
Supplementary Material: Probabilistic Linear Solvers for Machine Learning

Jonathan Wenger Philipp Hennig

University of Tübingen

Max Planck Institute for Intelligent Systems

Tübingen, Germany

{jonathan.wenger, philipp.hennig}@uni-tuebingen.de

This supplement complements the paper *Probabilistic Linear Solvers for Machine Learning* and is structured as follows. Section S1 explains the approach of probabilistic numerics to model (deterministic) numerical problems probabilistically in more depth. Section S2 introduces different variants of Kronecker products used to define matrix-variate normal distributions in Section S3. Section S4 details the matrix-based inference procedure of probabilistic linear solvers based on matrix-vector product observations. It also contains some more explanation regarding prior construction and stopping criteria. Section S5 and Section S6 outline theoretical results from the paper and properties of the proposed covariance class, in particular detailed proofs. Finally, Section S7 provides some background for the application of probabilistic linear solvers to the solution of discretized partial differential equations. To provide a clear exposition to the reader in some sections we restate results from the literature. References referring to sections, equations or theorem-type environments within this document are tagged with ‘S’, while references to, or results from the main paper are stated as is.

Preliminaries and Notation We consider the linear system $\mathbf{A}\mathbf{x}_* = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}_{\text{sym}}^{n \times n}$ is symmetric positive definite. The random variables \mathbf{A} , \mathbf{H} and \mathbf{x} model the linear operator \mathbf{A} , its inverse $\mathbf{H} = \mathbf{A}^{-1}$ and the solution \mathbf{x}_* . Algorithm 1 chooses actions $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_k] \in \mathbb{R}^{n \times k}$ given by its policy $\pi(\mathbf{s} \mid \mathbf{A}, \mathbf{H}, \mathbf{x}, \mathbf{A}, \mathbf{b})$ and computes observations $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_k] \in \mathbb{R}^{n \times k}$ given by a linear projection $\mathbf{y}_i = \mathbf{A}\mathbf{s}_i$ in each iteration $0 < i \leq k$.

S1 Probabilistic Modelling of Deterministic Problems

At first glance it might seem counterintuitive to frame a numerical problem in the language of probability theory. After all, when considering the exact problem $\mathbf{A}\mathbf{x}_* = \mathbf{b}$ all quantities involved \mathbf{A} , \mathbf{x}_* , and \mathbf{b} are deterministic. However, the distribution of the random variables \mathbf{A} , \mathbf{H} and \mathbf{x} represents *epistemic uncertainty* arising from finite computational resources. With a finite budget only a limited amount of information can be obtained about \mathbf{A} (e.g. via matrix-vector products). In particular, for a sufficiently large problem a priori the inverse $\mathbf{H} = \mathbf{A}^{-1}$ and the solution \mathbf{x}_* , while deterministic and computable in finite time, are not known. This uncertainty about the inverse is captured by the prior distribution of \mathbf{H} . In the Bayesian framework the belief about the inverse \mathbf{H} is then iteratively updated given new observations $\mathbf{y}_i = \mathbf{A}\mathbf{s}_i$.

The motivation for also estimating \mathbf{A} becomes clear if one considers the following. Usually in large-scale applications, the matrix \mathbf{A} is never actually formed in memory due to computational constraints. Instead only the matrix-vector product $\mathbf{v} \mapsto \mathbf{A}\mathbf{v}$ is available. Therefore without further computation, the value of any given matrix entry \mathbf{A}_{ij} is in fact uncertain. Further, generally other properties of the matrix \mathbf{A} such as its eigenspectrum are also not readily available. The probabilistic framework provides a principled way of incorporating prior knowledge about \mathbf{A} and makes assumptions about the problem explicit. Relating the prior model \mathbf{A} and \mathbf{H} is important here to allow Algorithm 1 to take such prior information into account in its policy. Finally, the strongest argument for a model \mathbf{A} may yet be the incorporation of noise. Suppose we only have access to $\mathbf{y}_i = (\mathbf{A} + \mathbf{E}_i)\mathbf{s}_i$ with

additive noise E_i . This is a common occurrence in application, where the linear system to be solved arises from an approximation itself or if A is constructed from data. Concrete examples are batched empirical risk minimization problems or stochastic quadratic optimization. In this setting the probabilistic linear solver must estimate the true A via its observations.

The application of probabilistic inference to numerical problems goes back well into the last century [1–3] and has recently seen a resurgence in research interest in the form of *probabilistic numerics*. Overviews discussing motivations and historical perspectives can be found in Hennig et al. [4] and Oates and Sullivan [5]. Hennig [6] gives additional insight into the statistical interpretation of linear systems.

S2 The Kronecker Product and its Variants

We will now introduce different types of Kronecker products needed for constructing covariances for matrix-variate distributions. In order to transfer results from probabilistic modelling of vector-variate random variables to the matrix-variate case, we need two types of vectorization operations, i.e. bijections between spaces of matrices and vector spaces.

Let $\text{vec} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{mn}$, denote the *column-wise stacking operator* [7], defined as

$$\text{vec}(\mathbf{X}) = (X_{11}, X_{21}, \dots, X_{m1}, X_{12}, \dots, X_{mn})^\top \in \mathbb{R}^{mn}.$$

Further, define $\text{svec} : \mathbb{R}_{\text{sym}}^{n \times n} \rightarrow \mathbb{R}^{\frac{1}{2}n(n+1)}$, the *column-wise symmetric stacking operator* [8] given by

$$\text{svec}(\mathbf{X}) = (X_{11}, \sqrt{2}X_{21}, \dots, \sqrt{2}X_{n1}, X_{22}, \sqrt{2}X_{32}, \dots, \sqrt{2}X_{n2}, \dots, X_{nn})^\top \in \mathbb{R}^{\frac{1}{2}n(n+1)}.$$

To translate between the two representations following Schacke [9] we also define the matrix $\mathbf{Q} \in \mathbb{R}^{\frac{1}{2}n(n+1) \times n^2}$ such that for all symmetric matrices $\mathbf{X} \in \mathbb{R}_{\text{sym}}^{n \times n}$, we have $\mathbf{Q} \text{vec}(\mathbf{X}) = \text{svec}(\mathbf{X})$ and $\text{vec}(\mathbf{X}) = \mathbf{Q}^\top \text{svec}(\mathbf{X})$. Note, that \mathbf{Q} has orthonormal rows, i.e. $\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$. For convenience we also name the inverse operations $\text{mat} := \text{vec}^{-1}$ and $\text{smat} := \text{svec}^{-1}$.

S2.1 Kronecker Product

We make extensive use of Kronecker-type structures for covariance matrices of matrix-variate distributions in this paper. The *Kronecker product* $\mathbf{A} \otimes \mathbf{B}$ [10] of two matrices $\mathbf{A} \in \mathbb{R}^{m_1 \times n_1}$ and $\mathbf{B} \in \mathbb{R}^{m_2 \times n_2}$ is given by

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} \mathbf{A}_{11}\mathbf{B} & \dots & \mathbf{A}_{1n_1}\mathbf{B} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{m_11}\mathbf{B} & \dots & \mathbf{A}_{m_1n_1}\mathbf{B} \end{pmatrix} \in \mathbb{R}^{(m_1m_2) \times (n_1n_2)}$$

The Kronecker product satisfies the characteristic property

$$(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{X}) = \text{vec}(\mathbf{B}\mathbf{X}\mathbf{A}^\top), \quad (\text{S1})$$

for $\mathbf{X} \in \mathbb{R}^{n_2 \times n_1}$. Characteristic properties of Kronecker-type products are useful to turn matrix equations into vector equations. We state a set of properties of the Kronecker product next without proof. More detail on Kronecker products can be found in Van Loan [10].

Proposition S1 (Properties of the Kronecker Product [10])

The Kronecker product satisfies the following identities:

$$\exists \mathbf{A}, \mathbf{B} : \mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A} \quad (\text{S2})$$

$$(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top \quad (\text{S3})$$

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} \quad (\text{S4})$$

$$(\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C} \quad (\text{S5})$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D}) \quad (\text{S6})$$

$$\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B}) \quad (\text{S7})$$

$$\mathbf{A} \in \mathbb{R}_{\text{sym}}^{m \times m}, \mathbf{B} \in \mathbb{R}_{\text{sym}}^{n \times n} \implies \mathbf{A} \otimes \mathbf{B} \in \mathbb{R}_{\text{sym}}^{mn \times mn} \quad (\text{S8})$$

$$\mathbf{A} \otimes \mathbf{B} = (\mathbf{L}_A \mathbf{L}_A^\top) \otimes (\mathbf{L}_B \mathbf{L}_B^\top) = (\mathbf{L}_A \otimes \mathbf{L}_B)(\mathbf{L}_A^\top \otimes \mathbf{L}_B^\top) \quad (\text{S9})$$

$$\mathbf{A} \otimes \mathbf{B} = (\mathbf{U}_A \mathbf{\Lambda}_A \mathbf{U}_A^\top) \otimes (\mathbf{U}_B \mathbf{\Lambda}_B \mathbf{U}_B^\top) = (\mathbf{U}_A \otimes \mathbf{U}_B)(\mathbf{\Lambda}_A \otimes \mathbf{\Lambda}_B)(\mathbf{U}_A^\top \otimes \mathbf{U}_B^\top) \quad (\text{S10})$$

S2.2 Box Product

The *box product* $A \boxtimes B \in \mathbb{R}^{(m_1 m_2) \times (n_1 n_2)}$ can be defined via its characteristic property

$$(A \boxtimes B) \text{vec}(Y) = \text{vec}(BY^T A^T) \quad (\text{S11})$$

for $Y \in \mathbb{R}^{n_1 \times n_2}$. See also Olsen et al. [11] for details.

Proposition S2 (Properties of the Box Product [11])

The box product satisfies the following identities:

$$\exists A, B : A \boxtimes B \neq B \boxtimes A \quad (\text{S12})$$

$$(A \boxtimes B)^T = B^T \boxtimes A^T \quad (\text{S13})$$

$$(A \boxtimes B)^{-1} = B^{-1} \boxtimes A^{-1} \quad (\text{S14})$$

$$(A + B) \boxtimes C = A \boxtimes C + B \boxtimes C \quad (\text{S15})$$

$$(A \boxtimes B)(C \boxtimes D) = (AD) \otimes (BC) \quad (\text{S16})$$

$$(A \boxtimes B)(C \otimes D) = (AD) \boxtimes (BC) \quad (\text{S17})$$

$$(A \otimes B)(C \boxtimes D) = (AC) \boxtimes (BD) \quad (\text{S18})$$

$$\text{tr}(A \boxtimes B) = \text{tr}(AB) \quad (\text{S19})$$

S2.3 Symmetric Kronecker Product

The *symmetric Kronecker product* $A \otimes B$ of two square matrices $A, B \in \mathbb{R}^{n \times n}$ is defined via its characteristic property for $X \in \mathbb{R}_{\text{sym}}^{n \times n}$ as

$$(A \otimes B) \text{svec}(X) = \frac{1}{2} \text{svec}(BXA^T + AXB^T) \quad (\text{S20})$$

or equivalently

$$A \otimes B = \frac{1}{2} Q(A \otimes B + B \otimes A) Q^T.$$

Proposition S3 (Properties of the Symmetric Kronecker Product [8, 9])

The symmetric Kronecker product satisfies the following identities:

$$A \otimes B = B \otimes A \quad (\text{S21})$$

$$(A \otimes B)^T = A^T \otimes B^T \quad (\text{S22})$$

$$(A \otimes A)^{-1} = A^{-1} \otimes A^{-1} \quad (\text{S23})$$

$$(A + B) \otimes C = A \otimes C + B \otimes C \quad (\text{S24})$$

$$(A \otimes B)(C \otimes D) = \frac{1}{2}(AC \otimes BD + AD \otimes BC) \quad (\text{S25})$$

$$A \in \mathbb{R}_{\text{sym}}^{n \times n}, B \in \mathbb{R}_{\text{sym}}^{n \times n} \implies A \otimes B \in \mathbb{R}_{\text{sym}}^{\frac{1}{2}n(n+1) \times \frac{1}{2}n(n+1)} \quad (\text{S26})$$

$$A \otimes A = (L_A L_A^T) \otimes (L_A L_A^T) = (L_A \otimes L_A)(L_A^T \otimes L_A^T) \quad (\text{S27})$$

$$A \otimes A = (U_A \Lambda_A U_A^T) \otimes (U_A \Lambda_A U_A^T) = (U_A \otimes U_A)(\Lambda_A \otimes \Lambda_A)(U_A^T \otimes U_A^T) \quad (\text{S28})$$

Note, that the symmetric Kronecker product represented as a $\frac{1}{2}n(n+1) \times \frac{1}{2}n(n+1)$ matrix is in general not symmetric.

Further properties can be found in Alizadeh et al. [8] and Schacke [9]. We prove the following technical results for mixed expressions of Kronecker-type products, which we will make use of later.

