text{vec}(\mathbf{Z}_i)^T \mathbf{\Omega}_i^+ \text{vec}(\mathbf{Z}_i),$$

where $\det^+ \boldsymbol{\Omega}_i$ denotes the pseudo-determinant of $\boldsymbol{\Omega}_i$ and $\boldsymbol{\Omega}_i^+$ the pseudo-inverse of $\boldsymbol{\Omega}_i$. To compute the surrogate function for the EM algorithm, one needs the conditional expectations

$$\mathbb{E}(\text{vec } \mathbf{Z}_i \mid \mathbf{Y}, \boldsymbol{\theta}^{(t)}) = \boldsymbol{\Omega}_i^{(t)} \boldsymbol{\Omega}^{-(t)} \text{vec}(\mathbf{Y} - \mathbf{X} \mathbf{B}^{(t)}) = \mathbf{E}_i^{(t)}$$

and the conditional covariances

$$\text{Cov}(\text{vec } \mathbf{Z}_i \mid \mathbf{Y}, \boldsymbol{\theta}^{(t)}) = \boldsymbol{\Omega}_i^{(t)} - \boldsymbol{\Omega}_i^{(t)} \boldsymbol{\Omega}^{-(t)} \boldsymbol{\Omega}_i^{(t)} = \mathbf{F}_i^{(t)},$$

where $\boldsymbol{\theta}$ is the parameter vector. These are employed to compute the conditional second moments

$$\mathbb{E}(\text{vec } \mathbf{Z}_i \text{ vec } \mathbf{Z}_i^T \mid \mathbf{Y}, \boldsymbol{\theta}^{(t)}) = \mathbf{F}_i^{(t)} + \mathbf{E}_i^{(t)} (\mathbf{E}_i^{(t)})^T = \mathbf{G}_i^{(t)}.$$

Here the random vector \mathbf{Z}_i should be replaced by $\mathbf{Z}_m - \mathbf{X} \mathbf{B}^{(t)}$ when $i = m$.

One can readily check that $\boldsymbol{\Omega}_i^+ = \boldsymbol{\Gamma}_i^+ \otimes \mathbf{V}_i^+ = \boldsymbol{\Gamma}_i^{-1} \otimes \mathbf{V}_i^+$ for $\boldsymbol{\Gamma}_i$ invertible. Since the pseudo-determinant of a positive semidefinite matrix equals the product of its positive eigenvalues, the formulas

$$\begin{aligned} \det^+ \boldsymbol{\Omega}_i &= (\det \boldsymbol{\Gamma}_i)^{r_i} (\det^+ \mathbf{V}_i)^{s_i} \\ \ln \det^+ \boldsymbol{\Omega}_i &= r_i \ln \det \boldsymbol{\Gamma}_i + s_i \ln \det^+ \mathbf{V}_i \end{aligned}$$

apply, where $r_i = \text{rank}(\mathbf{V}_i^+)$ and $s_i = \text{rank}(\boldsymbol{\Gamma}_i^+)$. In the M step of the EM algorithm, one maximizes the surrogate

$$-\frac{1}{2} \sum_{i=1}^m r_i \ln \det \boldsymbol{\Gamma}_i - \frac{1}{2} \sum_{i=1}^m \text{tr}[(\boldsymbol{\Gamma}_i^{-1} \otimes \mathbf{V}_i^+) \mathbf{G}_i^{(t)}]. \quad (3)$$

For $\boldsymbol{\Gamma}_i$ unstructured, we substitute $\boldsymbol{\Lambda}_i = \boldsymbol{\Gamma}_i^{-1}$ and maximize with respect to $\boldsymbol{\Lambda}_i$. Fortunately, the next lemma can be invoked.

Lemma S.4. *If the matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} are $d \times d$, $n \times n$, and $dn \times dn$ respectively, then*

$$\text{tr}[(\mathbf{A} \otimes \mathbf{B}) \mathbf{C}^T] = \text{tr}\{(\mathbf{I}_d \otimes \mathbf{1}_n)^T [(\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{B}) \odot \mathbf{C}] (\mathbf{I}_d \otimes \mathbf{1}_n) \mathbf{A}^T\}.$$

Proof. This trace identity is essentially proved in the text. □

Lemma S.4 yields

$$\text{tr}[(\boldsymbol{\Lambda}_i \otimes \mathbf{V}_i^+) \mathbf{G}_i^{(t)}] = \text{tr}\{(\mathbf{I}_d \otimes \mathbf{1}_n)^T [(\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{V}_i^+) \odot \mathbf{G}_i^{(t)}] (\mathbf{I}_d \otimes \mathbf{1}_n) \boldsymbol{\Lambda}_i\}.$$

The stationarity condition

$$\mathbf{0} = \frac{1}{2} r_i \boldsymbol{\Lambda}_i^{-1} - \frac{1}{2} (\mathbf{I}_d \otimes \mathbf{1}_n)^T [(\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{V}_i^+) \odot \mathbf{G}_i^{(t)}] (\mathbf{I}_d \otimes \mathbf{1}_n).$$

now entails the update

$$\boldsymbol{\Gamma}_i^{(t+1)} = \frac{1}{r_i} (\mathbf{I}_d \otimes \mathbf{1}_n)^T [(\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{V}_i^+) \odot \mathbf{G}_i^{(t)}] (\mathbf{I}_d \otimes \mathbf{1}_n).$$

