

Supplementary material for Paper “Matrix-normal models for fMRI analysis”

Anonymous Author(s)

1 Appendix A : Matrix-normal intersubject functional connectivity and simultaneous modeling

Here we derive matrix-normal variants of two additional models from the literature, intersubject functional connectivity (Simony et al., 2016), and simultaneous modeling (Turner et al., 2013, 2014, 2016).

1.1 Matrix-normal intersubject functional connectivity

The goal of the ISFC method is to estimate a “shared stimulus-induced covariance matrix” in fMRI data as a way to measure functional connectivity between brain regions while abstracting over subject-specific connectivity patterns and extracting only the patterns that are consistent across subjects. The intuition behind the method is simple: it computes pairwise correlations between each subject’s patterns and averages them. To prove that the method is indeed free of subject-specific bias, Simony and colleagues frame their model in terms of a gaussian generative model. Here is this generative model, rewritten in the matrix-normal formalism:

$$\mathbf{A} \mid \mathbf{C} \sim \mathcal{MN}(0, \mathbf{C}, \mathbf{I}) \tag{1}$$

$$\mathbf{D}_i \mid \sigma_{\mathbf{D}}^2 \sim \mathcal{MN}(0, \sigma_{\mathbf{D}}^2 \mathbf{I}, \mathbf{I}) \tag{2}$$

$$\mathbf{S} \sim \mathcal{MN}(0, \mathbf{I}, \mathbf{I}) \tag{3}$$

$$\mathbf{E} \mid \mathbf{Q} \sim \mathcal{MN}(0, \mathbf{Q}, \mathbf{I}) \tag{4}$$

$$\mathbf{X}_i = (\mathbf{A} + \mathbf{D}_i)\mathbf{S} + \mathbf{E}_i \tag{5}$$

The “shared stimulus-induced covariance matrix” the method is intended to estimate is \mathbf{C} , the row covariance of the projection matrix into latent space. The somewhat redundant formulation is needed to motivate the closed-form estimator used in the original method. However, the formulation required for the closed-form estimator places severe restrictions on the projection matrix \mathbf{S} , both in terms of its rank (which must be full) and distribution (which is independent standard normal). We instead simplify the model

and integrate out the projection. Let $\mathbf{W}_i = \mathbf{A} + \mathbf{D}_i$, and rewrite:

$$\mathbf{A} \mid \mathbf{C} \sim \mathcal{MN}(0, \mathbf{C}, \mathbf{I}) \quad (6)$$

$$\mathbf{W}_i \mid \mathbf{A}, \sigma_{\mathbf{D}}^2 \sim \mathcal{MN}(\mathbf{A}, \sigma_{\mathbf{D}}^2 \mathbf{I}, \mathbf{I}) \quad (7)$$

$$\mathbf{S} \sim \mathcal{MN}(0, \mathbf{I}, \mathbf{I}) \quad (8)$$

$$\mathbf{X}_i \mid \mathbf{W}, \mathbf{S}, \mathbf{Q} \sim \mathcal{MN}(\mathbf{W}_i \mathbf{S}, \mathbf{Q}, \mathbf{I}) \quad (9)$$

Then marginalize A :

$$\mathbf{W}_i \mid C, \sigma_{\mathbf{D}}^2 \sim \mathcal{MN}(0, \mathbf{C} + \sigma_{\mathbf{D}}^2 \mathbf{I}, \mathbf{I}) \quad (10)$$

$$\mathbf{S} \sim \mathcal{MN}(0, \mathbf{I}, \mathbf{I}) \quad (11)$$

$$\mathbf{X}_i \mid \mathbf{W}, \mathbf{S}, \mathbf{Q} \sim \mathcal{MN}(\mathbf{W}_i \mathbf{S}, \mathbf{Q}, \mathbf{I}) \quad (12)$$

The resultant model is remarkably similar to MN-SRM: ISFC models the row (spatial) noise covariance as full-rank whereas MN-SRM models it as diagonal. MN-SRM models the shared response covariance as full-rank but ISFC models it as diagonal. Finally, and most importantly, MN-SRM models the projection into latent space as orthonormal whereas ISFC is specifically interested in its covariance (which MN-SRM can in fact estimate).

1.2 Matrix-normal simultaneous modeling

The simultaneous modeling framework (Turner et al., 2015) is organized around attempts to estimate the joint covariance of the vector $\{\psi_1, \psi_2, \dots, \psi_p, \phi_1, \phi_2, \dots, \phi_k\}$, which is a combined vector of cognitive model parameters ψ and features extracted from fMRI signal ϕ . As it is a broad framework, a number of specific instances have been provided, with specific cognitive models including accumulator models and signal detection theory models, and feature extraction mechanisms including ICA, PCA, and other methods.

There are a number of challenges with the current formulation of simultaneous modeling that we address: first, while the formulation in terms of correlations between brain and behavior allows for intuitive interpretation, it makes it challenging to regularize the model, or place priors on brain-behavior relationships, except for the special case of complete independence. Second, by performing the feature extraction in an unsupervised way, there is no guarantee that the features extracted will be relevant to the behavior or cognitive model; on the other hand, applying the framework to whole-brain data is not generally tractable, as it involves estimating a sizable covariance matrix by MCMC.

We show how matrix-normal simultaneous modeling can address all of these challenges. Since SM is a framework rather than one specific model, and no public implementation is available, we focus on a toy example to illustrate our contribution. We choose factor analysis as our factor model, leave the cognitive model unspecified for the derivation, which is applicable to any cognitive model, and any *linear* factor model.

First, we can use properties of partitioned Gaussians to write the conditional distribution of $\phi \mid \psi$, which is a simple linear regression:

$$\phi_i \mid \mu_\phi, \boldsymbol{\ell}, \Sigma_{\phi|\psi}, \Psi \sim \mathcal{N}(\mu_\phi + \Psi \boldsymbol{\ell}, \Sigma_{\phi|\psi}) \quad (13)$$

As we show in the supplement, the intercepts and slopes here map directly to the full covariance of the simultaneous modeling framework. Now we add the cognitive model and factorization, and stack into matrix-variate form. We also add an additional design matrix for observed stimulus features \mathbf{X} and its coefficient matrix β :

$$\mathbf{H} \mid \Psi \sim \text{Cog.}(\Psi, \mathbf{S}) \quad (14)$$

$$\Phi \mid \beta, \ell, \Sigma_{\phi|\psi}, \mathbf{S}, \Psi \sim \mathcal{MN}(\Psi\ell + \mathbf{X}\beta, \Sigma_{\Phi t}, \Sigma_{\Phi s}) \quad (15)$$

$$\mathbf{Y}^\top \mid \Phi, \mathbf{W}, \Sigma_s, \Sigma_t \sim \mathcal{MN}(\mathbf{W}\Phi, \Sigma_s, \Sigma_t) \quad (16)$$