Corollary S1 (Mixed Kronecker Product Identities)

Let $A \in \mathbb{R}_{\text{sym}}^{n \times n}$, $B, C \in \mathbb{R}^{n \times k}$ and $X \in \mathbb{R}^{k \times k}$ such that $(CXB^T)^T = CXB^T$, then it holds that

$$Q^T(A \otimes A)Q(B \otimes C) \text{vec}(X) = \frac{1}{2}(AB \otimes AC + AC \boxtimes AB) \text{vec}(X) \quad (\text{S29})$$

$$(B^T \otimes C^T)Q^T(A \otimes A)Q = \frac{1}{2}(B^T A \otimes C^T A + B^T A \boxtimes C^T A). \quad (\text{S30})$$

$$(B^T \otimes C^T)Q^T(A \otimes A)Q(B \otimes C) \text{vec}(X) = \frac{1}{2}(B^T AB \otimes C^T AC + B^T AC \boxtimes C^T AB) \text{vec}(X). \quad (\text{S31})$$

Now, assume A to be invertible, $\text{rank}(C) = k$ and $Y \in \mathbb{R}^{k \times n}$ such that $(YC)^\top = YC$, then for

$$\begin{aligned} G &= (I_n \otimes C^\top)Q^\top(A \otimes A)Q(I_n \otimes C) \\ G_{\text{right}}^{-1} &= (2A^{-1} - C(C^\top AC)^{-1}C^\top) \otimes (C^\top AC)^{-1} \end{aligned}$$

we have $GG_{\text{right}}^{-1} \text{vec}(Y) = \text{vec}(Y)$, i.e. G_{right}^{-1} is the right inverse of G . Finally, for $D, E \in \mathbb{R}^{n \times n}$ and $Z \in \mathbb{R}_{\text{sym}}^{n \times n}$ such that $(EAZAD^\top)^\top = EAZAD^\top$, we have

$$(A^\top \otimes A^\top)Q(D \otimes E)Q^\top(A \otimes A) \text{svec}(Z) = (A^\top DA) \otimes (A^\top EA) \text{svec}(Z). \quad (\text{S32})$$

Proof. Let $X \in \mathbb{R}^{k \times k}$ such that $(CXB^\top)^\top = CXB^\top$, then

$$\begin{aligned} Q^\top(A \otimes A)Q(B \otimes C) \text{vec}(X) &= Q^\top(A \otimes A)Q \text{vec}(CXB^\top) \\ &= Q^\top(A \otimes A) \text{svec}(CXB^\top) \\ &= Q^\top \text{svec}(ACXB^\top A) \\ &= \frac{1}{2} \text{vec}(ACXB^\top A + ABX^\top C^\top A) \\ &= \frac{1}{2}(AB \otimes AC + AC \boxtimes AB), \end{aligned}$$

further it holds for $W \in \mathbb{R}_{\text{sym}}^{n \times n}$

$$\begin{aligned} (B^\top \otimes C^\top)Q^\top(A \otimes A)Q \text{vec}(W) &= (B^\top \otimes C^\top)Q^\top \text{svec}(AWA) \\ &= \text{vec}(C^\top AWAB) \\ &= \frac{1}{2}(C^\top AWAB + C^\top A^\top W^\top A^\top B) \\ &= \frac{1}{2}(B^\top A \otimes C^\top A + B^\top A \boxtimes C^\top A), \end{aligned}$$

and using the properties of the Kronecker and the Box product we obtain

$$\begin{aligned} (B^\top \otimes C^\top)Q^\top(A \otimes A)Q(B \otimes C) \text{vec}(X) &= (B^\top \otimes C^\top) \frac{1}{2}(B^\top A \otimes C^\top A + B^\top A \boxtimes C^\top A) \text{vec}(X) \\ &= \frac{1}{2}(B^\top A \otimes C^\top A + B^\top A \boxtimes C^\top A) \text{vec}(X). \end{aligned}$$

Now let A be invertible, let C have full rank and choose $Y \in \mathbb{R}^{k \times n}$ arbitrarily such that $(YC)^\top = YC$. Then using Proposition S1 and Proposition S2 we obtain

$$\begin{aligned} (I_n \otimes C^\top)Q^\top(A \otimes A)Q(I_n \otimes C)(2A^{-1} - C(C^\top AC)^{-1}C^\top) \otimes (C^\top AC)^{-1} \text{vec}(Y) \\ &= \frac{1}{2}(A \otimes C^\top AC + AC \boxtimes C^\top A)(2A^{-1} - C(C^\top AC)^{-1}C^\top) \otimes (C^\top AC)^{-1} \text{vec}(Y) \\ &= (I_n \otimes I_k - \frac{1}{2}AC(C^\top AC)^{-1}C^\top \otimes I_k + AC(C^\top AC)^{-1} \boxtimes C^\top - \frac{1}{2}AC(C^\top AC)^{-1} \boxtimes C^\top) \text{vec}(Y) \\ &= (I_n \otimes I_k - \frac{1}{2}AC(C^\top AC)^{-1}C^\top \otimes I_k + \frac{1}{2}AC(C^\top AC)^{-1} \boxtimes C^\top) \text{vec}(Y) \\ &= \text{vec}(Y) - \frac{1}{2}(YC(C^\top AC)^{-1}C^\top A - C^\top Y^\top (C^\top AC)^{-1}C^\top A) \\ &= \text{vec}(Y) \end{aligned}$$

Lastly, by assumption it holds that

$$\begin{aligned} (A^\top \otimes A^\top)Q(D \otimes E)Q^\top(A \otimes A) \text{svec}(Z) &= (A \otimes A)Q \text{vec}(EAZAD^\top) \\ &= \text{svec}(AEA ZAD^\top A) \\ &= \frac{1}{2}(AEA ZAD^\top A + ADA ZAE^\top A) \\ &= (ADA \otimes AEA) \text{svec}(Z). \end{aligned}$$

This concludes the proof. \square

S3 The Matrix-variate Normal Distribution

In order for our probabilistic linear solvers to infer the true latent \mathbf{A} or its inverse $\mathbf{H} = \mathbf{A}^{-1}$, we need a distribution expressing the belief of the solver over those latent quantities at any given point. A Gaussian distribution over matrices will play this role, motivated by the linear nature of the observations. This section closely follows Gupta and Nagar [12].

Definition S1 (Matrix-variate Normal Distribution [12])

Let $\mathbf{X}_0 \in \mathbb{R}^{m \times n}$ and let $\mathbf{V} \in \mathbb{R}_{\text{sym}}^m$ and $\mathbf{W} \in \mathbb{R}_{\text{sym}}^{n \times n}$ be positive-definite. We say a random matrix \mathbf{X} has a *matrix-variate normal distribution* with mean \mathbf{X}_0 and covariance $\mathbf{V} \otimes \mathbf{W}$, iff

$$\text{vec}(\mathbf{X}^\top) \sim \mathcal{N}_{mn}(\text{vec}(\mathbf{X}_0^\top), \mathbf{V} \otimes \mathbf{W}).$$

We write as a shorthand $\mathbf{X} \sim \mathcal{N}(\mathbf{X}_0, \mathbf{V} \otimes \mathbf{W})$.

Note, that the matrices \mathbf{V} and \mathbf{W} represent the covariance between rows and columns of \mathbf{X} , respectively. Since we model symmetric matrices in this work, we also introduce a Gaussian distribution over $\mathbb{R}_{\text{sym}}^{n \times n}$.

Definition S2 (Symmetric Matrix-variate Normal Distribution [12])

Let $\mathbf{X}_0, \mathbf{W} \in \mathbb{R}_{\text{sym}}^{n \times n}$ such that \mathbf{W} is positive-definite, then the random matrix \mathbf{X} has a *symmetric matrix-variate normal distribution*, iff

$$\text{svec}(\mathbf{X}) \sim \mathcal{N}_{\frac{1}{2}n(n+1)}(\text{svec}(\mathbf{X}_0), \mathbf{W} \otimes \mathbf{W}).$$

We write $\mathbf{X} \sim \mathcal{N}(\mathbf{X}_0, \mathbf{W} \otimes \mathbf{W})$.

It follows immediately from the definition that realizations of a symmetric matrix-variate normal distribution are symmetric matrices. This distribution also emerges naturally by conditioning a matrix-variate normal distribution on the linear constraint $\mathbf{X} = \mathbf{X}^\top$.

S4 Probabilistic Linear Solvers

Probabilistic linear solvers (PLS) [6, 13, 14] infer posterior beliefs over the matrix \mathbf{A} , its inverse \mathbf{H} or the solution $x_* = \mathbf{H}\mathbf{b}$ of a linear system via linear observations $\mathbf{Y} = \mathbf{A}\mathbf{S}$. We consider matrix-based inference [14] in this work. Assuming a prior $p(\mathbf{A})$ or $p(\mathbf{H})$, actions \mathbf{S} and linear observations \mathbf{Y} such methods return posterior distributions $p(\mathbf{A} \mid \mathbf{S}, \mathbf{Y})$ or $p(\mathbf{H} \mid \mathbf{S}, \mathbf{Y})$.

S4.1 Matrix-based Inference

The generic matrix-based inference procedure of probabilistic linear solvers is a consequence of the matrix-variate version of the following standard result for Gaussian inference under linear observations.

Theorem S1 (Linear Gaussian Inference [15])

Let $\mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\boldsymbol{\Sigma} \in \mathbb{R}_{\text{sym}}^{n \times n}$ positive-definite, and assume we are given observations of the form

$$\mathbf{B}\mathbf{v} + \mathbf{b} = \mathbf{y} \in \mathbb{R}^m,$$

where $\mathbf{B} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Assuming a Gaussian likelihood

$$p(\mathbf{y} \mid \mathbf{B}, \mathbf{v}, \mathbf{b}) = \mathcal{N}(\mathbf{y}; \mathbf{B}\mathbf{v} + \mathbf{b}, \boldsymbol{\Lambda}),$$

for $\boldsymbol{\Lambda} \in \mathbb{R}_{\text{sym}}^m$ positive definite, results in the posterior distribution

$$p(\mathbf{v} \mid \mathbf{y}, \mathbf{B}, \mathbf{b}) = \mathcal{N}(\mathbf{v}; \boldsymbol{\mu} + \boldsymbol{\Sigma}\mathbf{B}^\top(\mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^\top + \boldsymbol{\Lambda})^{-1}(\mathbf{y} - \mathbf{B}\boldsymbol{\mu} - \mathbf{b}), \\ \boldsymbol{\Sigma} - \boldsymbol{\Sigma}\mathbf{B}^\top(\mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^\top + \boldsymbol{\Lambda})^{-1}\mathbf{B}\boldsymbol{\Sigma}).$$

Further, the marginal distribution of \mathbf{y} is given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{B}\boldsymbol{\mu} + \mathbf{b}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^\top + \boldsymbol{\Lambda}).$$

S4.1.1 Asymmetric Model

Corollary S2 (Asymmetric matrix-based Gaussian Inference [16, 6, 14])

Assume a prior $p(\mathbf{A}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_0, \mathbf{V}_0 \otimes \mathbf{W}_0)$ and exact observations of the form $\mathbf{Y} = \mathbf{A}\mathbf{S}$, corresponding to a Dirac likelihood $p(\mathbf{Y} | \mathbf{A}, \mathbf{S}) = \delta(\mathbf{Y} - \mathbf{A}\mathbf{S})$, then the posterior $p(\mathbf{A} | \mathbf{S}, \mathbf{Y}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_k, \Sigma_k)$ is given by

$$\begin{aligned}\mathbf{A}_k &= \mathbf{A}_0 + \Delta_0 \mathbf{U}^\top \\ \Sigma_k &= \mathbf{V}_0 \otimes \mathbf{W}_0 (\mathbf{I}_n - \mathbf{S}\mathbf{U}^\top)\end{aligned}$$

where $\Delta_0 = \mathbf{Y} - \mathbf{A}_0\mathbf{S}$ and $\mathbf{U} = \mathbf{W}_0\mathbf{S}(\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}$.

Proof. In vectorized form the likelihood is given by

$$p(\text{vec}(\mathbf{Y}^\top) | \text{vec}(\mathbf{A}^\top), \text{vec}(\mathbf{S}^\top)) = \delta(\text{vec}(\mathbf{Y}^\top) - \text{vec}(\mathbf{S}^\top\mathbf{A}^\top)) = \delta(\text{vec}(\mathbf{Y}^\top) - (\mathbf{I} \otimes \mathbf{S}^\top) \text{vec}(\mathbf{A}^\top))$$

Using the Definition S1 of the matrix-variate normal distribution, applying Theorem S1 and using property (S6) of the Kronecker product in Proposition S1 leads to

$$\begin{aligned}\text{vec}(\mathbf{A}_k^\top) &= \text{vec}(\mathbf{A}_0^\top) + (\mathbf{V}_0 \otimes \mathbf{W}_0)(\mathbf{I} \otimes \mathbf{S})((\mathbf{I} \otimes \mathbf{S}^\top)(\mathbf{V}_0 \otimes \mathbf{W}_0)(\mathbf{I} \otimes \mathbf{S}))^{-1}(\text{vec}(\mathbf{Y}^\top) - (\mathbf{I} \otimes \mathbf{S}^\top) \text{vec}(\mathbf{A}_0^\top)) \\ &= \text{vec}(\mathbf{A}_0^\top) + (\mathbf{V}_0 \otimes \mathbf{W}_0\mathbf{S})(\mathbf{V}_0 \otimes \mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1} \text{vec}(\Delta_0^\top) \\ &= \text{vec}(\mathbf{A}_0^\top) + (\mathbf{I}_n \otimes \mathbf{W}_0\mathbf{S}(\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}) \text{vec}(\Delta_0^\top) \\ &= \text{vec}(\mathbf{A}_0^\top + \mathbf{U}\Delta_0^\top)\end{aligned}$$

and further analogously, additionally using bilinearity of the Kronecker product, we obtain

$$\begin{aligned}\Sigma_k &= \mathbf{V}_0 \otimes \mathbf{W}_0 - (\mathbf{V}_0 \otimes \mathbf{W}_0)(\mathbf{I} \otimes \mathbf{S})((\mathbf{I} \otimes \mathbf{S}^\top)(\mathbf{V}_0 \otimes \mathbf{W}_0)(\mathbf{I} \otimes \mathbf{S}))^{-1}(\mathbf{I} \otimes \mathbf{S}^\top)(\mathbf{V}_0 \otimes \mathbf{W}_0) \\ &= \mathbf{V}_0 \otimes \mathbf{W}_0 - (\mathbf{V}_0 \otimes \mathbf{W}_0\mathbf{S})(\mathbf{V}_0 \otimes \mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}(\mathbf{V}_0 \otimes \mathbf{S}^\top\mathbf{W}_0) \\ &= \mathbf{V}_0 \otimes \mathbf{W}_0 - \mathbf{V}_0 \otimes (\mathbf{W}_0\mathbf{S}(\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}\mathbf{S}^\top\mathbf{W}_0) \\ &= \mathbf{V}_0 \otimes \mathbf{W}_0(\mathbf{I} - \mathbf{S}\mathbf{U}^\top).\end{aligned}$$

This concludes the proof. \square

S4.1.2 Symmetric Model

Corollary S3 (Symmetric Matrix-based Gaussian Inference [16, 6, 14])

Assume a symmetric prior $p(\mathbf{A}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_0, \mathbf{W}_0 \otimes \mathbf{W}_0)$ and exact observations of the form $\mathbf{Y} = \mathbf{A}\mathbf{S}$, corresponding to a Dirac likelihood $p(\mathbf{Y} | \mathbf{A}, \mathbf{S}) = \delta(\mathbf{Y} - \mathbf{A}\mathbf{S})$, then the posterior $p(\mathbf{A} | \mathbf{S}, \mathbf{Y}) = \mathcal{N}(\mathbf{A}; \mathbf{A}_k, \Sigma_k)$ is given by

$$\begin{aligned}\mathbf{A}_k &= \mathbf{A}_0 + \Delta_0 \mathbf{U}^\top + \mathbf{U}\Delta_0^\top - \mathbf{U}\mathbf{S}^\top\Delta_0\mathbf{U}^\top = \mathbf{A}_0 + \mathbf{U}\mathbf{V}^\top + \mathbf{V}\mathbf{U}^\top \\ \Sigma_k &= \mathbf{W}_0(\mathbf{I}_n - \mathbf{S}\mathbf{U}^\top) \otimes \mathbf{W}_0(\mathbf{I}_n - \mathbf{S}\mathbf{U}^\top)\end{aligned}$$

where $\Delta_0 = \mathbf{Y} - \mathbf{A}_0\mathbf{S}$, $\mathbf{U} = \mathbf{W}_0\mathbf{S}(\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}$ and $\mathbf{V} = (\mathbf{I}_n - \frac{1}{2}\mathbf{U}\mathbf{S}^\top)\Delta_0$.

