
Supplemental material: Efficient inference in matrix-variate Gaussian models with iid observation noise

Oliver Stegle¹
Max Planck Institutes
Tübingen, Germany
stegle@tuebingen.mpg.de

Christoph Lippert¹
Max Planck Institutes
Tübingen, Germany
clippert@tuebingen.mpg.de

Joris Mooij
Institute for Computing and Information
Radboud University Nijmegen
Nijmegen, The Netherlands
j.mooij@cs.ru.nl

Neil Lawrence
Department of Computer Science
University of Sheffield
Sheffield, UK
N.Lawrence@sheffield.ac.uk

Karsten Borgwardt
Max Planck Institutes & Eberhard Karls Universität
Tübingen, Germany
karsten.borgwardt@tuebingen.mpg.de

1 Preliminaries

- $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{U}^T$ is the eigenvalue decomposition of the symmetric M -by- M matrix \mathbf{A} , where \mathbf{U} is an M -by- M orthonormal matrix, of the M eigenvectors of \mathbf{A} and \mathbf{S} is an M -by- M diagonal matrix, holding the corresponding eigenvalues of \mathbf{A} as diagonal entries.
- $\mathbf{A} \odot \mathbf{B}$ is the pointwise or Hadamard product of \mathbf{A} and \mathbf{B} .
- $\mathbf{A} \otimes \mathbf{B}$ is the Kronecker product of \mathbf{A} and \mathbf{B} .
- $\mathbf{Y} \in \mathbb{R}^{N \times D}$ is the matrix of all samples, having N rows and D columns.

- $\text{vec}(\mathbf{A}) = \begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{LM} \end{pmatrix}$ is the vec-operator, that concatenates the columns of a matrix \mathbf{A} in a column vector.

- $\text{vec}_{L \times M}^{-1} \begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{LM} \end{pmatrix} = (\mathbf{A})$ is the inverse vec-operator, that arranges the entries of a vector in an $L \times M$ matrix \mathbf{A} .

2 Covariance estimation in Kronecker Gaussian processes

Here, we give details of an efficient implementation of the tensor Gaussian process model derived in the main paper.

¹These authors contributed equally to this work.

The basic model

We start with a general form a Gaussian process model where the covariance has a Kronecker structure:

$$\ln p(\mathbf{Y} | \mathbf{R}, \mathbf{C}, \sigma^2) = \mathcal{N}(\text{vec}(\mathbf{Y}) | \mathbf{0}, \mathbf{C}(\boldsymbol{\Theta}_{\mathbf{C}}) \otimes \mathbf{R}(\boldsymbol{\Theta}_{\mathbf{R}}) + \sigma^2 \mathbf{I}). \quad (1)$$

Here, \mathbf{Y} is the data matrix with N rows (samples) and D columns (features). We defined $\mathbf{R}(\boldsymbol{\Theta}_{\mathbf{R}})$ as the row ‘‘row covariance’’ of the data matrix and $\mathbf{C}(\boldsymbol{\Theta}_{\mathbf{C}})$ corresponds to the ‘‘column covariance’’.

For notational convenience we will drop the dependence of the covariance matrices on additional hyperparameters $\boldsymbol{\Theta}_{\mathbf{R}}$ and $\boldsymbol{\Theta}_{\mathbf{C}}$, respectively. Furthermore, we will make the simplifying assumption that \mathbf{C} is kept constant, i.e. has no parameters that need to be adapted during learning. Importantly, this is no restriction for general solutions as all calculations can be performed with respect to other covariance as well, for example iteratively optimizing hyperparameters of \mathbf{R} and \mathbf{C} in turn.

To implement parameter optimization of the covariance parameters of the model in Equation (1), we require efficient evaluation of the marginal likelihood and the gradients with respect to hyperparameters.

2.1 Efficient evaluation of the marginal likelihood

The term we want to evaluate is the log-marginal-likelihood, given by the log of the multivariate Normal density

$$\mathcal{L}(\boldsymbol{\Theta}_{\mathbf{R}}, \boldsymbol{\Theta}_{\mathbf{C}}, \sigma^2) = -\frac{N \cdot D}{2} \ln 2\pi - \frac{1}{2} \underbrace{\ln |\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}_{N \cdot D}|}_{\text{log-det}} - \frac{1}{2} \underbrace{\text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}_{N \cdot D})^{-1} \text{vec}(\mathbf{Y})}_{\text{squared form}}. \quad (2)$$

In the following, we make heavy use of the eigenvalue decomposition of $\mathbf{C} \otimes \mathbf{R} = \mathbf{U} \mathbf{S} \mathbf{U}^T$, where \mathbf{U} is an $N \cdot D$ -by- $N \cdot D$ orthonormal matrix, holding the eigenvectors of \mathbf{C} and \mathbf{S} is an N -by- ND diagonal matrix, holding the corresponding eigenvalues on the diagonal. This decomposition can be efficiently obtained from the composition of the individual kronecker terms (after some reordering), i.e. $\mathbf{C} \otimes \mathbf{R} = (\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}})(\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}})(\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T)$.