This analysis combines the matrix-regression model for Φ and matrix-factor model for \mathbf{Y} . In this case since we only need the latent factors Φ to map to the cognitive parameters Ψ , we can marginalize over the factor mapping \mathbf{W} . For decoding cognitive parameters Ψ , we do not need the regression mapping either, so we marginalize over the coefficients, giving us a direct model from brain behavior via latent cognitive parameters and a neural factor space:

$$\mathbf{W} \sim \mathcal{MN}(0, \Sigma_s, \mathbf{I}) \quad (17)$$

$$\beta \sim \mathcal{MN}(0, \Sigma_{\Phi s}, \mathbf{U}) \quad (18)$$

$$\ell \sim \mathcal{MN}(0, \Sigma_{\Phi s}, \mathbf{V}) \quad (19)$$

$$\mathbf{H} \mid \Psi \sim \text{Cog.}(\Psi, \mathbf{X}) \quad (20)$$

$$\Phi \mid \Sigma_{\Phi s}, \Sigma_{\Phi t}, \mathbf{X}, \Psi \sim \mathcal{MN}(0, \Sigma_{\Phi s}, \Sigma_{\Phi t} + \mathbf{X}^\top \mathbf{U} \mathbf{X} + \Psi^\top \mathbf{V} \Psi) \quad (21)$$

$$\mathbf{Y}^\top \mid \Phi, \Sigma_s, \Sigma_t \sim \mathcal{MN}(0, \Sigma_s, \Sigma_t + \Phi^\top \Phi) \quad (22)$$

Given this marginalization, both the latent neural factors and the latent cognitive parameters appear in the model only as their inner products, and are perfectly nonidentifiable. Therefore, an equivalent model is a direct regression from voxels to cognitive parameters, marginalized over the mapping. This will be true for any linear factor model under marginalization:

$$\mathbf{H} \mid \Psi, \mathbf{X} \sim \text{Cog.}(\Psi, \mathbf{X}) \quad (23)$$

$$\mathbf{Y} \mid \Sigma_{\Phi s}, \Sigma_{\Phi t}, \mathbf{S}, \Psi, \mathbf{U}, \mathbf{V} \sim \mathcal{MN}(0, \Sigma_s, \Sigma_t + \mathbf{X} \mathbf{U}^\top \mathbf{X} + \Psi^\top \mathbf{V} \Psi) \quad (24)$$

In this view we have arrived again at an RSA-type intuition, namely that while it may very challenging to know the true projection from \mathbf{Y} to Ψ , mapping them on second-order statistics in time space can prove to be useful, especially as the dimensionality of \mathbf{Y} (and hence the marginalized-over mapping) grows.

With the mapping marginalized, we can still perform prediction from the model by maximizing the likelihood of the cognitive parameters corresponding to new data given parameters estimated previously:

$$\mathbf{Y}_{new} \mid \Psi \sim \mathcal{MN}(\mathbf{M}, \hat{\Sigma}_s, \mathbf{C}) \quad (25)$$

$$\mathbf{M} = \mathbf{Y}_{old} (\hat{\Sigma}_t + \hat{\Psi}^\top \hat{\Psi})^{-1} (\hat{\Psi}^\top \Psi) \quad (26)$$

$$\mathbf{C} = \hat{\Sigma}_t + \Psi^\top \Psi - (\Psi^\top \hat{\Psi}) (\hat{\Sigma}_t + \hat{\Psi}^\top \hat{\Psi})^{-1} (\hat{\Psi}^\top \Psi) \quad (27)$$

This maximization rotates the inner-products of the train and test sets into the same orientation. If the train and test sets have different numbers of TRs, we need to replace the temporal noise covariance matrix with a kernel function, but otherwise the derivation proceeds identically.

The resultant matrix-normal model mitigates the issues we identified previously: first, the only thing that scales with the number of voxels is the noise model rather than the mapping itself, allowing analysis to proceed using voxels directly assuming the noise model is efficient enough; second, it is targeted in that it automatically identifies the voxels most related to the cognitive model parameters; third, it is implicitly regularized via the priors on β and ℓ . As with all MNMs, it can also simultaneously handle both spatial and temporal noise in the fMRI signal.

2 Appendix B : Derivation of matrix normal identities

Consider the following three distributions:

$$\mathbf{X}_{ij} \sim \mathcal{MN}(\mathbf{A}_{ij}, \Sigma_{\mathbf{X}i}, \Sigma_{\mathbf{X}j}) \quad (28)$$

$$\mathbf{Y}_{jk} \sim \mathcal{MN}(\mathbf{B}_{jk}, \Sigma_{\mathbf{Y}j}, \Sigma_{\mathbf{Y}k}) \quad (29)$$

$$\mathbf{Z}_{ik} | \mathbf{X}_{ij}, \mathbf{Y}_{jk} \sim \mathcal{MN}(\mathbf{X}_{ij}\mathbf{Y}_{jk} + \mathbf{C}_{ik}, \Sigma_{\mathbf{Z}i}, \Sigma_{\mathbf{Z}k}) \quad (30)$$

We use lowercase subscripts to denote sizes, to make dimension constraints clearer. We first use the relationship between the matrix-normal and multivariate normal distribution to rewrite the densities in vectorized form. Next, we rewrite the vectorized product in the mean into kronecker form:

$$\text{vec}(\mathbf{Z}_{ik}) | \mathbf{X}_{ij}, \mathbf{Y}_{jk} \sim \mathcal{N}(\text{vec}(\mathbf{X}_{ij}\mathbf{Y}_{jk} + \mathbf{C}_{ik}), \Sigma_{\mathbf{Z}k} \otimes \Sigma_{\mathbf{Z}i}) \quad (31)$$

$$\text{vec}(\mathbf{Z}_{ik}) | \mathbf{X}_{ij}, \mathbf{Y}_{jk} \sim \mathcal{N}((\mathbf{I}_k \otimes \mathbf{X}_{ij})\text{vec}(\mathbf{Y}_{jk}) + \text{vec}(\mathbf{C}_{ik}), \Sigma_{\mathbf{Z}k} \otimes \Sigma_{\mathbf{Z}i}) \quad (32)$$

We recognize the resultant distribution as following into the form $y \sim \mathcal{N}(Mx + b, \Sigma)$. Now, the standard gaussian marginalization identity (e.g. Bishop et al. 2006) can be applied:

$$\text{vec}(\mathbf{Z}_{ik}) | \mathbf{X}_{ij} \sim \mathcal{N}((\mathbf{I}_k \otimes \mathbf{X}_{ij})\text{vec}(\mathbf{B}_{jk}) + \text{vec}(\mathbf{C}_{ik}), \Sigma_{\mathbf{Z}k} \otimes \Sigma_{\mathbf{Z}i} + (\mathbf{I}_k \otimes \mathbf{X}_{ij})(\Sigma_{\mathbf{Y}k} \otimes \Sigma_{\mathbf{Y}j})(\mathbf{I}_k \otimes \mathbf{X}_{ij})^\top) \quad (33)$$