Proof. A proof can be found in the appendix of Hennig [6]. We rederive it here in our notation. By assumption the likelihood takes the vectorized form

$$p(\text{vec}(\mathbf{Y}^\top) | \text{svec}(\mathbf{A}), \text{vec}(\mathbf{S}^\top)) = \delta(\text{vec}(\mathbf{Y}^\top) - \text{vec}(\mathbf{S}^\top\mathbf{A}^\top)) = \delta(\text{vec}(\mathbf{Y}^\top) - (\mathbf{I} \otimes \mathbf{S}^\top)\mathbf{Q}^\top \text{svec}(\mathbf{A}))$$

Applying Theorem S1 gives

$$\begin{aligned}\text{svec}(\mathbf{A}_k) &= \text{svec}(\mathbf{A}_0) + (\mathbf{W}_0 \otimes \mathbf{W}_0)\mathbf{Q}(\mathbf{I}_n \otimes \mathbf{S})\mathbf{G}^{-1}(\text{vec}(\mathbf{Y}^\top) - (\mathbf{I} \otimes \mathbf{S}^\top)\mathbf{Q}^\top \text{svec}(\mathbf{A}_0)) \\ &= \text{svec}(\mathbf{A}_0) + (\mathbf{W}_0 \otimes \mathbf{W}_0)\mathbf{Q}(\mathbf{I}_n \otimes \mathbf{S})\mathbf{G}^{-1} \text{vec}(\Delta_0^\top) \\ \Sigma_k &= \mathbf{W}_0 \otimes \mathbf{W}_0 - (\mathbf{W}_0 \otimes \mathbf{W}_0)\mathbf{Q}(\mathbf{I}_n \otimes \mathbf{S})\mathbf{G}^{-1}(\mathbf{I}_n \otimes \mathbf{S}^\top)\mathbf{Q}^\top(\mathbf{W}_0 \otimes \mathbf{W}_0),\end{aligned}$$

where $\Delta_0 = \mathbf{Y} - \mathbf{A}_0\mathbf{S}$ and the Gram matrix is given by

$$\mathbf{G} = (\mathbf{I}_n \otimes \mathbf{S}^\top)\mathbf{Q}^\top(\mathbf{W}_0 \otimes \mathbf{W}_0)\mathbf{Q}(\mathbf{I}_n \otimes \mathbf{S}) \in \mathbb{R}^{nk \times nk}.$$

Now since $(\Delta_0^\top\mathbf{S})^\top = \Delta_0^\top\mathbf{S}$, we have by Corollary S1 that the right inverse of \mathbf{G} is given by

$$\mathbf{G}_{\text{right}}^{-1} = (2\mathbf{W}_0^{-1} - \mathbf{S}(\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}\mathbf{S}^\top) \otimes (\mathbf{S}^\top\mathbf{W}_0\mathbf{S})^{-1}$$

and therefore using (S6) and (S29) we obtain

$$\begin{aligned}
\text{svec}(\mathbf{A}_k) &= \text{svec}(\mathbf{A}_0) + (\mathbf{W}_0 \otimes \mathbf{W}_0) \mathbf{Q} (\mathbf{I}_n \otimes \mathbf{S}) \mathbf{G}_{\text{right}}^{-1} \text{vec}(\mathbf{\Delta}_0^\top) \\
&= \text{svec}(\mathbf{A}_0) + \mathbf{Q} \mathbf{Q}^\top (\mathbf{W}_0 \otimes \mathbf{W}_0) \mathbf{Q} (2\mathbf{W}_0^{-1} - \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \mathbf{S}^\top) \otimes \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \text{vec}(\mathbf{\Delta}_0^\top) \\
&= \text{svec}(\mathbf{A}_0) + \mathbf{Q} \frac{1}{2} ((2\mathbf{I} - \mathbf{U} \mathbf{S}^\top) \otimes \mathbf{U} + \mathbf{U} \boxtimes (2\mathbf{I} - \mathbf{U} \mathbf{S}^\top)) \text{vec}(\mathbf{\Delta}_0^\top) \\
&= \text{svec}(\mathbf{A}_0) + \text{svec}(\mathbf{U} \mathbf{\Delta}_0^\top (\mathbf{I} - \frac{1}{2} \mathbf{U} \mathbf{S}^\top)^\top + (\mathbf{I} - \frac{1}{2} \mathbf{U} \mathbf{S}^\top) \mathbf{\Delta}_0 \mathbf{U}^\top) \\
&= \text{svec}(\mathbf{A}_0 + \mathbf{\Delta}_0 \mathbf{U}^\top + \mathbf{U} \mathbf{\Delta}_0^\top - \mathbf{U} \mathbf{S}^\top \mathbf{\Delta}_0 \mathbf{U}^\top).
\end{aligned}$$

Further by definition it holds that

$$\mathbf{U} \mathbf{V}^\top + \mathbf{V} \mathbf{U}^\top = \mathbf{U} \mathbf{\Delta}_0^\top (\mathbf{I}_n - \frac{1}{2} \mathbf{S} \mathbf{U}^\top) + (\mathbf{I}_n - \frac{1}{2} \mathbf{U} \mathbf{S}^\top) \mathbf{\Delta}_0 \mathbf{U}^\top = \mathbf{\Delta}_0 \mathbf{U}^\top + \mathbf{U} \mathbf{\Delta}_0^\top - \mathbf{U} \mathbf{S}^\top \mathbf{\Delta}_0 \mathbf{U}^\top.$$

For the covariance we obtain using the right inverse of the Gram matrix and (S32) that

$$\begin{aligned}
\mathbf{\Sigma}_k &= \mathbf{W}_0 \otimes \mathbf{W}_0 - (\mathbf{W}_0 \otimes \mathbf{W}_0) \mathbf{Q} (\mathbf{I}_n \otimes \mathbf{S}) \mathbf{G}^{-1} (\mathbf{I}_n \otimes \mathbf{S}^\top) \mathbf{Q}^\top (\mathbf{W}_0 \otimes \mathbf{W}_0) \\
&= \mathbf{W}_0 \otimes \mathbf{W}_0 - (2\mathbf{W}_0 - \mathbf{W}_0 \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \mathbf{S}^\top \mathbf{W}_0) \otimes (\mathbf{W}_0 \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \mathbf{S}^\top \mathbf{W}_0) \\
&= (\mathbf{W}_0 - \mathbf{W}_0 \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \mathbf{S}^\top \mathbf{W}_0) \otimes (\mathbf{W}_0 - \mathbf{W}_0 \mathbf{S} (\mathbf{S}^\top \mathbf{W}_0 \mathbf{S})^{-1} \mathbf{S}^\top \mathbf{W}_0) \\
&= \mathbf{W}_0 (\mathbf{I}_n - \mathbf{S} \mathbf{U}^\top) \otimes \mathbf{W}_0 (\mathbf{I}_n - \mathbf{S} \mathbf{U}^\top).
\end{aligned}$$

□

S4.2 Matrix-variate Prior Construction

From a practical point of view it is important to be able to construct a prior for \mathbf{A} and \mathbf{H} from an initial guess \mathbf{x}_0 for the solution. This reduces down to finding \mathbf{A}_0 and \mathbf{H}_0 symmetric positive definite, such that $\mathbf{A}_0 = \mathbf{H}_0^{-1}$ and $\mathbf{x}_0 = \mathbf{H}_0 \mathbf{b}$ for the covariance class derived in Section 3. We provide a computationally efficient construction of such a prior here.

Proposition S4

Let $\mathbf{x}_0 \in \mathbb{R}^n$ and $\mathbf{b} \in \mathbb{R}^n \setminus \{0\}$. Assume $\mathbf{x}_0^\top \mathbf{b} > 0$, then for $\alpha < \frac{\mathbf{b}^\top \mathbf{x}_0}{\mathbf{b}^\top \mathbf{b}}$,

$$\mathbf{H}_0 = \alpha \mathbf{I} + \frac{1}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{b}} (\mathbf{x}_0 - \alpha \mathbf{b})(\mathbf{x}_0 - \alpha \mathbf{b})^\top$$

is symmetric positive definite and $\mathbf{H}_0 \mathbf{b} = \mathbf{x}_0$. Further it holds that

$$\mathbf{A}_0 = \mathbf{H}_0^{-1} = \alpha^{-1} \mathbf{I} - \frac{\alpha^{-1}}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{x}_0} (\mathbf{x}_0 - \alpha \mathbf{b})(\mathbf{x}_0 - \alpha \mathbf{b})^\top.$$

If $\mathbf{x}_0^\top \mathbf{b} < 0$ or $\mathbf{x}_0^\top \mathbf{b} = 0$, then for $\mathbf{x}_1 = -\mathbf{x}_0$ or $\mathbf{x}_1 = \frac{\mathbf{b}^\top \mathbf{x}_0}{\mathbf{b}^\top \mathbf{A} \mathbf{b}} \mathbf{b}$ respectively, it holds that $\|\mathbf{x}_1 - \mathbf{x}_*\|_{\mathbf{A}}^2 < \|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}^2$, i.e. \mathbf{x}_1 is a strictly better initialization than \mathbf{x}_0 .

Proof. Let \mathbf{H}_0 as above. Then $\mathbf{H}_0 \mathbf{b} = \alpha \mathbf{b} + \mathbf{x}_0 - \alpha \mathbf{b} = \mathbf{x}_0$. The second term of the sum in the form of \mathbf{H}_0 is of rank 1. Its non-zero eigenvalue is given by

$$\lambda = \frac{1}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{b}} (\mathbf{x}_0 - \alpha \mathbf{b})^\top (\mathbf{x}_0 - \alpha \mathbf{b}) = \frac{1}{\mathbf{x}_0^\top \mathbf{b} - \alpha \mathbf{b}^\top \mathbf{b}} \|\mathbf{x}_0 - \alpha \mathbf{b}\|_2^2 \geq 0$$

since by assumption $\mathbf{x}_0^\top \mathbf{b} > 0$ and $\alpha < \frac{\mathbf{b}^\top \mathbf{x}_0}{\mathbf{b}^\top \mathbf{b}}$. Now by Weyl's theorem it holds that $\lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{E}) \leq \lambda_{\min}(\mathbf{A} + \mathbf{E})$ and therefore \mathbf{H}_0 is positive definite. By the matrix inversion lemma we have for $\gamma = \frac{\alpha^{-1}}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{b}}$ that

$$\begin{aligned}
\mathbf{A}_0 &= \mathbf{H}_0^{-1} = \alpha^{-1} \left(\mathbf{I} - \frac{\gamma}{1 + \gamma \|\mathbf{x}_0 - \alpha \mathbf{b}\|_2^2} (\mathbf{x}_0 - \alpha \mathbf{b})(\mathbf{x}_0 - \alpha \mathbf{b})^\top \right) \\
&= \alpha^{-1} \mathbf{I} - \frac{\alpha^{-2}}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{b} + \alpha^{-1} \|\mathbf{x}_0 - \alpha \mathbf{b}\|_2^2} (\mathbf{x}_0 - \alpha \mathbf{b})(\mathbf{x}_0 - \alpha \mathbf{b})^\top \\
&= \alpha^{-1} \mathbf{I} - \frac{\alpha^{-1}}{(\mathbf{x}_0 - \alpha \mathbf{b})^\top \mathbf{x}_0} (\mathbf{x}_0 - \alpha \mathbf{b})(\mathbf{x}_0 - \alpha \mathbf{b})^\top.
\end{aligned}$$

Finally, we obtain

$$\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}^2 = (\mathbf{x}_0 - \mathbf{A}^{-1}\mathbf{b})^\top \mathbf{A} (\mathbf{x}_0 - \mathbf{A}^{-1}\mathbf{b}) = \mathbf{x}_0^\top \mathbf{A} \mathbf{x}_0 + \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b} - 2\mathbf{b}^\top \mathbf{x}_0.$$

Therefore if either $\mathbf{x}_0^\top \mathbf{b} < 0$ or $\mathbf{x}_0^\top \mathbf{b} = 0$, then $\mathbf{x}_1 = -\mathbf{x}_0$ or $\mathbf{x}_1 = \frac{\mathbf{b}^\top \mathbf{b}}{\mathbf{b}^\top \mathbf{A} \mathbf{b}} \mathbf{b}$, respectively are closer to \mathbf{x}_* in \mathbf{A} norm by positive definiteness of \mathbf{A} . This concludes the proof. \square

S4.3 Stopping Criteria

In addition to the classic stopping criteria $\|\mathbf{A}\mathbf{x}_k - \mathbf{b}\|_2 \leq \max(\delta_{\text{rtol}}\|\mathbf{b}\|_2, \delta_{\text{atol}})$ it is natural from a probabilistic viewpoint to use the induced posterior covariance of \mathbf{x} . Let $\mathbf{M} \in \mathbb{R}_{\text{sym}}^{n \times n}$ be a positive-definite matrix, then by linearity and the cyclic property of the trace it holds that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_*} [\|\mathbf{x}_* - \mathbb{E}[\mathbf{x}]\|_{\mathbf{M}}^2] &= \mathbb{E}_{\mathbf{x}_*} [(\mathbf{x}_* - \mathbb{E}[\mathbf{x}])^\top \mathbf{M} (\mathbf{x}_* - \mathbb{E}[\mathbf{x}])] \\ &= \text{tr}(\mathbb{E}_{\mathbf{x}_*} [(\mathbf{x}_* - \mathbb{E}[\mathbf{x}])^\top \mathbf{M} (\mathbf{x}_* - \mathbb{E}[\mathbf{x}])]) \\ &= \mathbb{E}_{\mathbf{x}_*} [\text{tr}((\mathbf{x}_* - \mathbb{E}[\mathbf{x}])^\top \mathbf{M} (\mathbf{x}_* - \mathbb{E}[\mathbf{x}]))] \\ &= \mathbb{E}_{\mathbf{x}_*} [\mathbf{M} \text{tr}((\mathbf{x}_* - \mathbb{E}[\mathbf{x}]) (\mathbf{x}_* - \mathbb{E}[\mathbf{x}])^\top)] \\ &= \text{tr}(\mathbf{M} \mathbb{E}_{\mathbf{x}_*} [(\mathbf{x}_* - \mathbb{E}[\mathbf{x}]) (\mathbf{x}_* - \mathbb{E}[\mathbf{x}])^\top]) \\ &= \text{tr}(\mathbf{M} (\text{Cov}[\mathbf{x}_* - \mathbb{E}[\mathbf{x}]] + (\mathbb{E}_{\mathbf{x}_*}[\mathbf{x}_*] - \mathbb{E}[\mathbf{x}])^\top (\mathbb{E}_{\mathbf{x}_*}[\mathbf{x}_*] - \mathbb{E}[\mathbf{x}]))) \\ &= \text{tr}(\mathbf{M} \text{Cov}[\mathbf{x}_*]) + \|\mathbb{E}_{\mathbf{x}_*}[\mathbf{x}_*] - \mathbb{E}[\mathbf{x}]\|_{\mathbf{M}}^2. \end{aligned}$$

Assuming calibration holds, i.e. $\mathbf{x}_* \sim \mathcal{N}(\mathbb{E}[\mathbf{x}], \text{Cov}[\mathbf{x}])$, we can bound the (relative) error by terminating when $\text{tr}(\mathbf{M} \text{Cov}[\mathbf{x}]) \leq \max(\delta_{\text{rtol}}\|\mathbf{b}\|, \delta_{\text{atol}})$ either in l_2 -norm for $\mathbf{M} = \mathbf{I}$ or in \mathbf{A} -norm for $\mathbf{M} = \mathbf{A}$.