We derive efficient solutions for the logarithm of the determinant of $\mathbf{C} \otimes \mathbf{R}$ and the squared form separately.

2.1.1 Efficient evaluation of the log-det

Assuming that we have the eigenvalue decomposition for \mathbf{R} and \mathbf{C} , the logarithm of the determinant can be written as

$$\begin{aligned} \ln |\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}_{N \cdot D}| &= \ln |(\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}})(\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}})(\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T) + \sigma^2 \mathbf{I}_{N \cdot D}| \\ &= \ln |(\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}})| + \ln |(\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I}_{N \cdot D})| + \ln |\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T| \\ &= \ln |(\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I}_{N \cdot D})| \\ &= \sum_{r=1}^N \sum_{c=1}^D -\ln(\mathbf{S}_{\mathbf{R}}[r, r] \cdot \mathbf{S}_{\mathbf{C}}[c, c] + \sigma^2). \end{aligned} \quad (3)$$

This term can be evaluated in $O(N \cdot D)$.

2.1.2 Efficient evaluation of the squared form

Also the squared form in the log marginal likelihood can be evaluated efficiently as follows:

$$\begin{aligned}
\text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) &= \text{vec}(\mathbf{Y})^T \mathbf{U} (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} \mathbf{U}^T \text{vec}(\mathbf{Y}) \\
\text{Using orthogonality of } \mathbf{U} = (\mathbf{U}^T)^{-1}: & \\
&= \text{vec}(\mathbf{Y})^T (\mathbf{U}_C \otimes \mathbf{U}_R) (\mathbf{S}_C \otimes \mathbf{S}_R + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_C^T \otimes \mathbf{U}_R^T) \text{vec}(\mathbf{Y}) \\
\text{using the vectorization identities of kronecker structures} & \\
&= (\text{vec}(\mathbf{U}_R^T \mathbf{Y} \mathbf{U}_C))^T (\mathbf{S}_C \otimes \mathbf{S}_R + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_R^T \mathbf{Y} \mathbf{U}_C) \tag{4}
\end{aligned}$$

We only need to compute the diagonal of $(\mathbf{S}_C \otimes \mathbf{S}_R + \sigma^2 \mathbf{I})^{-1}$, which already has been computed before. As this term involves the rotation of the data matrix \mathbf{Y} , it can be evaluated in $O(N^2 D + ND^2)$.

2.2 Efficient evaluation of the gradients w.r.t. to covariance hyperparameters

Here, the aim is to evaluate the gradient of Equation (1) w.r.t. to a particular column covariance parameter $\theta_C \in \Theta_C$, row covariance parameter $\theta_R \in \Theta_R$ and the noise parameter σ^2 .

2.3 Derivatives w.r.t. noise variance σ^2

2.3.1 Derivatives of the log-det term w.r.t. noise variance σ^2

$$\begin{aligned}
\frac{d}{d\sigma^2} \ln |\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}| &= \text{Tr} \left[(\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \frac{d}{d\sigma^2} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}) \right] \\
&= \text{Tr} \left[(\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \frac{d}{d\sigma^2} (\sigma^2 \mathbf{I}) \right] \\
&= \text{Tr} [\mathbf{U} (\mathbf{S}_C \otimes \mathbf{S}_R + \sigma^2 \mathbf{I})^{-1} \mathbf{U}^T \mathbf{I}]
\end{aligned}$$

Using that the sum of the eigenvalues of a matrix equals its trace, the derivative becomes

$$\sum_{r=1}^N \sum_{c=1}^D \frac{1}{\mathbf{S}_R[r, r] \cdot \mathbf{S}_C[c, c] + \sigma^2}. \tag{5}$$

2.3.2 Squared form derivatives w.r.t. noise variance σ^2

$$\frac{d}{d\sigma^2} \text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) = \text{vec}(\mathbf{Y})^T \frac{d}{d\sigma^2} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}).$$

Using $\frac{d}{d\sigma^2} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \left[\frac{d}{d\sigma^2} \mathbf{A} \right] \mathbf{A}^{-1}$, this becomes

$$\begin{aligned}
& - \text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \left[\frac{d}{d\sigma^2} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}) \right] (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \\
&= - \text{vec}(\mathbf{Y})^T (\mathbf{U}_C \otimes \mathbf{U}_R) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_C^T \otimes \mathbf{U}_R^T) \left(\frac{d}{d\sigma^2} \sigma^2 \mathbf{I} \right) (\mathbf{U}_C \otimes \mathbf{U}_R) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_C^T \otimes \mathbf{U}_R^T) \text{vec}(\mathbf{Y}) \\
&= - \text{vec}(\mathbf{U}_R^T \mathbf{Y} \mathbf{U}_C)^T (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} \left(\frac{d}{d\sigma^2} \sigma^2 \mathbf{I} \right) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_R^T \mathbf{Y} \mathbf{U}_C) \tag{6}
\end{aligned}$$