We collect terms using the mixed-product property of kronecker products:

$$\text{vec}(\mathbf{Z}_{ik}) | \mathbf{X}_{ij} \sim \mathcal{N}(\text{vec}(\mathbf{X}_{ij}\mathbf{B}_{jk}) + \text{vec}(\mathbf{C}_{ik}), \Sigma_{\mathbf{Z}k} \otimes \Sigma_{\mathbf{Z}i} + \Sigma_{\mathbf{Y}k} \otimes \mathbf{X}_{ij}\Sigma_{\mathbf{Y}j}\mathbf{X}_{ij}^\top) \quad (34)$$

Now, we can see that the marginal density is a matrix-variate normal only if $\Sigma_{\mathbf{Z}k} = \Sigma_{\mathbf{Y}k}$ – that is, the variable we’re marginalizing over has the same covariance in the

dimension we are *not* marginalizing over as the marginal density. Otherwise the density is well-defined but not matrix-normal. If we let $\Sigma_k := \Sigma_{\mathbf{Z}_k} = \Sigma_{\mathbf{Y}_k}$, then we can factor out that term and rewrite the marginal density as a matrix normal:

$$\text{vec}(\mathbf{Z}_{ik}) \mid \mathbf{X}_{ij} \sim \mathcal{N}(\text{vec}(\mathbf{X}\mathbf{B}_{jk}) + \text{vec}(\mathbf{C}_{ik}), \Sigma_k \otimes \Sigma_{\mathbf{Z}_i} + \Sigma_k \otimes \mathbf{X}\Sigma_{\mathbf{Y}_j}\mathbf{X}^\top) \quad (35)$$

$$\text{vec}(\mathbf{Z}_{ik}) \mid \mathbf{X}_{ij} \sim \mathcal{N}(\text{vec}(\mathbf{X}\mathbf{B}_{jk}) + \text{vec}(\mathbf{C}_{ik}), \Sigma_k \otimes (\Sigma_{\mathbf{Z}_i} + \mathbf{X}\Sigma_{\mathbf{Y}_j}\mathbf{X}^\top)) \quad (36)$$

$$\mathbf{Z}_{ik} \mid \mathbf{X}_{ij} \sim \mathcal{MN}(\mathbf{X}\mathbf{B}_{jk} + \mathbf{C}_{ik}, \Sigma_{\mathbf{Z}_i} + \mathbf{X}\Sigma_{\mathbf{Y}_j}\mathbf{X}^\top, \Sigma_k) \quad (37)$$

Unlike the multivariate normal case, we can apply the same identity over either \mathbf{X} or \mathbf{Y} , since if $\mathbf{X} \sim \mathcal{MN}(M, U, V)$ then $\mathbf{X}^\top \sim \mathcal{MN}(M^\top, V, U)$. We write it directly below:

$$\mathbf{Z}_{ik}^\top \mid \mathbf{X}_{ij}, \mathbf{Y}_{jk} \sim \mathcal{MN}(\mathbf{Y}_{jk}^\top \mathbf{X}_{ij}^\top + \mathbf{C}_{ik}^\top, \Sigma_{\mathbf{Z}_k}, \Sigma_{\mathbf{Z}_i}) \quad (38)$$

$$\text{let } \Sigma_i := \Sigma_{\mathbf{Z}_i} = \Sigma_{\mathbf{X}_i} \quad (39)$$

$$\dots \quad (40)$$

$$\mathbf{Z}_{ik}^\top \mid \mathbf{Y}_{jk} \sim \mathcal{MN}(\mathbf{A}_{jk}^\top \mathbf{X}_{ij}^\top + \mathbf{C}_{ik}^\top, \Sigma_{\mathbf{Z}_k} + \mathbf{Y}^\top \Sigma_{\mathbf{Y}_j} \mathbf{Y}, \Sigma_{\mathbf{Z}_i}) \quad (41)$$

$$\mathbf{Z}_{ik} \mid \mathbf{Y}_{jk} \sim \mathcal{MN}(\mathbf{X}_{ij} \mathbf{A}_{jk} + \mathbf{C}_{ik}, \Sigma_{\mathbf{Z}_i}, \Sigma_{\mathbf{Z}_k} + \mathbf{Y}^\top \Sigma_{\mathbf{Y}_j} \mathbf{Y}) \quad (42)$$

Next, we do the same for the partitioned gaussian identity. First two vectorized matrix-normals that form our partition:

$$\mathbf{X}_{ij} \sim \mathcal{MN}(\mathbf{A}_{ij}, \Sigma_i, \Sigma_j) \rightarrow \text{vec}[\mathbf{X}_{ij}] \sim \mathcal{N}(\text{vec}[\mathbf{A}_{ij}], \Sigma_j \otimes \Sigma_i) \quad (43)$$

$$\mathbf{Y}_{ik} \sim \mathcal{MN}(\mathbf{B}_{ik}, \Sigma_i, \Sigma_k) \rightarrow \text{vec}[\mathbf{Y}_{ik}] \sim \mathcal{N}(\text{vec}[\mathbf{B}_{ik}], \Sigma_k \otimes \Sigma_i) \quad (44)$$

$$\begin{bmatrix} \text{vec}[\mathbf{X}_{ij}] \\ \text{vec}[\mathbf{Y}_{ik}] \end{bmatrix} \sim \mathcal{N} \left(\text{vec} \begin{bmatrix} \mathbf{A}_{ij} \\ \mathbf{B}_{ik} \end{bmatrix}, \begin{bmatrix} \Sigma_j \otimes \Sigma_i & \Sigma_{jk} \otimes \Sigma_i \\ \Sigma_{kj} \otimes \Sigma_i & \Sigma_k \otimes \Sigma_i \end{bmatrix} \right) \quad (45)$$

We apply the standard partitioned Gaussian identity and simplify using the properties of the vec operator and the mixed product property of kronecker products:

$$\text{vec}[\mathbf{X}_{ij}] \mid \text{vec}[\mathbf{Y}_{ik}] \sim \mathcal{N}(\text{vec}[\mathbf{A}_{ij}] + (\Sigma_{jk} \otimes \Sigma_i)(\Sigma_k^{-1} \otimes \Sigma_i^{-1})(\text{vec}[\mathbf{Y}_{ik}] - \text{vec}[\mathbf{B}_{ik}]), \quad (46)$$

$$\Sigma_j \otimes \Sigma_i - (\Sigma_{jk} \otimes \Sigma_i)(\Sigma_k^{-1} \otimes \Sigma_i^{-1})(\Sigma_{kj} \otimes \Sigma_i) \quad (47)$$

$$= \mathcal{N}(\text{vec}[\mathbf{A}_{ij}] + (\Sigma_{jk}\Sigma_k^{-1} \otimes \Sigma_i\Sigma_i^{-1})(\text{vec}[\mathbf{Y}_{ik}] - \text{vec}[\mathbf{B}_{ik}]), \quad (48)$$

$$\Sigma_j \otimes \Sigma_i - (\Sigma_{jk}\Sigma_k^{-1}\Sigma_{kj} \otimes \Sigma_i\Sigma_i^{-1}\Sigma_i) \quad (49)$$

$$= \mathcal{N}(\text{vec}[\mathbf{A}_{ij}] + (\Sigma_{jk}\Sigma_k^{-1} \otimes \mathbf{I})(\text{vec}[\mathbf{Y}_{ik}] - \text{vec}[\mathbf{B}_{ik}]), \quad (50)$$

$$\Sigma_j \otimes \Sigma_i - (\Sigma_{jk}\Sigma_k^{-1}\Sigma_{kj} \otimes \Sigma_i) \quad (51)$$

$$= \mathcal{N}(\text{vec}[\mathbf{A}_{ij}] + \text{vec}[\mathbf{Y}_{ik} - \mathbf{B}_{ik}\Sigma_k^{-1}\Sigma_{kj}], (\Sigma_j - \Sigma_{jk}\Sigma_k^{-1}\Sigma_{kj}) \otimes \Sigma_i) \quad (52)$$