We can efficiently evaluate the required $\text{tr}(\mathbf{M} \text{Cov}[\mathbf{x}])$ without ever forming $\text{Cov}[\mathbf{x}]$ in memory from already computed quantities. At iteration k we have $\text{Cov}[\mathbf{x}] = \text{Cov}[\mathbf{H}\mathbf{b}] = \frac{1}{2}(\mathbf{W}_k^\mathbf{H}(\mathbf{b}^\top \mathbf{W}_k^\mathbf{H} \mathbf{b}) + (\mathbf{W}_k^\mathbf{H} \mathbf{b})(\mathbf{W}_k^\mathbf{H} \mathbf{b})^\top)$ and therefore

$$\text{tr}(\mathbf{M} \text{Cov}[\mathbf{x}]) = \frac{1}{2}((\mathbf{b}^\top \mathbf{W}_k^\mathbf{H} \mathbf{b}) \text{tr}(\mathbf{M} \mathbf{W}_k^\mathbf{H}) + (\mathbf{W}_k^\mathbf{H} \mathbf{b})^\top \mathbf{M} (\mathbf{W}_k^\mathbf{H} \mathbf{b})).$$

Given the update for the covariance of the inverse view, we obtain the following recursion for its trace

$$\text{tr}(\mathbf{M} \mathbf{W}_k^\mathbf{H}) = \text{tr}(\mathbf{M} \mathbf{W}_{k-1}^\mathbf{H}) - \frac{1}{\mathbf{y}_k^\top \mathbf{W}_{k-1}^\mathbf{H} \mathbf{y}_k} \text{tr}((\mathbf{W}_{k-1}^\mathbf{H} \mathbf{y}_k)^\top \mathbf{M} (\mathbf{W}_{k-1}^\mathbf{H} \mathbf{y}_k)).$$

Computing the trace in this iterative fashion adds at most three matrix-vector products and three inner products for arbitrary \mathbf{M} all other quantities are computed for the covariance update anyhow.

For our proposed covariance class (3) we obtain for $\mathbf{M} = \mathbf{I}$ and $\Psi = \psi \mathbf{I}$ that

$$\begin{aligned} \text{tr}(\mathbf{W}_0^\mathbf{H}) &= \text{tr}(\mathbf{A}_0^{-1} \mathbf{Y} (\mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{A}_0^{-1} + (\mathbf{I} - \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top) \Psi (\mathbf{I} - \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top)) \\ &= \text{tr}((\mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{A}_0^{-1} \mathbf{Y}) + \psi \text{tr}((\mathbf{I} - \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top) (\mathbf{I} - \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top)) \\ &= \text{tr}((\mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{A}_0^{-1} \mathbf{Y}) + \psi \text{tr}(\mathbf{I} - \mathbf{Y} (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top) \\ &= \text{tr}((\mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{A}_0^{-1} \mathbf{A}_0^{-1} \mathbf{Y}) + \psi(n - k), \end{aligned}$$

which for a scalar prior mean $\mathbf{A}_0 = \alpha \mathbf{I}$ reduces to $\text{tr}(\mathbf{W}_0^\mathbf{H}) = \alpha^{-1}k + \psi(n - k)$.

S4.4 Implementation

In order to maintain numerical stability when performing low rank updates to symmetric positive definite matrices, as is the case in Algorithm 1 for the mean and covariance estimates, it is advantageous use a representation based on the Cholesky decomposition. One can perform the rank-2 update for the mean estimate and the rank-1 downdate for the covariance in Corollary S3 in each iteration of the algorithm for their respective Cholesky factors instead (see also Seeger [17]). The rank-2 update can be seen as a combination of a rank-1 up- and downdate by recognizing that

$$\mathbf{u}\mathbf{v}^\top + \mathbf{v}\mathbf{u}^\top = \frac{1}{2}((\mathbf{u} + \mathbf{v})(\mathbf{u} + \mathbf{v})^\top - (\mathbf{u} - \mathbf{v})(\mathbf{u} - \mathbf{v})^\top).$$

Similar updates arise in Quasi-Newton methods for the approximate (inverse) Hessian [18]. Having Cholesky factors of the mean and covariance available has the additional advantage that downstream sampling or the evaluation of the probability density function is computationally cheap.

S5 Theoretical Properties: Proofs for Section 2.3

In this section we provide detailed proofs for the theoretical results on convergence and the connection of Algorithm 1 to the method of conjugate gradients. We restate each theorem here as a reference to the reader. We begin by proving an intermediate result giving an interpretation to the posterior mean of \mathbf{A} and \mathbf{H} at each step of the method.

Proposition S5 (Subspace Equivalency)

Let \mathbf{A}_k and \mathbf{H}_k be the posterior means defined as in Section 2.1 and assume \mathbf{A}_0 and \mathbf{H}_0 are symmetric. Then for $1 \leq k \leq n$ it holds that

$$\mathbf{A}_k \mathbf{S} = \mathbf{Y} \quad \text{and} \quad \mathbf{H}_k \mathbf{Y} = \mathbf{S}, \quad (\text{S33})$$

i.e. \mathbf{A}_k and \mathbf{H}_k act like \mathbf{A} and \mathbf{A}^{-1} on the spaces spanned by the actions \mathbf{S} , respectively the observations \mathbf{Y} .

Proof. Since \mathbf{A}_0 and \mathbf{H}_0 are symmetric so are the expressions $\Delta_{\mathbf{A}} \mathbf{S}$ and $\Delta_{\mathbf{H}}^{\top} \mathbf{Y}$. We have that

$$\begin{aligned} \mathbf{A}_k \mathbf{S} &= (\mathbf{A}_0 + \Delta_{\mathbf{A}} \mathbf{U}_{\mathbf{A}}^{\top} + \mathbf{U}_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} - \mathbf{U}_{\mathbf{A}} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \mathbf{U}_{\mathbf{A}}^{\top}) \mathbf{S} \\ &= \mathbf{A}_0 \mathbf{S} + \Delta_{\mathbf{A}} \mathbf{I} + \mathbf{U}_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} - \mathbf{U}_{\mathbf{A}} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \mathbf{I} \\ &= \mathbf{A}_0 \mathbf{S} + \mathbf{Y} - \mathbf{A}_0 \mathbf{S} \\ &= \mathbf{Y}. \end{aligned}$$

In the case of the inverse model we obtain

$$\begin{aligned} \mathbf{H}_k \mathbf{Y} &= (\mathbf{H}_0 + \Delta_{\mathbf{H}} \mathbf{U}_{\mathbf{H}}^{\top} + \mathbf{U}_{\mathbf{H}} \Delta_{\mathbf{H}}^{\top} - \mathbf{U}_{\mathbf{H}} \mathbf{Y}^{\top} \Delta_{\mathbf{H}} \mathbf{U}_{\mathbf{H}}^{\top}) \mathbf{Y} \\ &= \mathbf{H}_0 \mathbf{Y} + \Delta_{\mathbf{H}} \mathbf{I} + \mathbf{U}_{\mathbf{H}} \Delta_{\mathbf{H}}^{\top} \mathbf{Y} - \mathbf{U}_{\mathbf{H}} \mathbf{Y}^{\top} \Delta_{\mathbf{H}} \mathbf{I} \\ &= \mathbf{H}_0 \mathbf{Y} + \mathbf{S} - \mathbf{H}_0 \mathbf{Y} \\ &= \mathbf{S} \end{aligned}$$

□

S5.1 Conjugate Directions Method

Theorem 1 (Conjugate Directions Method)

Given a prior $p(\mathbf{H}) = \mathcal{N}(\mathbf{H}; \mathbf{H}_0, \mathbf{W}_0^{\mathbf{H}} \otimes \mathbf{W}_0^{\mathbf{H}})$ such that $\mathbf{H}_0, \mathbf{W}_0^{\mathbf{H}} \in \mathbb{R}^{n \times n}$ positive definite, then actions \mathbf{s}_i of Algorithm 1 are \mathbf{A} -conjugate, i.e. for $0 \leq i, j \leq k$ with $i \neq j$ it holds that $\mathbf{s}_i^{\top} \mathbf{A} \mathbf{s}_j = 0$.

Proof. Since \mathbf{H}_0 is assumed to be symmetric, the form of the posterior mean in Section 2.1 implies that \mathbf{H}_k is symmetric for all $1 \leq k \leq n$. Now conjugacy is shown by induction. To that end, first consider the base case $k = 2$. We have

$$\begin{aligned} \mathbf{s}_2^{\top} \mathbf{A} \mathbf{s}_1 &= -\mathbf{r}_1^{\top} \mathbf{H}_1 \mathbf{A} \mathbf{s}_1 = -(\mathbf{r}_0^{\top} + \alpha_1 \mathbf{y}_1^{\top}) \mathbf{H}_1 \mathbf{A} \mathbf{s}_1 = -\left(\mathbf{r}_0^{\top} \mathbf{H}_1 - \frac{\mathbf{s}_1^{\top} \mathbf{r}_0}{\mathbf{s}_1^{\top} \mathbf{y}_1} \mathbf{y}_1^{\top} \mathbf{H}_1 \right) \mathbf{y}_1 \\ &= -\mathbf{r}_0^{\top} \mathbf{s}_1 + \mathbf{s}_1^{\top} \mathbf{r}_0 = 0 \end{aligned}$$

where we used (S33) and the definition of α_i in Algorithm 1. Now for the induction step, assume that $\mathbf{s}_i^{\top} \mathbf{A} \mathbf{s}_j = 0$ for all $i \neq j$ such that $1 \leq i, j \leq k$. We obtain for $1 \leq j \leq k$ that

$$\begin{aligned} \mathbf{s}_{k+1}^{\top} \mathbf{A} \mathbf{s}_j &= -\mathbf{r}_k^{\top} \mathbf{H}_k \mathbf{A} \mathbf{s}_j = -\left(\sum_{1 \leq l \leq k} \alpha_l \mathbf{y}_l + \mathbf{r}_0 \right)^{\top} \mathbf{H}_k \mathbf{y}_j = -\sum_{1 \leq l \leq k} \alpha_l \mathbf{y}_l^{\top} \mathbf{s}_j - \mathbf{r}_0^{\top} \mathbf{s}_j \\ &= -\alpha_j \mathbf{y}_j^{\top} \mathbf{s}_j - \mathbf{r}_0^{\top} \mathbf{s}_j = \mathbf{s}_j^{\top} \mathbf{r}_{j-1} - \mathbf{r}_0^{\top} \mathbf{s}_j = \mathbf{s}_j^{\top} \left(\sum_{1 \leq l < j} \alpha_l \mathbf{y}_l + \mathbf{r}_0 \right) - \mathbf{r}_0^{\top} \mathbf{s}_j \\ &= \mathbf{s}_j^{\top} \mathbf{r}_0 - \mathbf{r}_0^{\top} \mathbf{s}_j = 0 \end{aligned}$$

where we used the update equation of the residual \mathbf{r}_i in Algorithm 1, the definition of α_i , the induction hypothesis and (S33). This proves the statement. □

S5.2 Relationship to the Conjugate Gradient Method

Theorem 2 (Connection to the Conjugate Gradient Method)

Given a scalar prior mean $\mathbf{A}_0 = \mathbf{H}_0^{-1} = \alpha \mathbf{I}$ with $\alpha > 0$, assume (1) and (2) hold, then the iterates \mathbf{x}_i of Algorithm 1 are identical to the ones produced by the conjugate gradient method.

Proof. The proof outlined here is closely related to the proofs connecting Quasi-Newton methods to the conjugate gradient method [19, 6], but makes different assumptions on the prior distribution.

We begin by recognizing that the choice of step length α_i in Algorithm 1 is identical to the one in the conjugate gradient method [18]. Hence, it suffices to show that $\mathbf{s}_i \propto \mathbf{s}_i^{\text{CG}}$. Theorem 1 established that Algorithm 1 is a conjugate directions method. Now by assumption $\mathbf{A}_0 = \alpha \mathbf{I}$ and $\mathbf{H}_0 = \mathbf{A}_0^{-1}$, therefore $\mathbf{s}_1 = -\alpha \mathbf{I} \mathbf{r}_0 \propto -\mathbf{r}_0 = \mathbf{s}_1^{\text{CG}}$. It suffices show that \mathbf{s}_i lies in the Krylov space $\mathcal{K}_i(\mathbf{A}, \mathbf{r}_0) = \{\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{i-1}\mathbf{r}_0\}$ for all $0 < i \leq n$. This completes the argument, since $\mathcal{K}_i(\mathbf{A}, \mathbf{r}_0)$ is an i -dimensional subspace of \mathbb{R}^n and thus \mathbf{A} -conjugacy uniquely determines the search directions up to scaling, as \mathbf{A} is positive definite.