2.4 Derivatives w.r.t a row covariance parameter $\theta_{\mathbf{R}}$

2.4.1 Derivatives of the determinant w.r.t. a row covariance parameter $\theta_{\mathbf{R}}$

$$\begin{aligned}
\frac{d}{d\theta_{\mathbf{R}}} \ln |\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}| &= \text{Tr} \left[(\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \frac{d}{d\theta_{\mathbf{R}}} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}) \right] \\
&= \text{Tr} \left[(\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^{\text{T}} \otimes \mathbf{U}_{\mathbf{R}}^{\text{T}}) (\mathbf{C} \otimes \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}) \right]. \\
\text{Using the identity } (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) &= \mathbf{AC} \otimes \mathbf{BD} \text{ this equals} \\
\text{Tr} \left[(\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^{\text{T}} \mathbf{C} \otimes \mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}) \right] \\
\text{Using } \text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}), \text{ this becomes} \\
\text{Tr} \left[(\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^{\text{T}} \mathbf{C} \otimes \mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}) (\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) \right] \\
&= \text{Tr} \left[(\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^{\text{T}} \mathbf{C} \mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \right] \\
&= \text{Tr} \left[(\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} (\mathbf{S}_{\mathbf{C}} \otimes (\mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}})) \right] \\
&= \text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \text{diag}(\mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \right) \quad (7)
\end{aligned}$$

As this derivation only involves the trace, we just need the diagonal of the Kronecker product, which only involves the diagonal of $(\mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}})$.

For the special case, where $\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}$ has only one non-zero row:

$$\begin{aligned}
&\text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \text{diag}(\mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \right) \\
&= \text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \text{diag}([\mathbf{U}_{\mathbf{R}}]_i^{\text{T}} [\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \mathbf{U}_{\mathbf{R}}) \right) \\
&= \text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \left([\mathbf{U}_{\mathbf{R}}]_i \odot \left([\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \mathbf{U}_{\mathbf{R}} \right) \right)^{\text{T}} \right) \quad (8)
\end{aligned}$$

For the special case, where $\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}$ has only one non-zero column:

$$\begin{aligned}
&\text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \text{diag}(\mathbf{U}_{\mathbf{R}}^{\text{T}} \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \right) \\
&= \text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \left(\left(\mathbf{U}_{\mathbf{R}}^{\text{T}} [\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \right) \odot [\mathbf{U}_{\mathbf{R}}]_i^{\text{T}} \right) \right) \quad (9)
\end{aligned}$$

So for the case, where $\mathbf{K}_{\mathbf{R}} = \mathbf{X}\mathbf{X}^{\text{T}}$, the derivative w.r.t. an entry of \mathbf{X} is:

$$\begin{aligned}
&\text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \left([\mathbf{U}_{\mathbf{R}}]_i \odot \left([\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \mathbf{U}_{\mathbf{R}} \right) \right)^{\text{T}} \right) + \\
&\text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \left(\left(\mathbf{U}_{\mathbf{R}}^{\text{T}} [\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \right) \odot [\mathbf{U}_{\mathbf{R}}]_i^{\text{T}} \right) \right) \\
&= 2 \text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right)^{\text{T}} \left(\text{diag}(\mathbf{S}_{\mathbf{C}}) \otimes \left(\left(\mathbf{U}_{\mathbf{R}}^{\text{T}} [\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \right) \odot [\mathbf{U}_{\mathbf{R}}]_i^{\text{T}} \right) \right) \\
&= 2 \text{vec} \left(\left(\left(\mathbf{U}_{\mathbf{R}}^{\text{T}} [\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}]_i \right)^{\text{T}} \odot [\mathbf{U}_{\mathbf{R}}]_i \right) \text{vec}_{N \times D}^{-1} \left(\text{diag} \left((\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \right) \text{diag}(\mathbf{S}_{\mathbf{C}}) \right) \right), \quad (10)
\end{aligned}$$

where we rearranged the transposes and used $\mathbf{A} \odot \mathbf{B} = \mathbf{B} \odot \mathbf{A}$.