Next, we recognize that this multivariate gaussian is equivalent to the following matrix variate gaussian:

$$\mathbf{X}_{ij} | \mathbf{Y}_{ik} \sim \mathcal{MN}(\mathbf{A}_{ij} + (\mathbf{Y}_{ik} - \mathbf{B}_{ik})\Sigma_k^{-1}\Sigma_{kj}, \Sigma_i, \Sigma_j - \Sigma_{jk}\Sigma_k^{-1}\Sigma_{kj}) \quad (53)$$

The conditional in the other direction can be written by working through the same algebra:

$$\mathbf{Y}_{ik} | \mathbf{X}_{ij} \sim \mathcal{MN}(\mathbf{B}_{ik} + (\mathbf{X}_{ij} - \mathbf{A}_{ij})\Sigma_j^{-1}\Sigma_{jk}, \Sigma_i, \Sigma_k - \Sigma_{kj}\Sigma_j^{-1}\Sigma_{jk}) \quad (54)$$

Finally, vertical rather than horizontal concatenation (yielding a partitioned row rather than column covariance) can be written by recognizing the behavior of the matrix normal under transposition:

$$\mathbf{X}_{ji}^\top | \mathbf{Y}_{ki}^\top \sim \mathcal{MN}(\mathbf{A}_{ji}^\top + \Sigma_{jk}\Sigma_k^{-1}(\mathbf{Y}_{ki}^\top - \mathbf{B}_{ki}^\top), \Sigma_j - \Sigma_{jk}\Sigma_k^{-1}\Sigma_{kj}, \Sigma_i) \quad (55)$$

$$\mathbf{Y}_{ki}^\top | \mathbf{X}_{ji}^\top \sim \mathcal{MN}(\mathbf{B}_{ki}^\top + \Sigma_{kj}\Sigma_j^{-1}(\mathbf{X}_{ji}^\top - \mathbf{A}_{ji}^\top), \Sigma_k - \Sigma_{kj}\Sigma_j^{-1}\Sigma_{jk}, \Sigma_i) \quad (56)$$

3 Appendix C : Expectation Conditional Maximization (ECM) derivation for Matrix-Normal Shared Response Model

The Q function, marginalized \mathbf{W}

$$\mathbf{X} \sim \mathcal{MN}(\mathbf{WS} + \mathbf{b}\mathbf{1}^\top, \rho \otimes \Sigma_v, \Sigma_t) \quad (57)$$

$$\mathbf{S} \sim \mathcal{MN}(0, \mathbf{I}, \Sigma_t) \quad (58)$$

$$\mathbf{W} \sim \mathcal{MN}(0, \rho \otimes \Sigma_v, \mathbf{I}) \quad (59)$$

$$\begin{aligned} \mathcal{L} &:= \mathbb{E}_{p(\mathbf{W}|\mathbf{X},\theta')} \log p(\mathbf{X}, \mathbf{W} | \theta) = \frac{1}{2} \mathbb{E} \left[nv \log |\Sigma_t^{-1}| + tv \log |\rho^{-1}| + tn \log |\Sigma_v^{-1}| \right. \\ &\quad - \text{Tr} \left[\Sigma_t^{-1} (\mathbf{X} - \mathbf{WS} - \mathbf{b}\mathbf{1}^\top)^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{WS} - \mathbf{b}\mathbf{1}^\top) \right] \\ &\quad + kv \log |\rho^{-1}| + kn \log |\Sigma_v^{-1}| - \text{Tr} \left[\Sigma_w^{-1} \mathbf{W}^\top (\rho \otimes \Sigma_v)^{-1} \mathbf{W} \right] \\ &\quad \left. + k \log |\Sigma_t^{-1}| - \text{Tr} [\Sigma_t^{-1} \mathbf{S}^\top \mathbf{S}] \right] + \text{const.}_\theta \quad (60) \\ &= \frac{1}{2} \left[(nv + k) \log |\Sigma_t^{-1}| + v(k + t) \log |\rho^{-1}| + n(k + t) \log |\Sigma_v^{-1}| \right. \\ &\quad - \text{Tr} \left[\Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \\ &\quad - \text{Tr} \left[\Sigma_w^{-1} \mathbf{W}'^\top (\rho \otimes \Sigma_v)^{-1} \mathbf{W}' \right] - \text{Tr} [\rho^{-1} \rho'] \text{Tr} [\Sigma_v^{-1} \Sigma'_v] \text{Tr} [\Sigma_w' \mathbf{S}^\top \Sigma_t^{-1} \mathbf{S}] \\ &\quad \left. - \text{Tr} [\rho^{-1} \rho'] \text{Tr} [\Sigma_v^{-1} \Sigma'_v] \text{Tr} [\Sigma_w^{-1} \Sigma'_w] - \text{Tr} [\Sigma_t^{-1} \mathbf{S}^\top \mathbf{S}] \right] + \text{const.}_\theta \quad (61) \end{aligned}$$

The sufficient statistics are:

$$\mathbf{W} | \mathbf{X}, \theta \sim \mathcal{MN}(\mathbf{W}', \rho'_w \otimes \Sigma'_{vw}, \Sigma'_w) \quad (62)$$

$$\Sigma'_w := \mathbf{I} - S(\Sigma_t + \mathbf{S}^\top \mathbf{S})^{-1} \mathbf{S}^\top = (\mathbf{I} + \mathbf{S} \Sigma_t^{-1} \mathbf{S}^\top)^{-1} \quad (63)$$

$$\rho'_w := \rho_w \quad (64)$$

$$\Sigma'_{vw} := \Sigma_{vw} \quad (65)$$

$$\mathbf{W}' = (\mathbf{X} - \mathbf{b}\mathbf{1}^\top) \mathbf{S}^\top (\Sigma_t + \mathbf{S}^\top \mathbf{S})^{-1} = (\mathbf{X} - \mathbf{b}\mathbf{1}^\top) \Sigma_t^{-1} \mathbf{S}^\top \Sigma'_w \quad (66)$$