In the case $m = 1$, the single update reduces to

$$\mathbf{\Gamma}^{(t+1)} = \frac{1}{r}(\mathbf{Y} - \mathbf{X}\mathbf{B}^{(t)})^T \mathbf{V}^+ (\mathbf{Y} - \mathbf{X}\mathbf{B}^{(t)}),$$

which matches the earlier result of Glanz and Carvalho (2013). When $\mathbf{\Gamma}_i$ is the scalar σ_i^2 ,

$$\begin{aligned} \mathbf{E}_i^{(t)} &= \sigma_i^{2(t)} \mathbf{V}_i \mathbf{\Omega}^{-(t)} (\mathbf{y} - \mathbf{X}\mathbf{\beta}^{(t)}) \\ \mathbf{F}_i^{(t)} &= \sigma_i^{2(t)} \mathbf{V}_i - \sigma_i^{2(t)} \mathbf{V}_i \mathbf{\Omega}^{-(t)} \sigma_i^{2(t)} \mathbf{V}_i. \end{aligned}$$

One recovers the representation (3) by substituting these quantities in equation (3) and invoking the identities $\mathbf{V}_i \mathbf{V}_i^+ \mathbf{V}_i = \mathbf{V}_i$ and $\text{tr}(\mathbf{V}_i \mathbf{V}_i^+) = \text{rank}(\mathbf{V}_i)$ and the cyclic permutation property of the trace.

S.7 Proof of Lemma 3

Proof. Direct substitution shows that \mathbf{Y} solves the equivalent equation $\mathbf{X}\mathbf{B}\mathbf{X} = \mathbf{A}$. To show uniqueness, suppose $\mathbf{Y}^{-1}\mathbf{A}\mathbf{Y}^{-1} = \mathbf{B}$ and $\mathbf{Z}^{-1}\mathbf{A}\mathbf{Z}^{-1} = \mathbf{B}$. The equations

$$\begin{aligned} (\mathbf{B}^{1/2}\mathbf{Y}\mathbf{B}^{1/2})^2 &= \mathbf{B}^{1/2}\mathbf{Y}\mathbf{B}\mathbf{Y}\mathbf{B}^{1/2} = \mathbf{B}^{1/2}\mathbf{A}\mathbf{B}^{1/2} \\ (\mathbf{B}^{1/2}\mathbf{Z}\mathbf{B}^{1/2})^2 &= \mathbf{B}^{1/2}\mathbf{Z}\mathbf{B}\mathbf{Z}\mathbf{B}^{1/2} = \mathbf{B}^{1/2}\mathbf{A}\mathbf{B}^{1/2} \end{aligned}$$

imply $\mathbf{B}^{1/2}\mathbf{Y}\mathbf{B}^{1/2} = \mathbf{B}^{1/2}\mathbf{Z}\mathbf{B}^{1/2}$ by virtue of the uniqueness of symmetric square root. Since $\mathbf{B}^{-1/2}$ is positive definite, $\mathbf{Y} = \mathbf{Z}$. \square

S.8 Proof of Proposition 2

Proof. If \mathbf{V}_i has strictly positive diagonal entries, then so does $\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{V}_i$, and the Hadamard product $(\mathbf{1}_d \mathbf{1}_d^T \otimes \mathbf{V}_i) \odot \mathbf{\Omega}^{-(t)}$ is positive definite by Schur's lemma. Since the matrix $\mathbf{I}_d \otimes \mathbf{1}_n$ has full column rank d , the matrix \mathbf{M}_i is also positive definite. Finally, if no column of $\mathbf{R}^{(t)}$ lies in the null space of \mathbf{V}_i , and $\mathbf{\Gamma}^{(t)}$ is positive definite, then $\mathbf{\Gamma}_i^{(t)} \mathbf{R}^{(t)T} \mathbf{V}_i \mathbf{R}^{(t)} \mathbf{\Gamma}_i^{(t)}$ is positive definite. The second claim follows by induction and Lemma 3. \square

S.9 Proof of Lemma 4

Proof. Under the hypotheses, the representations $\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^+ \mathbf{A}^T = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ and $\mathbf{B}^+ = \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1}$ are well known. The choice $\mathbf{B}^+ \mathbf{A}^+ = \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ satisfies the four equations characterizing the pseudo-inverse of $\mathbf{A}\mathbf{B}$. \square