To complete the proof we proceed as follows. The posterior mean of the inverse model \mathbf{H}_{i-1} at step $i-1$ maps an arbitrary vector $\mathbf{v} \in \mathbb{R}^n$ to $\text{span}(\mathbf{H}_0 \mathbf{v}, \mathbf{H}_0 \mathbf{Y}_{1:i-1}, \mathbf{S}_{1:i-1}, \mathbf{W}_0^{\text{H}} \mathbf{Y}_{1:i-1})$. This follows directly from its form in given in Section 2.1. By assumption $\mathbf{H}_0 = \mathbf{A}_0^{-1} = \alpha^{-1} \mathbf{I}$, therefore using (1) and (2) we have $\text{span}(\mathbf{W}_0^{\text{H}} \mathbf{Y}_{1:i-1}) = \text{span}(\mathbf{Y}_{1:i-1})$. This implies \mathbf{H}_{i-1} maps to $\text{span}(\mathbf{v}, \mathbf{S}_{1:i-1}, \mathbf{Y}_{1:i-1})$ and thus $\mathbf{s}_i \in \text{span}(\mathbf{r}_{i-1}, \mathbf{S}_{1:i-1}, \mathbf{Y}_{1:i-1})$. We will now show that $\text{span}(\mathbf{r}_{i-1}, \mathbf{S}_{1:i-1}, \mathbf{Y}_{1:i-1}) \subset \mathcal{K}_i(\mathbf{A}, \mathbf{r}_0)$ by induction, completing the argument.

We begin with the base case. Since \mathbf{H}_0 is assumed to be scalar, we have $\mathbf{s}_1 \propto \mathbf{r}_0 \in \mathcal{K}_0(\mathbf{A}, \mathbf{r}_0)$ and therefore $\mathbf{y}_1 = \mathbf{A}\mathbf{s}_1$ and $\mathbf{r}_1 = \mathbf{r}_0 + \alpha_1 \mathbf{y}_1$ are in $\mathcal{K}_1(\mathbf{A}, \mathbf{r}_0)$. For the induction step assume $\text{span}(\mathbf{r}_{i-1}, \mathbf{S}_{1:i-1}, \mathbf{Y}_{1:i-1}) \subset \mathcal{K}_i(\mathbf{A}, \mathbf{r}_0)$. The definition of the policy of Algorithm 1 gives

$$\mathbf{s}_i = -\mathbb{E}[\mathbf{H}] \mathbf{r}_{i-1} \propto \mathbf{H}_{i-1} \mathbf{r}_{i-1} \in \text{span}(\mathbf{r}_{i-1}, \mathbf{S}_{1:i-1}, \mathbf{Y}_{1:i-1}) \subset \mathcal{K}_i(\mathbf{A}, \mathbf{r}_0),$$

where we used the induction hypothesis. This implies that $\mathbf{y}_i = \mathbf{A}\mathbf{s}_i \in \mathcal{K}_{i+1}(\mathbf{A}, \mathbf{r}_0)$ and $\mathbf{r}_i = \mathbf{r}_{i-1} + \alpha_i \mathbf{y}_i \in \mathcal{K}_{i+1}(\mathbf{A}, \mathbf{r}_0)$ by the definition of the Krylov space. Therefore, $\text{span}(\mathbf{r}_i, \mathbf{S}_{1:i}, \mathbf{Y}_{1:i}) \subset \mathcal{K}_{i+1}(\mathbf{A}, \mathbf{r}_0)$. This completes the proof. \square

S6 Prior Covariance Class: Proofs for Section 3

S6.1 Hereditary Positive-Definiteness

Proposition 1 (Hereditary Positive Definiteness [20, 16])

Let $\mathbf{A}_0 \in \mathbb{R}_{\text{sym}}^{n \times n}$ be positive definite. Assume the actions \mathbf{S} are \mathbf{A} -conjugate and $\mathbf{W}_0^{\text{A}} \mathbf{S} = \mathbf{Y}$, then for $i \in \{0, \dots, k-1\}$ it holds that \mathbf{A}_{i+1} is symmetric positive definite.

Proof. This is shown in Hennig and Kiefel [16]. We give an identical proof in our notation as a reference to the reader. By Theorem 7.5 in Dennis and Moré [20] it holds that if \mathbf{A}_i is positive definite and $\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1} \neq 0$, then \mathbf{A}_{i+1} is positive definite if and only if $\det(\mathbf{A}_{i+1}) > 0$. By the matrix determinant lemma and the recursive formulation of the posterior we have

$$\det(\mathbf{A}_{i+1}) = \det(\mathbf{A}_i) \left(\frac{1}{(\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})^2} \left((\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})^2 - (\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{y}_{i+1})(\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1}) + (\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})(\mathbf{y}_{i+1}^{\text{T}} \mathbf{s}_{i+1}) \right) \right)$$

Hence it suffices to show that

$$0 < (\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})^2 - (\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{y}_{i+1})(\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1}) + (\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})(\mathbf{y}_{i+1}^{\text{T}} \mathbf{s}_{i+1}),$$

which simplifies to

$$\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{y}_{i+1} - \frac{(\mathbf{y}_{i+1}^{\text{T}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1})^2}{\mathbf{s}_{i+1}^{\text{T}} \mathbf{W}_i^{\text{A}} \mathbf{A}_i^{-1} \mathbf{W}_i^{\text{A}} \mathbf{s}_{i+1}} < \mathbf{y}_{i+1}^{\text{T}} \mathbf{s}_{i+1}$$

Now by $\mathbf{W}_0^{\mathbf{A}}\mathbf{S} = \mathbf{Y}$, we have $\mathbf{W}_i^{\mathbf{A}}\mathbf{s}_{i+1} = \mathbf{W}_0^{\mathbf{A}}\mathbf{s}_{i+1} = \mathbf{y}_{i+1}$ and the above reduces to

$$0 < \mathbf{s}_{i+1}^{\top} \mathbf{A} \mathbf{s}_{i+1},$$

which is fulfilled by the assumption that \mathbf{A} is positive definite. Thus \mathbf{A}_{i+1} is positive definite. Symmetry follows immediately from the form of the posterior mean. \square

S6.2 Posterior Correspondence

Definition 1

Let \mathbf{A}_i and \mathbf{H}_i be the means of \mathbf{A} and \mathbf{H} at step i . We say a prior induces *posterior correspondence* if

$$\mathbf{A}_i^{-1} = \mathbf{H}_i \quad (\text{S34})$$

for all steps $0 \leq i \leq k$ of the solver. If only

$$\mathbf{A}_i^{-1} \mathbf{Y} = \mathbf{H}_i \mathbf{Y}, \quad (\text{S35})$$

we say that *weak posterior correspondence* holds.

S6.2.1 Matrix-variate Normal Prior

We begin by establishing posterior correspondence in the case of general matrix-variate normal priors, i.e. the inference setting detailed in Corollary S2. We begin by proving a general non-constructive condition and close with a sufficient condition for correspondence with limits the possible choices of covariance factors to a specific class.

Lemma S1 (General Correspondence)

Let $1 \leq k \leq n$, $\mathbf{W}_0^{\mathbf{A}}, \mathbf{W}_0^{\mathbf{H}}$ symmetric positive-definite and assume $\mathbf{A}_0^{-1} = \mathbf{H}_0$, then (S34) holds if and only if

$$0 = (\mathbf{A}\mathbf{S} - \mathbf{A}_0\mathbf{S}) [(\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \mathbf{A}\mathbf{S})^{-1} \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - (\mathbf{S}^{\top} \mathbf{A}^{\top} \mathbf{W}_0^{\mathbf{H}} \mathbf{A}\mathbf{S})^{-1} \mathbf{S}^{\top} \mathbf{A}^{\top} \mathbf{W}_0^{\mathbf{H}}]. \quad (\text{S36})$$

Proof. By the matrix inversion lemma we have

$$\begin{aligned} 0 &= \mathbf{A}_k^{-1} - \mathbf{H}_k \\ &= (\mathbf{A}_0 + (\mathbf{Y} - \mathbf{A}_0\mathbf{S})(\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{S})^{-1} \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}})^{-1} - \mathbf{H}_0 - (\mathbf{S} - \mathbf{H}_0\mathbf{Y})(\mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}} \mathbf{Y})^{-1} \mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}} \\ &= \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} (\mathbf{Y} - \mathbf{A}_0\mathbf{S})(\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{S} + \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} (\mathbf{Y} - \mathbf{A}_0\mathbf{S}))^{-1} \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \\ &\quad - \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} (\mathbf{A}_0\mathbf{S} - \mathbf{Y})(\mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}} \mathbf{Y})^{-1} \mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}} \\ &= -\mathbf{A}_0^{-1} (\mathbf{Y} - \mathbf{A}_0\mathbf{S}) [(\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \mathbf{Y})^{-1} \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - (\mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}} \mathbf{Y})^{-1} \mathbf{Y}^{\top} \mathbf{W}_0^{\mathbf{H}}], \end{aligned}$$

where we used the assumption $\mathbf{H}_0 = \mathbf{A}_0^{-1}$. Left-multiplying with $-\mathbf{A}_0$ and using $\mathbf{Y} = \mathbf{A}\mathbf{S}$ completes the proof. \square

Corollary S4 (Correspondence at Convergence)

Let $k = n$, $\mathbf{H}_0 = \mathbf{A}_0^{-1}$ and assume \mathbf{S} has full rank, i.e. the linear solver has performed n linearly independent actions, then (S34) holds for any symmetric positive-definite choice of $\mathbf{W}_0^{\mathbf{A}}$ and $\mathbf{W}_0^{\mathbf{H}}$.

Proof. By assumption, $\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1}$ and $\mathbf{S}^{\top} \mathbf{A}^{\top} \mathbf{W}_0^{\mathbf{H}}$ are invertible. Then by Lemma S1 the correspondence condition (S34) holds. \square

Theorem S2 (Sufficient Condition for Correspondence)

Let $1 \leq k \leq n$ arbitrary and assume $\mathbf{H}_0 = \mathbf{A}_0^{-1}$. Assume $\mathbf{W}_0^{\mathbf{A}}, \mathbf{A}_0, \mathbf{W}_0^{\mathbf{H}}$ satisfy

$$0 = \mathbf{S}^{\top} (\mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - \mathbf{A}^{\top} \mathbf{W}_0^{\mathbf{H}}) \quad (\text{S37})$$

or equivalently let $\mathbf{B}_{\langle \mathbf{S} \rangle^{\perp}} \in \mathbb{R}^{n \times k}$ be a basis of the orthogonal space $\langle \mathbf{S} \rangle^{\perp}$ spanned by the actions. For $\Phi \in \mathbb{R}^{(n-k) \times n}$ arbitrary, if

$$\mathbf{W}_0^{\mathbf{H}} = \mathbf{A}^{-\top} (\mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - \mathbf{B}_{\langle \mathbf{S} \rangle^{\perp}} \Phi) \quad (\text{S38})$$

and the commutation relations

$$[\mathbf{A}_0, \mathbf{A}] = \mathbf{0} \quad (\text{S39})$$

$$[\mathbf{W}_0^{\mathbf{A}}, \mathbf{A}] = \mathbf{0} \quad (\text{S40})$$

$$[\mathbf{B}_{\langle \mathcal{S} \rangle^\perp} \Phi, \mathbf{A}] = \mathbf{0} \quad (\text{S41})$$

are fulfilled, then $\mathbf{W}_0^{\mathbf{H}}$ is symmetric and (S34) holds.

Proof. By assumption $\mathbf{W}_0^{\mathbf{A}}$ is symmetric positive-definite and (S37) is equivalent to $\mathbf{S}^\top \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} = \mathbf{S}^\top \mathbf{A}^\top \mathbf{W}_0^{\mathbf{H}}$, which implies (S36). Now, assumption (S37) is equivalent to columns of the difference $\mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - \mathbf{A}^\top \mathbf{W}_0^{\mathbf{H}}$ lying in L , i.e. we can choose a basis $\mathbf{B}_{\langle \mathcal{S} \rangle^\perp}$ and coefficient matrix Φ such that

$$\mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - \mathbf{A}^\top \mathbf{W}_0^{\mathbf{H}} = \mathbf{B}_{\langle \mathcal{S} \rangle^\perp} \Phi.$$

Rearranging the above gives (S38). With the commutation relations and

$$[\mathbf{A}, \mathbf{B}] = \mathbf{0} \iff [\mathbf{A}^{-1}, \mathbf{B}] = \mathbf{0} \iff [\mathbf{A}, \mathbf{B}^{-1}] = \mathbf{0} \iff [\mathbf{A}^{-1}, \mathbf{B}^{-1}] = \mathbf{0}$$

it holds that

$$(\mathbf{W}_0^{\mathbf{H}})^\top = \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \mathbf{A}^{-1} - \mathbf{B}_{\langle \mathcal{S} \rangle^\perp} \Phi \mathbf{A}^{-1} = \mathbf{A}^{-\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} - \mathbf{A}^{-\top} \mathbf{B}_{\langle \mathcal{S} \rangle^\perp} \Phi = \mathbf{W}_0^{\mathbf{H}}$$

hence $\mathbf{W}_0^{\mathbf{H}}$ is symmetric. Finally, by Lemma S1 posterior mean correspondence (S34) holds. \square

If we want to ensure correspondence for all iterations, (S41) is trivially satisfied. The question now becomes what form can \mathbf{A}_0 and $\mathbf{W}_0^{\mathbf{A}}$ take in order to ensure symmetric $\mathbf{W}_0^{\mathbf{H}}$. This comes down to finding matrices which commute with \mathbf{A} .

Lemma S2 (Commuting Matrices of a Symmetric Matrix)

Let $r \in \mathbb{N}$, $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ symmetric. Assume \mathbf{M} has the form

$$\mathbf{M} = \mathfrak{p}_r(\mathbf{A}) = \sum_{i=0}^r c_i \mathbf{A}^i$$

for a set of coefficients $c_i \in \mathbb{R}$, then \mathbf{M} and \mathbf{A} commute. If \mathbf{A} has n distinct eigenvalues, \mathbf{M} is diagonalizable and $[\mathbf{M}, \mathbf{A}] = \mathbf{0}$, then

$$\mathbf{M} = \mathfrak{p}_{n-1}(\mathbf{A}),$$

i.e. \mathbf{M} is a polynomial in \mathbf{A} of degree at most $n - 1$.