3.1 Gradients for \mathbf{S}

$$d_{\mathbf{S}}\mathcal{L} = \frac{1}{2}d \left[-\text{Tr} \left[\Sigma_t^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \right] \quad (67)$$

$$- \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}] \text{Tr}[\Sigma'_w\mathbf{S}\Sigma_t^{-1}\mathbf{S}] - \text{Tr}[\Sigma_t^{-1}\mathbf{S}^\top\mathbf{S}] \quad (68)$$

$$= \frac{1}{2} \left[-2 \text{Tr} \left[\Sigma_t^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1}(\mathbf{W}'d\mathbf{S}) \right] \right] \quad (69)$$

$$- 2 \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}] \text{Tr}[\Sigma'_w d\mathbf{S}\Sigma_t^{-1}\mathbf{S}^\top]^{-1}\Sigma'_w] - 2 \text{Tr}[\Sigma_t^{-1}\mathbf{S}^\top d\mathbf{S}] \quad (70)$$

$$= - \text{Tr} \left[\Sigma_t^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1}\mathbf{W}'d\mathbf{S} \right] \quad (71)$$

$$- \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}] \text{Tr}[\Sigma_t^{-1}\mathbf{S}^\top\Sigma'_w d\mathbf{S}] - \text{Tr}[\Sigma_t^{-1}\mathbf{S}^\top d\mathbf{S}] \quad (72)$$

$$(73)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{S}} = \mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})\Sigma_t^{-1} - \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}]\Sigma'_w\mathbf{S}\Sigma_t^{-1} - \mathbf{S}\Sigma_t^{-1} \quad (74)$$

$$= \mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})(\mathbf{X} - \mathbf{b}\mathbf{1}^\top) - (\mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})\mathbf{W}'\mathbf{S} - \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}]\Sigma'_w\mathbf{S} - \mathbf{S} \quad (75)$$

$$= \mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})(\mathbf{X} - \mathbf{b}\mathbf{1}^\top) - (\mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})\mathbf{W}' + \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}]\Sigma'_w + 1)\mathbf{S} \quad (76)$$

$$(77)$$

$$\hat{\mathbf{S}} = (\mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})\mathbf{W}' + \text{Tr}[\rho^{-1}\rho'_w] \text{Tr}[\Sigma_v^{-1}\Sigma'_{vw}]\Sigma'_w + 1)^{-1}\mathbf{W}'^\top (\rho^{-1} \otimes \Sigma_v^{-1})(\mathbf{X} - \mathbf{b}\mathbf{1}^\top) \quad (78)$$

3.2 Gradients for \mathbf{b}

$$d_{\mathbf{b}}\mathcal{L} = -\frac{1}{2} \text{Tr} \left[\Sigma_t^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \quad (79)$$

$$= - \text{Tr} \left[\mathbf{1}^\top \Sigma_t^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1}d\mathbf{b} \right] \quad (80)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = (\rho \otimes \Sigma_v)^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})\Sigma_t^{-1}\mathbf{1} \quad (81)$$

$$(82)$$

$$0 = (\rho \otimes \Sigma_v)^{-1}(\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})\Sigma_t^{-1}\mathbf{1} \quad (83)$$

$$\mathbf{b}\mathbf{1}^\top \Sigma_t^{-1}\mathbf{1} = (\mathbf{X} - \mathbf{W}'\mathbf{S})\Sigma_t^{-1}\mathbf{1} \quad (84)$$

$$\hat{\mathbf{b}} = \frac{(\mathbf{X} - \mathbf{W}'\mathbf{S})\Sigma_t^{-1}\mathbf{1}}{\sum \Sigma_t^{-1}} \quad (85)$$

3.3 Gradients for Σ_t

$$d_{\Sigma_t^{-1}} \mathcal{L} = \frac{1}{2} d \left[(nv + k) \log |\Sigma_t^{-1}| + v(k + t) \log |\rho^{-1}| + n(k + t) \log |\Sigma_v^{-1}| \right] \quad (86)$$

$$- \text{Tr} \left[\Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \quad (87)$$

$$- \text{Tr} \left[\Sigma_w^{-1} \mathbf{W}'^\top (\rho \otimes \Sigma_v)^{-1} \mathbf{W}' \right] - \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \text{Tr}[\Sigma'_w \mathbf{S} \Sigma_t^{-1} \mathbf{S}^\top] \quad (88)$$

$$- \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \text{Tr}[\Sigma_w^{-1} \Sigma'_w] - \text{Tr}[\Sigma_t^{-1} \mathbf{S}^\top \mathbf{S}] \quad (89)$$

$$= \frac{1}{2} \left[(nv + k) \text{Tr}[\Sigma_t d\Sigma_t^{-1}] - \text{Tr} \left[d\Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \right] \quad (90)$$

$$- \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \text{Tr}[\Sigma'_w \mathbf{S} d\Sigma_t^{-1} \mathbf{S}^\top] - \text{Tr}[d\Sigma_t^{-1} \mathbf{S}^\top \mathbf{S}] \quad (91)$$

$$(92)$$

$$\frac{\partial \mathcal{L}}{\partial \Sigma_t^{-1}} = \frac{1}{2} \left[(nv + k) \Sigma_t - (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \quad (93)$$

$$- \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \mathbf{S}^\top \Sigma'_w \mathbf{S} - \mathbf{S}^\top \mathbf{S} \quad (94)$$

$$(95)$$

$$\widehat{\Sigma_t^{-1}} = \left(\frac{1}{nv + k} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right) \quad (96)$$

$$- \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \mathbf{S}^\top \Sigma'_w \mathbf{S} - \mathbf{S}^\top \mathbf{S} \Big)^{-1} \quad (97)$$

3.4 Gradients for Σ_v

Here again we assume ρ is diagonal, in which case:

$$d_{\Sigma_v^{-1}} \mathcal{L} = d \frac{1}{2} [(nv + k) \log |\Sigma_t^{-1}| + v(k + t) \log |\rho^{-1}| + n(k + t) \log |\Sigma_v^{-1}|] \quad (98)$$

$$- \text{Tr} \left[\Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top (\rho \otimes \Sigma_v)^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \quad (99)$$

$$- \text{Tr} \left[\Sigma_w^{-1} \mathbf{W}'^\top (\rho \otimes \Sigma_v)^{-1} \mathbf{W}' \right] - \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \text{Tr}[\Sigma'_w \mathbf{S}^\top \Sigma_t^{-1} \mathbf{S}] \quad (100)$$

$$- \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma_v^{-1} \Sigma'_v] \text{Tr}[\Sigma_w^{-1} \Sigma'_w] - \text{Tr}[\Sigma_t^{-1} \mathbf{S}^\top \mathbf{S}] \quad (101)$$

$$= \frac{1}{2} [n(k + t) \Sigma_v d \Sigma_v^{-1}] \quad (102)$$

$$- \sum_j \tau_j \text{Tr} \left[\Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top d \Sigma_v^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \right] \quad (103)$$

$$- \sum_j \tau_j \text{Tr} \left[\Sigma_w^{-1} \mathbf{W}'^\top d \Sigma_v^{-1} \mathbf{W}' \right] - \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma'_v d \Sigma_v^{-1}] \text{Tr}[\Sigma'_w (\mathbf{I} + \mathbf{S}^\top \Sigma_t^{-1} \mathbf{S})] \quad (104)$$

$$(105)$$

$$\frac{\partial \mathcal{L}}{\partial \Sigma_v^{-1}} = \frac{1}{2} [n(k + t) \Sigma_v] \quad (106)$$