2.4.2 Derivatives of the squared form w.r.t. a row covariance parameter $\theta_{\mathbf{R}}$

$$\begin{aligned}
& \frac{d}{d\theta_{\mathbf{R}}} \text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) = \text{vec}(\mathbf{Y})^T \frac{d}{d\theta_{\mathbf{R}}} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \\
& \text{using } \frac{d}{d\theta_{\mathbf{R}}} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{A} \right] \mathbf{A}^{-1} \\
& = -\text{vec}(\mathbf{Y})^T (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \left[\frac{d}{d\theta_{\mathbf{R}}} (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I}) \right] (\mathbf{C} \otimes \mathbf{R} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \\
& = -\text{vec}(\mathbf{Y})^T (\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T) (\mathbf{C} \otimes \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}) (\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T) \text{vec}(\mathbf{Y}) \\
& = -\text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}})^T (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^T \otimes \mathbf{U}_{\mathbf{R}}^T) (\mathbf{C} \otimes \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}) (\mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}) (\mathbf{S} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}}) \\
& = -(\text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}}))^T (\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} (\mathbf{U}_{\mathbf{C}}^T \mathbf{C} \mathbf{U}_{\mathbf{C}} \otimes \mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) (\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}}) \\
& = -\text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}})^T (\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} (\mathbf{S}_{\mathbf{C}} \otimes (\mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}})) (\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}}) \\
& = -\text{vec}(\tilde{\mathbf{Y}})^T (\mathbf{S}_{\mathbf{C}} \otimes (\mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}})) \text{vec}(\tilde{\mathbf{Y}}) \\
& = -\text{vec}(\tilde{\mathbf{Y}})^T \text{vec}((\mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}), \tag{11}
\end{aligned}$$

where $\text{vec}(\tilde{\mathbf{Y}}) = (\mathbf{S}_{\mathbf{C}} \otimes \mathbf{S}_{\mathbf{R}} + \sigma^2 \mathbf{I})^{-1} \text{vec}(\mathbf{U}_{\mathbf{R}}^T \mathbf{Y} \mathbf{U}_{\mathbf{C}})$.

For the special case, where the derivative of \mathbf{R} is a row matrix:

$$\begin{aligned}
& -\text{vec}(\tilde{\mathbf{Y}})^T \text{vec}((\mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}) \\
& = -\text{vec}(\tilde{\mathbf{Y}})^T \text{vec}\left(\left([\mathbf{U}_{\mathbf{R}}]_{i:}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{i:} \mathbf{U}_{\mathbf{R}}\right) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}\right) \\
& = -\text{Tr}\left(\left(\tilde{\mathbf{Y}}^T [\mathbf{U}_{\mathbf{R}}]_{i:}^T\right) \left(\left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{i:} \mathbf{U}_{\mathbf{R}} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}\right)\right) \\
& = -\left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{i:} \mathbf{U}_{\mathbf{R}} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}} \tilde{\mathbf{Y}}^T [\mathbf{U}_{\mathbf{R}}]_{i:}^T, \tag{12}
\end{aligned}$$

where we use the fact that the trace of an outer product between vectors equals to the inner product between the vectors.

For the special case, where the derivative of \mathbf{R} is a column matrix:

$$\begin{aligned}
& = -\text{vec}(\tilde{\mathbf{Y}})^T \text{vec}((\mathbf{U}_{\mathbf{R}}^T \frac{d}{d\theta_{\mathbf{R}}} \mathbf{R} \mathbf{U}_{\mathbf{R}}) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}) \\
& = -\text{vec}(\tilde{\mathbf{Y}})^T \text{vec}\left(\left(\mathbf{U}_{\mathbf{R}}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i} [\mathbf{U}_{\mathbf{R}}]_{i:}\right) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}\right) \\
& = -\text{Tr}\left(\tilde{\mathbf{Y}}^T (\mathbf{U}_{\mathbf{R}}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i} [\mathbf{U}_{\mathbf{R}}]_{i:}) \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}}\right) \\
& = -[\mathbf{U}_{\mathbf{R}}]_{i:} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}} \tilde{\mathbf{Y}}^T \mathbf{U}_{\mathbf{R}}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i} \tag{13}
\end{aligned}$$

So for the derivative w.r.t. entries in hidden factors \mathbf{X} in a linear kernel $\mathbf{K}_{\mathbf{R}} = \mathbf{X}\mathbf{X}^T$ we have to compute the following term:

$$\begin{aligned}
& -\left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i} \mathbf{U}_{\mathbf{R}} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}} \tilde{\mathbf{Y}}^T [\mathbf{U}_{\mathbf{R}}]_{i:}^T - [\mathbf{U}_{\mathbf{R}}]_{i:} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}} \tilde{\mathbf{Y}}^T \mathbf{U}_{\mathbf{R}}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i} \\
& - 2[\mathbf{U}_{\mathbf{R}}]_{i:} \tilde{\mathbf{Y}} \mathbf{S}_{\mathbf{C}} \tilde{\mathbf{Y}}^T \mathbf{U}_{\mathbf{R}}^T \left[\frac{d}{d\theta_{\mathbf{R}}} \mathbf{R}\right]_{:i}, \tag{14}
\end{aligned}$$

where we observe that the second term is the transpose of the first one.