Proof. The first result follows immediately since

$$\mathbf{W}_0^{\mathbf{A}} \mathbf{A} = \mathfrak{p}_r(\mathbf{A}) \mathbf{A} = \sum_{i=0}^r c_i \mathbf{A}^{i+1} = \mathbf{A} \mathfrak{p}_r(\mathbf{A}) = \mathbf{A} \mathbf{W}_0^{\mathbf{A}}.$$

Assume now that \mathbf{A} has n distinct eigenvalues $\lambda_0, \dots, \lambda_{n-1}$, \mathbf{M} is diagonalizable and \mathbf{M} and \mathbf{A} commute. Now, if and only if $[\mathbf{A}, \mathbf{M}] = \mathbf{0}$, then \mathbf{A} and \mathbf{M} are simultaneously diagonalizable by Theorem 5.2 in Conrad [21], i.e. we can find a common basis in which both \mathbf{A} and \mathbf{M} are represented by diagonal matrices. Hence, the set of matrices commuting with \mathbf{A} forms an n -dimensional subspace $\mathcal{U}_n \subset \mathbb{R}^{n \times n}$. Now, by the first part of this proof $\{\mathbf{I}, \mathbf{A}, \dots, \mathbf{A}^{n-1}\} \subset \mathcal{U}_n$. It remains to be shown, that this set forms a basis of \mathcal{U}_n . By isomorphism of finite dimensional vector spaces this is equivalent to proving that

$$\{\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_{n-1}\} := \left\{ \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} \lambda_0 \\ \vdots \\ \lambda_{n-1} \end{pmatrix}, \dots, \begin{pmatrix} \lambda_0^{n-1} \\ \vdots \\ \lambda_{n-1}^{n-1} \end{pmatrix} \right\}$$

forms a basis of \mathbb{R}^n . It suffices to show that all \mathbf{b}_i are independent. Assume the contrary, then $\sum_{i=0}^{n-1} \alpha_i \mathbf{b}_i = \mathbf{0}$ for some $\alpha_0, \dots, \alpha_{n-1} \in \mathbb{R}$, such that not all $\alpha_i = 0$. This implies that the polynomial $\sum_{i=0}^{n-1} \alpha_i x^i$ has n zeros $\lambda_0, \dots, \lambda_{n-1}$. This contradicts the fundamental theorem of algebra, concluding the proof. \square

The above suggests that tractable choices of \mathbf{A}_0 and $\mathbf{W}_0^{\mathbf{A}}$ for the non-symmetric matrix-variate prior, which imply symmetric $\mathbf{W}_0^{\mathbf{H}}$, are of polynomial form in \mathbf{A} .

Example S1 (Posterior Correspondence Covariance Class)

Tractable choices of the prior parameters in the \mathbf{A} view, which satisfy posterior correspondence and the commutation relations are for example

$$\mathbf{A}_0 = c_0 \mathbf{I} \quad \text{and} \quad \mathbf{W}_0^{\mathbf{A}} = \sum_{i=1}^{n-1} c_i \mathbf{A}^i,$$

where $\mathbf{H}_0 = \mathbf{A}_0^{-1}$ with $c_i \in \mathbb{R}$. Motivated by $\text{tr}(\mathbf{A}) \stackrel{!}{=} \text{tr}(\mathbf{A}_0)$ an initial choice could be $c_0 = n^{-1} \text{tr}(\mathbf{A})$.

Finally, note that in practice we do not actually require $\mathbf{W}_0^{\mathbf{A}}$. We only ever need access to $\mathbf{W}_0^{\mathbf{A}} \mathbf{S}$.

S6.2.2 Symmetric Matrix-variate Normal Prior

We now turn to the symmetric model, which we assumed throughout the paper, given in Corollary S3. We prove Theorem 3, the main result of this section demonstrating *weak posterior correspondence* for the symmetric Kronecker covariance, by employing the matrix inversion lemma for the posterior mean \mathbf{A}_k . We begin by establishing a set of technical lemmata first, which mainly expand terms appearing during matrix block inversion.

Lemma S3 (Symmetric Posterior Inverse)

Under the assumptions of Corollary S3, the inverse of the posterior mean is given by

$$\mathbf{A}_k^{-1} = \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} [\mathbf{U}_{\mathbf{A}} \quad \mathbf{V}_{\mathbf{A}}] \begin{bmatrix} \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{I} + \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \\ \mathbf{I} + \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{U}_{\mathbf{A}}^{\top} \\ \mathbf{V}_{\mathbf{A}}^{\top} \end{bmatrix} \mathbf{A}_0^{-1}$$

where

$$\begin{aligned} \mathbf{U}_{\mathbf{A}} &:= \mathbf{W}_0^{\mathbf{A}} \mathbf{S} (\mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{S})^{-1} \in \mathbb{R}^{n \times k}, \\ \mathbf{V}_{\mathbf{A}} &:= (\mathbf{I} - \frac{1}{2} \mathbf{U}_{\mathbf{A}} \mathbf{S}^{\top}) (\mathbf{Y} - \mathbf{A}_0 \mathbf{S}) = (\mathbf{I} - \frac{1}{2} \mathbf{U}_{\mathbf{A}} \mathbf{S}^{\top}) \Delta_{\mathbf{A}} \in \mathbb{R}^{n \times k}. \end{aligned}$$

Proof. We rewrite the rank-2 update in Section 2.1 as follows

$$\mathbf{A}_k = \mathbf{A}_0 + \mathbf{U}_{\mathbf{A}} \mathbf{V}_{\mathbf{A}}^{\top} + \mathbf{V}_{\mathbf{A}} \mathbf{U}_{\mathbf{A}}^{\top} = \mathbf{A}_0 + [\mathbf{U}_{\mathbf{A}} \quad \mathbf{V}_{\mathbf{A}}] \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{\mathbf{A}}^{\top} \\ \mathbf{V}_{\mathbf{A}}^{\top} \end{bmatrix}.$$

Then the statement follows directly from the matrix inversion lemma. \square

Next, we expand the terms inside the blocks of the matrix to be inverted in Lemma S3. This leads to the following lemma.

Lemma S4

Given the assumptions of Corollary S3, let $\mathbf{W}_0^{\mathbf{A}}$ and \mathbf{A}_0 be symmetric and assume (2) and (1) hold. Define

$$\begin{aligned} \mathbf{\Lambda} &= \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{S} \\ \mathbf{\Pi} &= \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \Delta_{\mathbf{A}}, \end{aligned}$$

then $\mathbf{\Lambda} \in \mathbb{R}^{m \times m}$ and $\mathbf{\Lambda} + \mathbf{\Pi} \in \mathbb{R}^{m \times m}$ are symmetric and invertible and we obtain

$$\mathbf{\Lambda} + \mathbf{\Pi} = \mathbf{S}^{\top} \mathbf{W}_0^{\mathbf{A}} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} = \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} = \mathbf{S}^{\top} \mathbf{A} \mathbf{W}_0^{\mathbf{H}} \mathbf{A} \mathbf{S} \quad (\text{S42})$$

$$\mathbf{\Pi} = \Delta_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \quad (\text{S43})$$

$$\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \Delta_{\mathbf{A}} = \mathbf{\Lambda}^{-1} \mathbf{\Pi} \quad (\text{S44})$$

$$\Delta_{\mathbf{A}}^{\top} \mathbf{S} = \mathbf{S}^{\top} \Delta_{\mathbf{A}} \quad (\text{S45})$$

$$\mathbf{U}_{\mathbf{A}} = \mathbf{A} \mathbf{S} \mathbf{\Lambda}^{-1} \quad (\text{S46})$$

$$\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} = \mathbf{\Lambda}^{-1} (\mathbf{\Lambda} + \mathbf{\Pi}) \mathbf{\Lambda}^{-1} \quad (\text{S47})$$

$$I + U_{\mathbf{A}}^{\top} A_0^{-1} V_{\mathbf{A}} = \Lambda^{-1}(\Lambda + \Pi)(I - \frac{1}{2}\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}}) \quad (\text{S48})$$

$$I + V_{\mathbf{A}}^{\top} A_0^{-1} U_{\mathbf{A}} = (I - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1})(\Lambda + \Pi)\Lambda^{-1} \quad (\text{S49})$$

$$V_{\mathbf{A}}^{\top} A_0^{-1} V_{\mathbf{A}} = \Pi - \frac{1}{2}((\Lambda + \Pi)\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} + \Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}(\Lambda + \Pi)) \quad (\text{S50})$$

$$+ \frac{1}{4}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}(\Lambda + \Pi)\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \quad (\text{S51})$$

Proof. We begin by proving that Λ and $\Lambda + \Pi$ are symmetric and invertible. We have by Sylvester's rank inequality that Λ is invertible. For symmetric $W_0^{\mathbf{A}}$, Λ is symmetric by definition. We have that

$$\begin{aligned} \Lambda + \Pi &= S^{\top}W_0^{\mathbf{A}}S + S^{\top}W_0^{\mathbf{A}}A_0^{-1}(AS - A_0S) = S^{\top}W_0^{\mathbf{A}}A_0^{-1}AS = S^{\top}AA_0^{-1}AS \\ &= S^{\top}W_0^{\mathbf{A}}A_0^{-1}AS = S^{\top}AW_0^{\mathbf{H}}AS \end{aligned}$$

Thus, by Sylvester's rank inequality $\Lambda + \Pi$ is invertible. Given symmetric A_0 , it is symmetric. Further, it holds that

$$\begin{aligned} \Pi &= \Lambda + \Pi - \Lambda = S^{\top}AA_0^{-1}AS - S^{\top}AS = \Delta_{\mathbf{A}}^{\top}A_0^{-1}AS \\ U_{\mathbf{A}}^{\top}A_0^{-1}\Delta_{\mathbf{A}} &= (S^{\top}W_0^{\mathbf{A}}S)^{-1}S^{\top}W_0^{\mathbf{A}}A_0^{-1}\Delta_{\mathbf{A}} = \Lambda^{-1}\Pi \\ \Delta_{\mathbf{A}}^{\top}S &= (AS - A_0S)^{\top}S = S^{\top}AS - S^{\top}A_0S \\ U_{\mathbf{A}} &= W_0^{\mathbf{A}}S(S^{\top}W_0^{\mathbf{A}}S)^{-1} = AS\Lambda^{-1} \\ U_{\mathbf{A}}^{\top}A_0^{-1}U_{\mathbf{A}} &= \Lambda^{-1}S^{\top}AA_0^{-1}AS\Lambda^{-1} = \Lambda^{-1}(\Lambda + \Pi)\Lambda^{-1} \\ I + U_{\mathbf{A}}^{\top}A_0^{-1}V_{\mathbf{A}} &= I + \Lambda^{-1}S^{\top}AA_0^{-1}(I - \frac{1}{2}U_{\mathbf{A}}S^{\top})\Delta_{\mathbf{A}} = I + \Lambda^{-1}S^{\top}AA_0^{-1}(I - \frac{1}{2}AS\Lambda^{-1}S^{\top})\Delta_{\mathbf{A}} \\ &= I + \Lambda^{-1}S^{\top}AA_0^{-1}(AS - A_0S) - \frac{1}{2}\Lambda^{-1}S^{\top}AA_0^{-1}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \\ &= \Lambda^{-1}(\Lambda + \Pi) - \frac{1}{2}\Lambda^{-1}(\Lambda + \Pi)\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} = \Lambda^{-1}(\Lambda + \Pi)(I - \frac{1}{2}\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}}) \\ I + V_{\mathbf{A}}^{\top}A_0^{-1}U_{\mathbf{A}} &= (I + U_{\mathbf{A}}^{\top}A_0^{-1}V_{\mathbf{A}})^{\top} = (\Lambda^{-1}(\Lambda + \Pi)(I - \frac{1}{2}\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}}))^{\top} \\ &= (I - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1})(\Lambda + \Pi)\Lambda^{-1}, \end{aligned}$$

where we used that Λ and $\Lambda + \Pi$ are symmetric. Finally, we have that

$$\begin{aligned} V_{\mathbf{A}}^{\top}A_0^{-1}V_{\mathbf{A}} &= \Delta_{\mathbf{A}}^{\top}(I - \frac{1}{2}SU_{\mathbf{A}}^{\top})A_0^{-1}(I - \frac{1}{2}U_{\mathbf{A}}S^{\top})\Delta_{\mathbf{A}} \\ &= \Delta_{\mathbf{A}}^{\top}(I - \frac{1}{2}S\Lambda^{-1}S^{\top}A)A_0^{-1}(I - \frac{1}{2}AS\Lambda^{-1}S^{\top})\Delta_{\mathbf{A}} \\ &= \Delta_{\mathbf{A}}^{\top}A_0^{-1}(I - \frac{1}{2}AS\Lambda^{-1}S^{\top})\Delta_{\mathbf{A}} - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}(I - \frac{1}{2}AS\Lambda^{-1}S^{\top})\Delta_{\mathbf{A}} \\ &= (S^{\top}AA_0^{-1} - S^{\top})(I - \frac{1}{2}AS\Lambda^{-1}S^{\top})\Delta_{\mathbf{A}} - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}\Delta_{\mathbf{A}} \\ &\quad + \frac{1}{4}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \\ &= S^{\top}AA_0^{-1}\Delta_{\mathbf{A}} - S^{\top}\Delta_{\mathbf{A}} - \frac{1}{2}S^{\top}AA_0^{-1}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} + \frac{1}{2}S^{\top}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \\ &\quad - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}\Delta_{\mathbf{A}} + \frac{1}{4}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \\ &= S^{\top}AA_0^{-1}AS - S^{\top}AS - S^{\top}\Delta_{\mathbf{A}} - \frac{1}{2}(\Lambda + \Pi)\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} + \frac{1}{2}S^{\top}\Delta_{\mathbf{A}} \\ &\quad - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}\Delta_{\mathbf{A}} + \frac{1}{4}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}AS\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} \\ &= \Pi - \frac{1}{2}S^{\top}\Delta_{\mathbf{A}} - \frac{1}{2}(\Lambda + \Pi)\Lambda^{-1}S^{\top}\Delta_{\mathbf{A}} - \frac{1}{2}\Delta_{\mathbf{A}}^{\top}S\Lambda^{-1}S^{\top}AA_0^{-1}(AS - A_0S) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{4} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A A_0^{-1} A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} \\
= & \Pi - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} - \frac{1}{2} (\Lambda + \Pi) \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} (\Lambda + \Pi) + \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} \Lambda \\
& + \frac{1}{4} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A A_0^{-1} A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} \\
= & \Pi - \frac{1}{2} ((\Lambda + \Pi) \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} + \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} (\Lambda + \Pi)) + \frac{1}{4} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} (\Lambda + \Pi) \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}},
\end{aligned}$$

where we dropped some of the terms temporarily for clarity of exposition. \square

We will now use these intermediate results to perform block inversion on the $2k \times 2k$ matrix to be inverted in Lemma S3.