$$- \sum_j \tau_j (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top \quad (107)$$

$$- \sum_j \tau_j \left[\mathbf{W}' \Sigma_w^{-1} \mathbf{W}'^\top - \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma'_w (\mathbf{I} + \mathbf{S} \Sigma_t^{-1} \mathbf{S}^\top)] \Sigma'_v \right] \quad (108)$$

$$(109)$$

$$\widehat{\Sigma_v^{-1}} = \left(\frac{1}{n(k + t)} \sum_j \tau_j (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S}) \Sigma_t^{-1} (\mathbf{X} - \mathbf{b}\mathbf{1}^\top - \mathbf{W}'\mathbf{S})^\top \right) \quad (110)$$

$$+ \left(\sum_j \tau_j \left[\mathbf{W}' \Sigma_w^{-1} \mathbf{W}'^\top - \text{Tr}[\rho^{-1} \rho'] \text{Tr}[\Sigma'_w (\mathbf{I} + \mathbf{S} \Sigma_t^{-1} \mathbf{S}^\top)] \Sigma'_v \right] \right)^{-1} \quad (111)$$

3.5 Constrained covariances

For template constraints (e.g. diagonal, blocked, banded), we can elementwise-multiply the gradient by a template matrix, and construct the constrained update.

Algorithm 1 Solve $x = (L_0 \otimes L_1 \otimes \cdots \otimes L_{n-1}) \setminus y$

```

1: Input: vector  $y$ , matrices  $L_0, L_1, \dots, L_{n-1}$ 
2: Output: vector  $x$ 
3: if  $n == 1$  then
4:   return matrix_triangular_solve( $L_0, y$ )
5: else
6:    $x = y$ 
7:    $na = \dim(L_0)$ 
8:    $nb = \dim(L_1) \times \dim(L_2) \times \cdots \times \dim(L_{n-1})$ 
9:   for  $i = 0$  to  $na - 1$  do
10:     $t = x[i * nb : (i + 1) * nb] / L_0[i, i]$ 
11:     $x[i * nb : (i + 1) * nb] = (L_1 \otimes \cdots \otimes L_{n-1}) \setminus t$ 
12:    for  $j = i + 1$  to  $na - 1$  do
13:       $x[j * nb : (j + 1) * nb] - = L_0[j, i] * t$ 
14:    end for
15:  end for
16:  return  $x$ 
17: end if

```

4 Appendix D : Algorithm for solving kronecker factored matrices

In algorithm 1, we show how to efficiently solve for a lower triangular matrix that is the kronecker product of smaller lower triangular matrices.

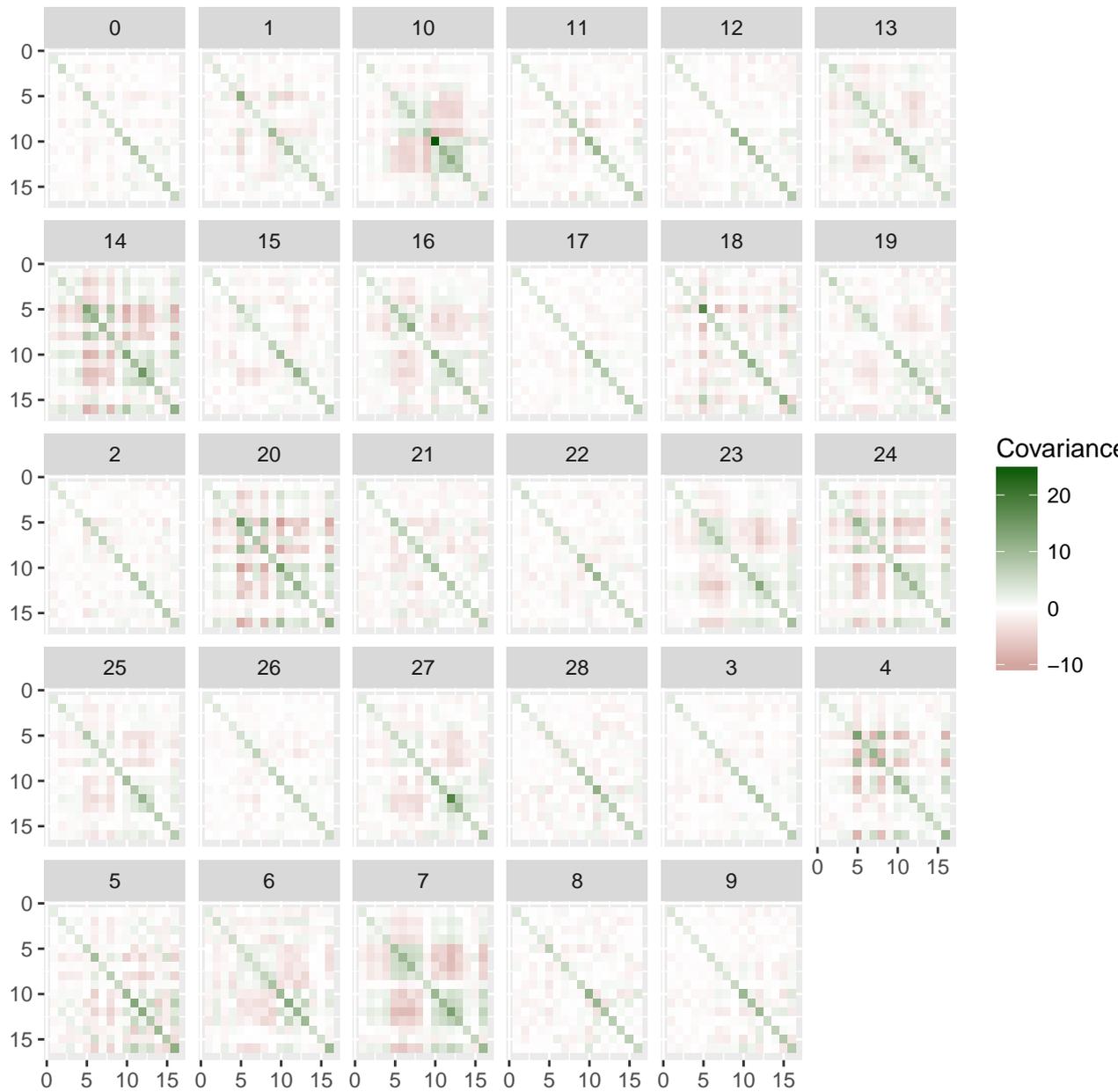
Since the cholesky of a kronecker product is the kronecker product of its cholesky factors, we avoid computing the cholesky factorization of a large matrix and instead only cholesky factorize the individual factors. Algorithm 1 is recursive: line 11 calls the same function but with one less kronecker factor. The masked variant of the algorithm is similar except for lines 4, 11 and 13. Lines 4 and 11 now perform matrix solves with a mask. Line 13 multiplies $L_0[j, i]$ not with t but with $t' = (L_1 \otimes \cdots \otimes L_{n-1}) \cdot x[i * nb : (i + 1) * nb]$. t and t' are identical when no rows and columns are masked, but differ when some of them are masked. Solving $\Sigma^{-1} \mathbf{X}$ now involves the following steps - (1) Cholesky factorize the kronecker factor matrices. (2) Use algorithm 1 to solve $Z = (L_0 \otimes L_1 \otimes \cdots \otimes L_{n-1}) \setminus X$. (3) Apply the corresponding upper triangular variant to solve $(L_0 \otimes L_1 \otimes \cdots \otimes L_{n-1})^T \setminus Z$.