S.10 Proof of Lemma 5

Proof. Suppose \mathbf{A} has spectral decomposition $\sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T$. The matrix $\mathbf{P} = \sum_{\lambda_i > 0} \mathbf{u}_i \mathbf{u}_i^T$ projects onto the range of \mathbf{A} and therefore also projects onto the range of \mathbf{B} . It follows that $\mathbf{P}\mathbf{B} = \mathbf{B}$ and by symmetry that $\mathbf{B}\mathbf{P} = \mathbf{B}$. This allows us to write

$$\begin{aligned} &(\mathbf{B} + \epsilon \mathbf{I})(\mathbf{A} + \epsilon \mathbf{I})^{-1}(\mathbf{B} + \epsilon \mathbf{I}) \\ &= \mathbf{B}\mathbf{P}(\mathbf{A} + \epsilon \mathbf{I})^{-1}\mathbf{P}\mathbf{B} + \epsilon \mathbf{B}\mathbf{P}(\mathbf{A} + \epsilon \mathbf{I})^{-1} + \epsilon(\mathbf{A} + \epsilon \mathbf{I})^{-1}\mathbf{P}\mathbf{B} + \epsilon^2(\mathbf{A} + \epsilon \mathbf{I})^{-1}. \end{aligned}$$

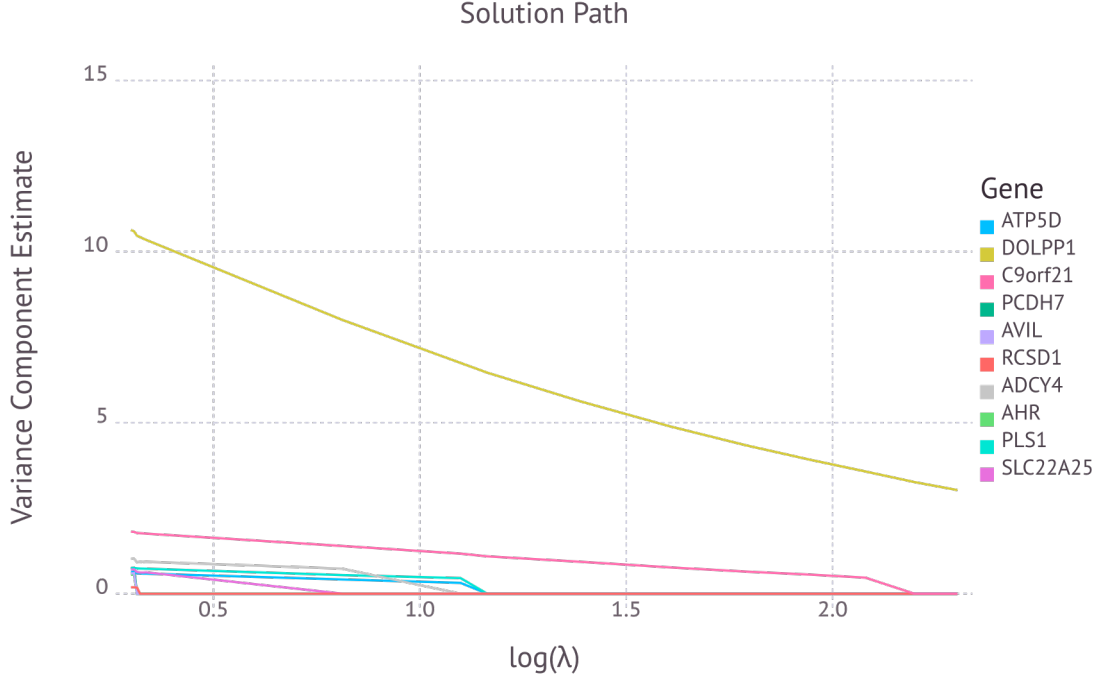


Figure 1: Solution path of the lasso penalized variance component model (17) indicates the top genes in an association study of 200 genes and the complex trait **height**.

The last three of these terms vanish as $\epsilon \downarrow 0$; the first term tends to the claimed limit. These assertions follow from the expressions

$$P(A + \epsilon I)^{-1}P = P(A + \epsilon I)^{-1} = (A + \epsilon I)^{-1}P = \sum_{\lambda_i > 0} \frac{1}{\lambda_i + \epsilon} \mathbf{u}_i \mathbf{u}_i^T$$

and $\epsilon^2(A + \epsilon I)^{-1} = \sum_i \frac{\epsilon^2}{\lambda_i + \epsilon} \mathbf{u}_i \mathbf{u}_i^T$. □

S.11 Proof of Lemma 6

Proof. In fact, both matrices have range equal to the range of \mathbf{Z} . The matrices \mathbf{Z} and $\mathbf{Z}\mathbf{R}^{1/2}$ clearly have the same range. Furthermore, the matrices $\mathbf{Z}\mathbf{R}^{1/2}$ and $\mathbf{Z}\mathbf{R}^{1/2}\mathbf{R}^{1/2}\mathbf{Z}^T$ also have the same range. □

S.12 Lasso solution path

Figure 1 displays the lasso solution path for the QTL example in Section 7.

References

Glanz, H. and Carvalho, L. (2013). An expectation-maximization algorithm for the matrix normal distribution. *arXiv*, 1309.6609.

- Lange, K. (2010). *Numerical Analysis for Statisticians*. Statistics and Computing. Springer, New York, second edition.
- Reinsel, G. (1984). Estimation and prediction in a multivariate random effects generalized linear model. *J. Amer. Statist. Assoc.*, 79(386):406–414.