Lemma S5

Given the assumptions of Corollary S3, additionally assume (1) and (2) hold. Let

$$\mathbf{T} = \begin{bmatrix} \mathbf{T}_{11} & \mathbf{T}_{12} \\ \mathbf{T}_{21} & \mathbf{T}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{I} + \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \\ \mathbf{I} + \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \end{bmatrix}^{-1},$$

then the block matrices $\mathbf{T}_{ij} \in \mathbb{R}^{m \times m}$ are given by

$$\begin{aligned}
\mathbf{T}_{11} &= \Lambda (\Lambda + \Pi)^{-1} \Lambda - (\mathbf{I} - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) (\mathbf{I} - \frac{1}{2} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) \\
\mathbf{T}_{12} &= (\mathbf{I} - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) \\
\mathbf{T}_{21} &= \mathbf{T}_{12}^{\top} = (\mathbf{I} - \frac{1}{2} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) \\
\mathbf{T}_{22} &= -\Lambda^{-1}.
\end{aligned}$$

Proof. Let

$$\mathbf{K} = \mathbf{T}^{-1} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{I} + \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \\ \mathbf{I} + \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}} & \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} \end{bmatrix},$$

then the inverse of the Schur complement $\mathbf{D} = \mathbf{K} / (\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}})$ is given by

$$\begin{aligned}
\mathbf{D}^{-1} &= (\mathbf{K}_{22} - \mathbf{K}_{21} \mathbf{K}_{11}^{-1} \mathbf{K}_{12})^{-1} \\
&= (\mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} - (\mathbf{I} + \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}}) (\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}})^{-1} (\mathbf{I} + \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}}))^{-1} \\
&= (\mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} - (\mathbf{I} - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1}) (\Lambda + \Pi) (\mathbf{I} - \frac{1}{2} \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}}))^{-1} \\
&= (\mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} - (\Lambda - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S) \Lambda^{-1} (\Lambda + \Pi) \Lambda^{-1} (\Lambda - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}}))^{-1} \\
&= (\mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}} - (\Lambda + \Pi) + \frac{1}{2} (\Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} (\Lambda + \Pi) + (\Lambda + \Pi) \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}}) \\
&\quad - \frac{1}{4} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} (\Lambda + \Pi) \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}})^{-1} \\
&= (\Pi - \Lambda - \Pi)^{-1} \\
&= -\Lambda^{-1},
\end{aligned}$$

where we used Lemma S4. By block matrix inversion and again with Lemma S4 we obtain

$$\begin{aligned}
\mathbf{T}_{11} &= (\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}})^{-1} + (\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}})^{-1} (\mathbf{I} + \mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{V}_{\mathbf{A}}) \mathbf{D}^{-1} (\mathbf{I} + \mathbf{V}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}}) (\mathbf{U}_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} \mathbf{U}_{\mathbf{A}})^{-1} \\
&= \Lambda (\Lambda + \Pi)^{-1} \Lambda + \Lambda (\mathbf{I} - \frac{1}{2} \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}}) \mathbf{D}^{-1} (\mathbf{I} - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1}) \Lambda \\
&= \Lambda (\Lambda + \Pi)^{-1} \Lambda + (\Lambda - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}}) \mathbf{D}^{-1} (\Lambda - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S)
\end{aligned}$$

as well as

$$\begin{aligned} T_{12} &= -(U_{\mathbf{A}}^{\top} A_0^{-1} U_{\mathbf{A}})^{-1} (I + U_{\mathbf{A}}^{\top} A_0^{-1} V_{\mathbf{A}}) D^{-1} \\ &= -\Lambda (\Lambda + \Pi)^{-1} \Lambda \Lambda^{-1} (\Lambda + \Pi) (I - \frac{1}{2} \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}}) D^{-1} \\ &= -(\Lambda - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}}) D^{-1} \end{aligned}$$

$$T_{21} = T_{12}^{\top} = -D^{-\top} (\Lambda - \frac{1}{2} \Delta_{\mathbf{A}}^{\top} S)$$

and finally $T_{22} = D^{-1} = -\Lambda^{-1}$. \square

Lemma S6

Given the assumptions of Corollary S3, additionally assume (1) and (2) hold. Let

$$F = A_0^{-1} [U_{\mathbf{A}} \quad V_{\mathbf{A}}] \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix} \begin{bmatrix} U_{\mathbf{A}}^{\top} \\ V_{\mathbf{A}}^{\top} \end{bmatrix} A_0^{-1},$$

where T is chosen as in Lemma S5, then if $S^{\top} A S = I$, we have

$$F = A_0^{-1} A S (I + \Pi)^{-1} S^{\top} A A_0^{-1} - S S^{\top}.$$

Proof. By expanding the quadratic and using Lemma S5, we obtain the terms

$$\begin{aligned} F_{11} &:= A_0^{-1} U_{\mathbf{A}} T_{11} U_{\mathbf{A}}^{\top} A_0^{-1} \\ &= A_0^{-1} U_{\mathbf{A}} \Lambda (\Lambda + \Pi)^{-1} \Lambda U_{\mathbf{A}}^{\top} A_0^{-1} - A_0^{-1} U_{\mathbf{A}} (I - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) (I - \frac{1}{2} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) U_{\mathbf{A}}^{\top} A_0^{-1} \\ &= A_0^{-1} A S (\Lambda + \Pi)^{-1} S^{\top} A A_0^{-1} - A_0^{-1} A S \Lambda^{-1} (I - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) (I - \frac{1}{2} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) \Lambda^{-1} S^{\top} A A_0^{-1} \\ &= A_0^{-1} A S (\Lambda + \Pi)^{-1} S^{\top} A A_0^{-1} - A_0^{-1} A S \Lambda^{-2} S^{\top} A A_0^{-1} \\ &\quad + \frac{1}{2} A_0^{-1} A S \Lambda^{-1} (S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} + \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) \Lambda^{-1} S^{\top} A A_0^{-1} \\ &\quad - \frac{1}{4} A_0^{-1} A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-2} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A A_0^{-1} \end{aligned}$$

$$\begin{aligned} F_{12} &:= A_0^{-1} U_{\mathbf{A}} T_{12} V_{\mathbf{A}}^{\top} A_0^{-1} \\ &= A_0^{-1} U_{\mathbf{A}} (I - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) V_{\mathbf{A}}^{\top} A_0^{-1} \\ &= A_0^{-1} A S \Lambda^{-1} (I - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) \Delta_{\mathbf{A}}^{\top} (I - \frac{1}{2} S U_{\mathbf{A}}^{\top}) A_0^{-1} \\ &= A_0^{-1} A S \Lambda^{-1} (I - \frac{1}{2} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1}) \Delta_{\mathbf{A}}^{\top} (I - \frac{1}{2} S \Lambda^{-1} S^{\top} A) A_0^{-1} \\ &= A_0^{-1} A S \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} A_0^{-1} - \frac{1}{2} A_0^{-1} A S \Lambda^{-1} (S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} + \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A) A_0^{-1} \\ &\quad + \frac{1}{4} A_0^{-1} A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A A_0^{-1} \end{aligned}$$

$$\begin{aligned} F_{21} &:= F_{12}^{\top} = A_0^{-1} (I - \frac{1}{2} A S \Lambda^{-1} S^{\top}) \Delta_{\mathbf{A}} (I - \frac{1}{2} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S) \Lambda^{-1} S A A_0^{-1} \\ &= A_0^{-1} \Delta_{\mathbf{A}} \Lambda^{-1} S^{\top} A A_0^{-1} - \frac{1}{2} A_0^{-1} (\Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S + A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}}) \Lambda^{-1} S^{\top} A A_0^{-1} \\ &\quad + \frac{1}{4} A_0^{-1} A S \Lambda^{-1} S^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} S \Lambda^{-1} S^{\top} A A_0^{-1} \end{aligned}$$

$$\begin{aligned} F_{22} &:= A_0^{-1} V_{\mathbf{A}} T_{22} V_{\mathbf{A}}^{\top} A_0^{-1} \\ &= -A_0^{-1} (I - \frac{1}{2} U_{\mathbf{A}} S^{\top}) \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} (I - \frac{1}{2} S U_{\mathbf{A}}^{\top}) A_0^{-1} \\ &= -A_0^{-1} (I - \frac{1}{2} A S \Lambda^{-1} S^{\top}) \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} (I - \frac{1}{2} S \Lambda^{-1} S^{\top} A) A_0^{-1} \end{aligned}$$

$$\begin{aligned}
&= -\mathbf{A}_0^{-1} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} \mathbf{A}_0^{-1} + \frac{1}{2} \mathbf{A}_0^{-1} (\mathbf{A} \mathbf{S} \Lambda^{-1} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} + \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \Lambda^{-1} \mathbf{S}^{\top} \mathbf{A}) \mathbf{A}_0^{-1} \\
&\quad - \frac{1}{4} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \Lambda^{-1} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Lambda^{-1} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \Lambda^{-1} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1}
\end{aligned}$$

Assuming $\mathbf{S}^{\top} \mathbf{A} \mathbf{S} = \mathbf{I}$, it holds that

$$\begin{aligned}
\mathbf{F}_{11} &= \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{I} + \mathbf{\Pi})^{-1} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} + \frac{1}{2} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{S}^{\top} \Delta_{\mathbf{A}} + \Delta_{\mathbf{A}}^{\top} \mathbf{S}) \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} \\
&\quad - \frac{1}{4} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1}
\end{aligned}$$

$$\begin{aligned}
\mathbf{F}_{12} &= \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} - \frac{1}{2} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} + \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A}) \mathbf{A}_0^{-1} \\
&\quad + \frac{1}{4} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1}
\end{aligned}$$

$$\begin{aligned}
\mathbf{F}_{21} &= \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \frac{1}{2} \mathbf{A}_0^{-1} (\Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} + \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}}) \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} \\
&\quad + \frac{1}{4} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1}
\end{aligned}$$

$$\begin{aligned}
\mathbf{F}_{22} &= \mathbf{A}_0^{-1} \Delta_{\mathbf{A}} \mathbf{S}^{\top} - \mathbf{A}_0^{-1} \Delta_{\mathbf{A}} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} + \frac{1}{2} (\mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} + \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A}) \mathbf{A}_0^{-1} \\
&\quad - \frac{1}{4} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1},
\end{aligned}$$

which leads to

$$\mathbf{F}_{11} + \mathbf{F}_{12} = \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{I} + \mathbf{\Pi})^{-1} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} + \frac{1}{2} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{S}^{\top} \Delta_{\mathbf{A}} \mathbf{S}^{\top} \mathbf{A} - \mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top}) \mathbf{A}_0^{-1}$$

$$\begin{aligned}
\mathbf{F}_{21} + \mathbf{F}_{22} &= \mathbf{A}_0^{-1} \Delta_{\mathbf{A}} \mathbf{S}^{\top} + \frac{1}{2} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} - \mathbf{S}^{\top} \Delta_{\mathbf{A}} \mathbf{S}^{\top} \mathbf{A}) \mathbf{A}_0^{-1} \\
&= \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} \mathbf{S}^{\top} - \mathbf{S} \mathbf{S}^{\top} + \frac{1}{2} \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{S}^{\top} \Delta_{\mathbf{A}} \Delta_{\mathbf{A}}^{\top} - \mathbf{S}^{\top} \Delta_{\mathbf{A}} \mathbf{S}^{\top} \mathbf{A}) \mathbf{A}_0^{-1}.
\end{aligned}$$

Finally, adding up the individual terms we obtain

$$\mathbf{F} = \mathbf{F}_{11} + \mathbf{F}_{12} + \mathbf{F}_{21} + \mathbf{F}_{22} = \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{I} + \mathbf{\Pi})^{-1} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} - \mathbf{S} \mathbf{S}^{\top}.$$

□

Theorem 2 (Weak Posterior Correspondence)

Let $\mathbf{W}_0^{\mathbf{H}} \in \mathbb{R}_{\text{sym}}^{n \times n}$ be positive definite. Assume $\mathbf{H}_0 = \mathbf{A}_0^{-1}$, and that $\mathbf{W}_0^{\mathbf{A}}, \mathbf{A}_0, \mathbf{W}_0^{\mathbf{H}}$ satisfy (1) and (2), then weak posterior correspondence holds for the symmetric Kronecker covariance.

Proof. First note that without loss of generality $\mathbf{S}^{\top} \mathbf{A} \mathbf{S} = \mathbf{I}$, i.e. only the direction of the action matters in Algorithm 1 not its magnitude. This can be seen from the forms of \mathbf{A}_k and \mathbf{H}_k in Section 2.1. Any positive factor $\alpha > 0$ of s_k cancels in the update expressions. Expanding the right hand side we have using (S33), that $\mathbf{H}_k \mathbf{Y} = \mathbf{S}$. Then by Lemma S3, Lemma S6 and $\mathbf{S}^{\top} \mathbf{A} \mathbf{S} = \mathbf{I}$, the left hand side evaluates to

$$\begin{aligned}
\mathbf{A}_k^{-1} \mathbf{Y} &= (\mathbf{A}_0^{-1} - \mathbf{F}) \mathbf{Y} \\
&= (\mathbf{A}_0^{-1} - \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} (\mathbf{I} + \mathbf{\Pi})^{-1} \mathbf{S}^{\top} \mathbf{A} \mathbf{A}_0^{-1} + \mathbf{S} \mathbf{S}^{\top}) \mathbf{A} \mathbf{S} \\
&= \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} - \mathbf{A}_0^{-1} \mathbf{A} \mathbf{S} + \mathbf{S} \\
&= \mathbf{S} \\
&= \mathbf{H}_k \mathbf{Y}.
\end{aligned}$$

This concludes the proof. □

This theorem shows that for a certain choice of symmetric matrix-variate normal prior the estimated inverse of the matrix \mathbf{H}_k corresponds to the inverse of the estimated matrix \mathbf{A}_k^{-1} . It also shows that both act like \mathbf{A}^{-1} on the space spanned by \mathbf{Y} , consistent with the interpretation of the two being the best guess for the inverse \mathbf{A}^{-1} .