We can calculate log-determinant for kronecker products as follows. After cholesky factorization, $\log |\Sigma| = 2 \cdot \sum_i ((\log |L_i|) (\prod_{j, j \neq i} \dim(L_j)))$. $\log |L_i|$ is easy to calculate for a triangular matrix L_i . For masked kronecker product, the latter product term in the previous expression is replaced by counting the number of valid rows/columns corresponding to that element in the mask.

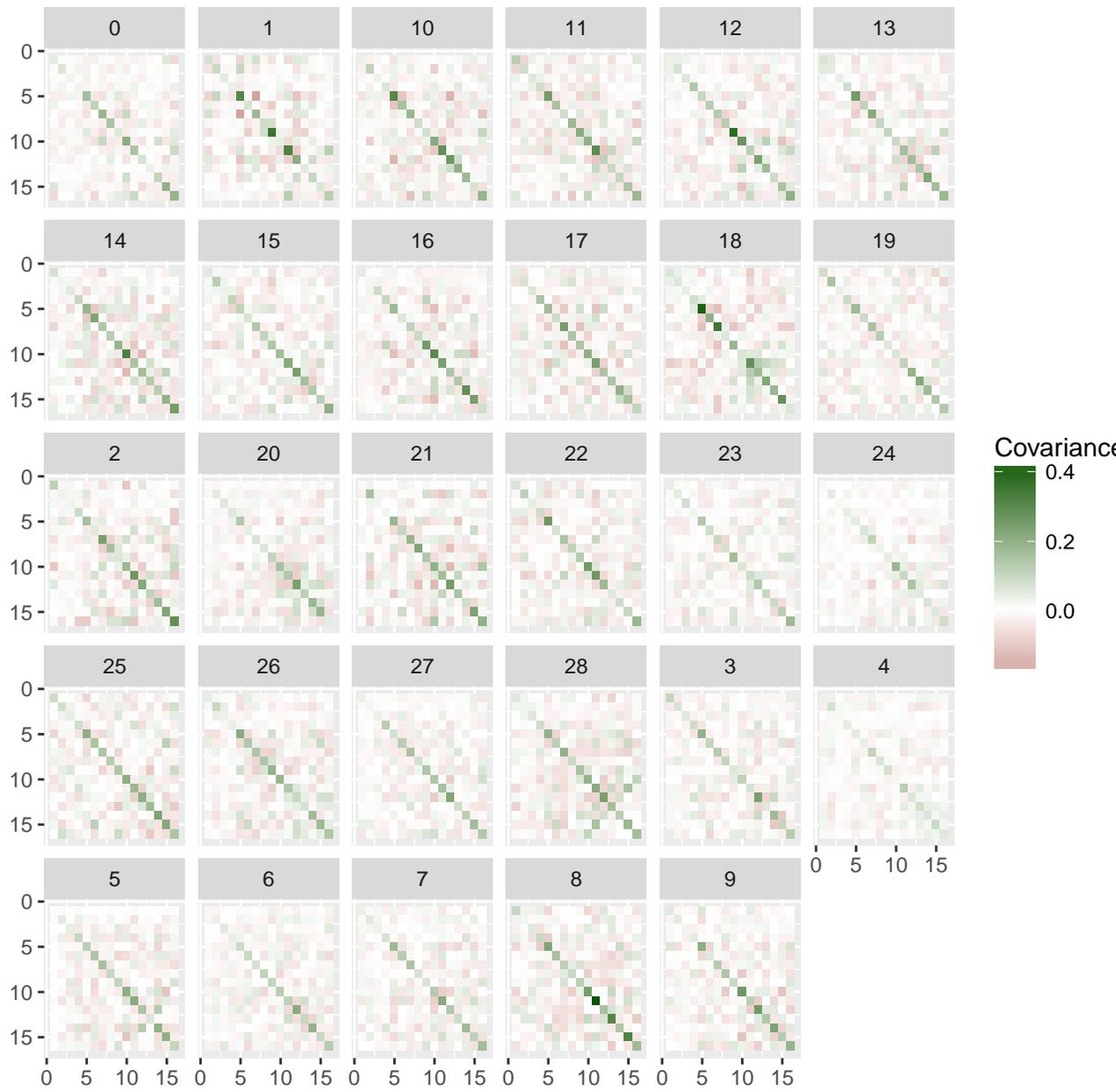
5 Appendix E : Additional null hypothesis RSA results

First, we show RSA matrices under the null hypothesis for all subjects and methods:

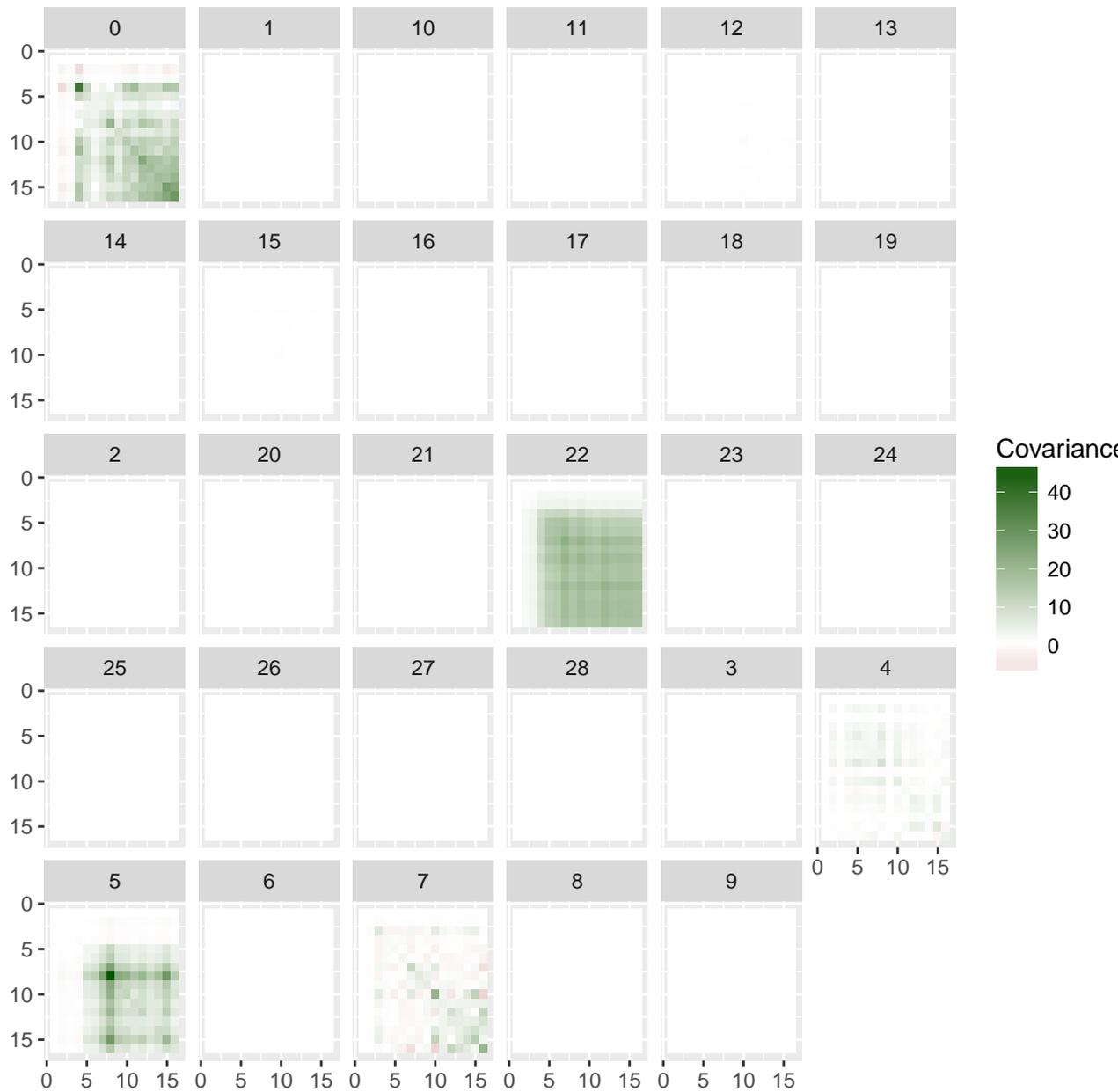
Naive RSA average matrix, null



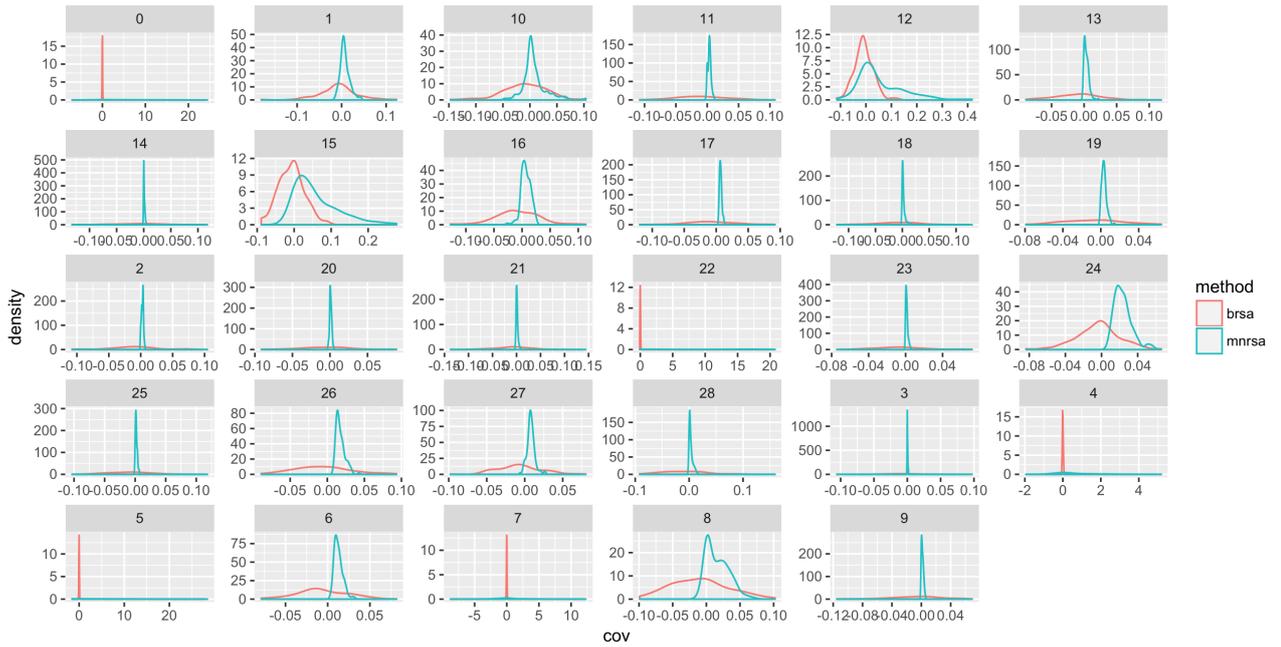
BRSA average matrix, null



MN-RSA average matrix, null

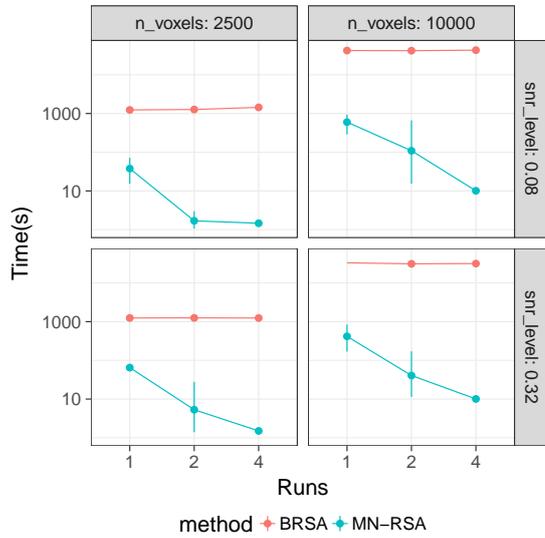


Notice that only for MN-RSA most of the covariances are noticeably degenerate. This is not a scaling effect on the figure driven by the color bar, but an effect on the underlying data, as we can see in the distribution of values in the covariance matrix for BRSA and MN-RSA:



6 Appendix F : timing figures for BRSA and MN-RSA

Experiment details mentioned in main text. Note time on the log scale.



Bibliography

E. Simony, C. J. Honey, J. Chen, O. Lositsky, Y. Yeshurun, A. Wiesel, and U. Hasson.
Dynamic reconfiguration of the default mode network during narrative comprehension.

Nature Communications, 7(May 2015):12141, jul 2016.

- B. M. Turner, B. U. Forstmann, E.-J. Wagenmakers, S. D. Brown, P. B. Sederberg, and M. Steyvers. A Bayesian framework for simultaneously modeling neural and behavioral data. *NeuroImage*, 72:193–206, may 2013.
- B. M. Turner, P. B. Sederberg, and J. L. McClelland. Bayesian analysis of simulation-based models. *Journal of Mathematical Psychology*, 2014.
- B. M. Turner, L. van Maanen, and B. U. Forstmann. Informing cognitive abstractions through neuroimaging: The neural drift diffusion model. *Psychological Review*, 122(2):312–336, 2015.
- B. M. Turner, C. A. Rodriguez, T. M. Norcia, S. M. McClure, and M. Steyvers. Why more is better: Simultaneous modeling of EEG, fMRI, and behavioral data. *NeuroImage*, 128:96–115, mar 2016.