S7 Galerkin's Method for PDEs

In the spirit of applying machine learning in the sciences [22], we briefly outlined an application of Algorithm 1 to the solution of partial differential equations in Section 4. As an example we considered the Dirichlet problem for the Poisson equation given by

$$\begin{cases} -\Delta u(x, y) = f(x, y) & (x, y) \in \text{int } \Omega \\ u(x, y) = u_{\partial\Omega}(x, y) & (x, y) \in \partial\Omega \end{cases} \quad (\text{S52})$$

where Ω is a connected open region with sufficiently regular boundary and $u_{\partial\Omega} : \partial\Omega \rightarrow \mathbb{R}$ defines the boundary conditions. The corresponding weak solution of (S52) is given by $u \in V$ such that for all test functions $v \in V$

$$a(u, v) := \int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx =: f(v), \quad (\text{S53})$$

where $a(\cdot, \cdot)$ is a bilinear form. Next, one derives the *Galerkin equation* by choosing a finite-dimensional subspace $V_{\square} \subset V$ and corresponding basis $e_1^{\square}, \dots, e_n^{\square}$. Then (S53) reduces to finding $u \in V_{\square}$ such that for all $i \in \{1, \dots, n\}$ it holds that $a(u, e_i^{\square}) = \sum_{j=1}^n u_j a(e_j^{\square}, e_i^{\square}) = f(e_i^{\square})$ which is a linear system $\mathbf{A}u = \mathbf{f}$ with the entries of the Gram matrix given by $\mathbf{A}_{ij} = a(e_j^{\square}, e_i^{\square})$ and $f_i = f(e_i^{\square})$.

S7.1 Operator View

The operator view provides another motivation for placing a distribution over the matrix \mathbf{A} of a linear system. When approximating the solution to a PDE, as we do here, then solution-based inference for linear systems [13, 14] can be viewed as placing a Gaussian process prior over the solution $u : \Omega \rightarrow \mathbb{R}$ [23]. The matrix-based approach [6] instead can be interpreted as placing a Gaussian measure [24] on the infinite-dimensional space of the differential operator instead. This induces a Gaussian distribution on the Gram matrix \mathbf{A} modelling the uncertainty about the actions of the (discretized) differential operator.

Definition S3 (Infinite-dimensional Gaussian Measures [24])

Let W be a topological vector space with Borel probability measure μ , then μ is Gaussian, iff for each continuous linear functional $f \in W^*$, the pushforward $\mu \circ f^{-1}$ is a Gaussian measure on \mathbb{R} , i.e. f is a Gaussian random variable on (W, \mathcal{B}_W, μ) .

This definition and further detail on Gaussian measures in infinite-dimensional spaces can be found in the book by Bogachev [24]. We now model the differential operator as a random variable on the space of bounded linear operators and show that this induces a distribution on the Gram matrix arising from discretization via Galerkin's method.

Theorem S3 (Gaussian Measures on the Space of Bounded Linear Operators)

Let V be a Hilbert space and let $W = B(V, V)$ be the space of bounded linear operators from V to V with Borel probability measure μ and let \mathbf{A} be a Gaussian random variable on (W, \mathcal{B}_W, μ) . Consider the operator equation

$$\mathbf{A}u = \mathbf{f}$$

and let $a : V \times V \rightarrow \mathbb{R}$, $(u, v) \mapsto \langle \mathbf{A}u, v \rangle_V = \langle \mathbf{f}, v \rangle_V$ be its corresponding bilinear form. Let V_{\square} be an n -dimensional subspace of V , then the resulting Gram matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is matrix-variate Gaussian.

Proof. Since V is Banach, so is W . Define the functional $a_W : W \rightarrow \mathbb{R}$ given by $a_W(\mathbf{A}, u, v) = a(u, v)$ for fixed $u, v \in V$. The map $a_W(\cdot, u, v)$ is linear by linearity of the inner product and bounded since using the Cauchy-Schwarz inequality, it holds that

$$|a_W(\mathbf{A}, u, v)| = |\langle \mathbf{A}u, v \rangle_V| \leq \|\mathbf{A}u\|_V \|v\|_V \leq \|\mathbf{A}\|_W \|u\|_V \|v\|_V = C \|\mathbf{A}\|_W.$$

Therefore $a_W(\cdot, u, v) \in W^*$ for all $u, v \in V$. By Definition S3 of a Gaussian measure the push forward $\mu \circ a_W^{-1}$ is a Gaussian measure on \mathbb{R} for all $u, v \in V$, in particular also for a basis $\{v_i\}_{i=1}^n$ of V_{\square} . Therefore the Gram matrix \mathbf{A} given by $\mathbf{A}_{ij} = a(v_i, v_j) = a_W(\mathbf{A}, v_i, v_j)$ is matrix-variate Gaussian since its components are Gaussian. \square

Remark S1

The Laplacian $\Delta : H^2(\Omega) \rightarrow L^2(\Omega)$ is a bounded linear operator on the Sobolev space $H^2(\Omega)$. Note, that in general differential operators are in fact *not bounded*. Hence, the simple argument above does not generalize to arbitrary differential operators.

Remark S2

If the bilinear form a in addition to being continuous is also weakly coercive, then by the Lax-Milgram theorem the operator equation has a unique solution. A symmetric and weakly coercive operator implies a symmetric positive-definite Gram matrix.

S7.2 Discretization Refinement

The linear system $\mathbf{A}\mathbf{u} = \mathbf{f}$ arises from discretizing (S52) using Galerkin's method on a given mesh \square defined via a finite-dimensional subspace $V_\square \subset V$ such that $\mathbf{u} \in V_\square$. By solving this problem using a probabilistic linear solver we obtain a posterior distribution over the inverse \mathbf{H} of the discretized differential operator \mathbf{A} . Our goal is to leverage the obtained information about the solution on the coarse mesh to extrapolate to a refined discretization, similar in spirit to multi-grid methods [25]. This approach can be seen as an instance of transfer learning and could be used for adaptive probabilistic mesh refinement strategies based on the uncertainty about the solution in a certain region of the mesh.

Consider a fine mesh \boxplus given by V_{\boxplus} , where $n_{\boxplus} = \dim(V_{\boxplus}) > \dim(V_\square) = n_\square$ such that $V_\square \subset V_{\boxplus} \subset V$. We would like to transfer information from solving the problem on the coarse mesh V_\square to the solution of the discretized PDE on the fine mesh V_{\boxplus} . To do so we compute the predictive distribution on the fine mesh, given the belief over the inverse differential operator on the coarse mesh, i.e.

$$p(\mathbf{H}_{\boxplus}) = \int p(\mathbf{H}_{\boxplus} \mid \mathbf{H}_\square) p(\mathbf{H}_\square) d\mathbf{H}_\square.$$

Define the *prolongation operator* $\mathbf{P} : \mathbb{R}^{n_\square} \rightarrow \mathbb{R}^{n_{\boxplus}}$ given by $\mathbf{P}_{ij} = \langle \mathbf{e}_i^{\boxplus}, \mathbf{e}_j^\square \rangle$ satisfying $\mathbf{P}^\top \mathbf{P} = \mathbf{I} \in \mathbb{R}^{n_\square \times n_\square}$, implying it is injective. The distribution over the inverse operator on the fine mesh given the inverse operator on the coarse mesh is given by

$$p(\mathbf{H}_{\boxplus} \mid \mathbf{H}_\square) = \mathcal{N}(\mathbf{H}_{\boxplus}; \mathbf{P}\mathbf{H}_\square\mathbf{P}^\top, \mathbf{\Lambda}) \quad (\text{S54})$$

where $\mathbf{\Lambda} \in \mathbb{R}_{\text{sym}}^{n_{\boxplus} \times n_{\boxplus}}$ positive definite models the numerical uncertainty induced by the coarser discretization. This corresponds to the assumption that solving the problem on a coarser grid approximates the solution on a fine grid projected to the coarse grid.

Now assume we have a posterior distribution over the inverse differential operator on the coarse grid from a solve of the coarse problem using Algorithm 1, given by

$$p(\mathbf{H}_\square) = \mathcal{N}(\mathbf{H}_\square; \mathbf{H}_\square^k, \mathbf{W}_\square^k \otimes \mathbf{W}_\square^k).$$

The projection in (S54) is a linear map, since by the characteristic property of the Kronecker product (S1) we have

$$\text{svec}(\mathbf{P}\mathbf{H}_\square\mathbf{P}^\top) = \mathbf{Q}(\mathbf{P} \otimes \mathbf{P})\mathbf{Q}^\top \text{svec}(\mathbf{H}_\square).$$

Therefore by Theorem S1 the predictive distribution is also closed-form and Gaussian.

Proposition S6 (Predictive Distribution on Fine Mesh)

Let $p(\mathbf{H}_\square) = \mathcal{N}(\mathbf{H}_\square; \mathbf{H}_\square^k, \mathbf{W}_\square^k \otimes \mathbf{W}_\square^k)$ be a prior on \mathbf{H}_\square and assume a likelihood of the form (S54). Then the predictive distribution is given by $p(\mathbf{H}_{\boxplus}) = \mathcal{N}(\mathbf{H}_{\boxplus}; \mathbf{H}_{\boxplus}^0, \mathbf{\Sigma}_{\boxplus}^0)$, where

$$\begin{aligned} \mathbf{H}_{\boxplus}^0 &= \mathbf{P}\mathbf{H}_\square^k\mathbf{P}^\top, \\ \mathbf{\Sigma}_{\boxplus}^0 &= \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top \otimes \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top + \mathbf{\Lambda}. \end{aligned}$$

Proof. By Theorem S1 we obtain for the mean and covariance of the predictive distribution

$$\begin{aligned} \mathbf{H}_{\boxplus}^0 &= \mathbf{P}\mathbf{H}_\square^k\mathbf{P}^\top \\ \mathbf{\Sigma}_{\boxplus}^0 &= \mathbf{Q}(\mathbf{P} \otimes \mathbf{P})\mathbf{Q}^\top (\mathbf{W}_\square^k \otimes \mathbf{W}_\square^k) \mathbf{Q}(\mathbf{P}^\top \otimes \mathbf{P}^\top) \mathbf{Q}^\top + \mathbf{\Lambda} \\ &= \frac{1}{2} \mathbf{Q}(\mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top \otimes \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top + \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top \boxtimes \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top) \mathbf{Q}^\top + \mathbf{\Lambda} \\ &= \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top \otimes \mathbf{P}\mathbf{W}_\square^k\mathbf{P}^\top + \mathbf{\Lambda} \end{aligned}$$

where we used (S31) and the symmetry of \mathbf{W}_\square^k . □

For general Λ the covariance of the predictive distribution does not have symmetric Kronecker form, making its use as a prior for a new solve on the fine mesh challenging. We aim to exploit structural assumptions on Λ and results on nearest Kronecker products to a sum of Kronecker products to remedy this shortcoming in the future.

References

- [1] Frederick Michael Larkin. Estimation of a non-negative function. *BIT Numerical Mathematics*, 9(1):30–52, 1969.
- [2] Persi Diaconis. Bayesian numerical analysis. *Statistical Decision Theory and Related Topics IV*, 1:163–175, 1988.
- [3] Anthony O’Hagan. Some Bayesian numerical analysis. *Bayesian Statistics*, 4:345–363, 1992.
- [4] Philipp Hennig, Mike A. Osborne, and Mark Girolami. Probabilistic numerics and uncertainty in computations. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 471(2179), 2015.
- [5] Chris Oates and Tim J. Sullivan. A modern retrospective on probabilistic numerics. *Statistics and Computing*, 10 2019.
- [6] Philipp Hennig. Probabilistic interpretation of linear solvers. *SIAM Journal on Optimization*, 25(1):234–260, 2015.
- [7] Harold V Henderson and Shayle R Searle. The vec-permutation matrix, the vec operator and Kronecker products: A review. *Linear and multilinear algebra*, 9(4):271–288, 1981.
- [8] Farid Alizadeh, Jean-Pierre A. Haeberly, and Michael L. Overton. Primal-dual interior-point methods for semidefinite programming: convergence rates, stability and numerical results. *SIAM Journal on Optimization*, 8(3):746–768, 1998.
- [9] Kathrin Schacke. On the Kronecker product. Technical report, University of Waterloo, 2013.
- [10] Charles F. Van Loan. The ubiquitous Kronecker product. *Journal of Computational and Applied Mathematics*, 123(1-2):85–100, 2000.
- [11] Peder A. Olsen, Steven J. Rennie, and Vaibhava Goel. Efficient automatic differentiation of matrix functions. In *Recent Advances in Algorithmic Differentiation*, pages 71–81. Springer, 2012.
- [12] Arjun K. Gupta and Daya K. Nagar. *Matrix-variate distributions*. Chapman and Hall/CRC, 2000.
- [13] Jon Cockayne, Chris Oates, Ilse C. Ipsen, and Mark Girolami. A Bayesian conjugate gradient method. *Bayesian Analysis*, 14(3):937–1012, 2019.
- [14] Simon Bartels, Jon Cockayne, Ilse C. Ipsen, and Philipp Hennig. Probabilistic linear solvers: A unifying view. *Statistics and Computing*, 29(6):1249–1263, 2019.
- [15] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, 2006.
- [16] Philipp Hennig and Martin Kiefel. Quasi-Newton method: A new direction. *Journal of Machine Learning Research*, 14(Mar):843–865, 2013.
- [17] Matthias Seeger. Low rank updates for the Cholesky decomposition. Technical report, University of California at Berkeley, 2008.
- [18] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- [19] Larry Nazareth. A relationship between the BFGS and conjugate gradient algorithms and its implications for new algorithms. *SIAM Journal on Numerical Analysis*, 16(5):794–800, 1979.

- [20] John E. Dennis, Jr and Jorge J. Moré. Quasi-Newton methods, motivation and theory. *SIAM review*, 19(1):46–89, 1977.
- [21] Keith T. Conrad. The minimal polynomial and some applications. Technical report, University of Connecticut, 2008.
- [22] Giuseppe Carleo, Ignacio Cirac, Kyle Cranmer, Laurent Daudet, Maria Schuld, Naftali Tishby, Leslie Vogt-Maranto, and Lenka Zdeborová. Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4):045002, 2019.
- [23] Mark Girolami, Eky Febrianto, Ge Yin, and Fehmi Cirak. The statistical finite element method (statFEM) for coherent synthesis of observation data and model predictions. *arXiv pre-print*, art. arXiv:1905.06391, May 2020.
- [24] Vladimir Igorevich Bogachev. *Gaussian measures*, volume 62 of *Mathematical Surveys and Monographs*. American Mathematical Society, 1998.
- [25] Pieter Wesseling. *An Introduction to Multigrid Methods*. R.T. Edwards, 2004.