

# Derivatives of Gaussian KL-Divergence for some parameterizations of the posterior covariance for variational Gaussian-process inference

M. Rule

March 25, 2020

*These notes provide the derivatives of the KL-divergence  $D_{\text{KL}} [Q(\mathbf{z})||P(\mathbf{z})]$  between two multivariate Gaussian distributions  $Q(\mathbf{z})$  and  $P(\mathbf{z})$  with respect to a few parameterizations  $\theta$  of the covariance matrix  $\Sigma(\theta)$  of  $Q$ . This is useful for variational Gaussian process inference, where clever parameterizations of the posterior covariance are required to make the problem tractable. Tables for differentiating matrix-valued functions can be found in [The Matrix Cookbook](#).*

Consider two multivariate Gaussian distributions  $Q(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_q, \Sigma(\theta))$  and  $P(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0 = \Lambda^{-1})$  with dimension  $L$ . The KL divergence  $D_{\text{KL}} [Q(\mathbf{z})||P(\mathbf{z})]$  [has the closed form](#)

$$\begin{aligned} \mathcal{D} &:= D_{\text{KL}} [Q(\mathbf{z})||\text{Pr}(\mathbf{z})] \\ &= \frac{1}{2} \left\{ (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_q)^\top \Lambda (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_q) \right. \\ &\quad \left. + \text{tr}(\Lambda \Sigma) - \ln |\Sigma| - \ln |\Lambda| \right\} + \text{constant}. \end{aligned} \tag{1}$$

In variational Bayesian inference, we minimize  $\mathcal{D}$  while maximizing the expected log-probability of some observations with respect to  $Q(\mathbf{z})$ . Closed-form derivatives of  $\mathcal{D}$  in terms of the parameters of  $Q$  are useful for manually optimizing code for larger problems. The derivatives of  $\mathcal{D}$  in terms of  $\boldsymbol{\mu}_q$  are straightforward:  $\partial_{\boldsymbol{\mu}_q} \mathcal{D} = \Lambda(\boldsymbol{\mu}_q - \boldsymbol{\mu}_z)$  and  $\text{H}_{\boldsymbol{\mu}_q} \mathcal{D} = \Lambda$ . In these notes, we explore derivatives of  $\mathcal{D}$  with respect to a few different parameterizations (“ $\theta$ ”) of  $\Sigma(\theta)$ .

We evaluate the following parameterizations for  $\Sigma$ : 1. Optimizing the full  $\Sigma$  directly 2.  $\Sigma \approx \mathbf{X}\mathbf{X}^\top$  3.  $\Sigma \approx \mathbf{A}^\top \text{diag}[\mathbf{v}]\mathbf{A}$  4.  $\Sigma \approx [\Lambda + \text{diag}[\mathbf{p}]]^{-1}$  5.  $\mathbf{F}^\top \mathbf{Q}\mathbf{Q}^\top \mathbf{F}$ , where  $\mathbf{Q} \in \mathbb{R}^{K \times K}$ ,  $K < L$  and  $\mathbf{F} \in \mathbb{R}^{K \times L}$ ,  $\mathbf{F}\mathbf{F}^\top = \mathbf{I}$ .

## 0.1 $\Sigma$

We first obtain gradients of  $\mathcal{D}$  in  $\Sigma$  (assuming  $\Sigma$  is full-rank). These can be used to derive gradients in  $\theta$  for some parameterizations  $\Sigma(\theta)$  using the chain rule. The gradient of  $\mathcal{D}$  in  $\Sigma$  can be obtained using identities (57) and (100) in [The Matrix Cookbook](#):

$$\begin{aligned} \partial_\Sigma \mathcal{D} &= \partial_\Sigma \{ \text{tr}(\Lambda \Sigma) - \ln |\Sigma| \} \\ &= \frac{1}{2} (\Lambda - \Sigma^{-1}). \end{aligned} \tag{2}$$

The Hessian in  $\Sigma$  is a fourth-order tensor. It’s simpler to express the Hessian in terms of a Hessian-vector product, which can be used with [Krylov subspace](#) solvers to efficiently compute the update in Newton’s method. Considering an  $L \times L$  matrix  $\mathbf{M}$ , the Hessian-vector product is given by

$$[\mathbf{H}_\Sigma \mathcal{D}] \mathbf{M} = \partial_\Sigma \langle \partial_\Sigma \mathcal{D}, \mathbf{M} \rangle = \partial_\Sigma \text{tr} [(\partial_\Sigma \mathcal{D})^\top \mathbf{M}], \quad (3)$$

where  $\langle \cdot, \cdot \rangle$  denotes the scalar (Frobenius) product. This is given by identity (124) in the Matrix Cookbook:

$$\partial_\Sigma \text{tr} \left[ \frac{1}{2} (\mathbf{\Lambda} - \Sigma^{-1})^\top \mathbf{M} \right] = -\frac{1}{2} \partial_\Sigma \text{tr} [\Sigma^{-1} \mathbf{M}] = \frac{1}{2} \Sigma^{-1} \mathbf{M}^\top \Sigma^{-1}. \quad (4)$$

## 0.2 $\Sigma \approx \mathbf{X}\mathbf{X}^\top$

We consider an approximate posterior covariance of the form

$$\Sigma \approx \mathbf{X}\mathbf{X}^\top, \quad \mathbf{X} \in \mathbb{R}^{L \times K} \quad (5)$$

where  $\mathbf{X}$  is a rank- $K < L$  matrix with  $L$  rows and  $K$  columns.

Since  $\mathbf{X}$  is not full rank, the log-determinant  $\ln |\Sigma| = \ln |\mathbf{X}\mathbf{X}^\top|$  in (1) diverges, due to the zero eigenvalues in the null space of  $\mathbf{X}$ . However, since this null-space is not being optimized, it does not affect our gradient. It is sufficient to replace the log-determinant with that of the reduced-rank representation,  $\ln |\mathbf{X}^\top \mathbf{X}|$ . Identity (55) in The Matrix Cookbook provides the derivative of this,  $\partial_{\mathbf{X}} \ln |\mathbf{X}^\top \mathbf{X}| = 2\mathbf{X}^{+\top}$ , where  $(\cdot)^+$  is the pseudoinverse. Combined with identity (112), this gives the following gradient of  $\mathcal{D}(\mathbf{X})$ :

$$\partial_{\mathbf{X}} \mathcal{D} = \partial_{\mathbf{X}} \frac{1}{2} \{ \text{tr} [\mathbf{\Lambda} \mathbf{X} \mathbf{X}^\top] - \ln |\mathbf{X}^\top \mathbf{X}| \} = \mathbf{\Lambda} \mathbf{X} - \mathbf{X}^{+\top}. \quad (6)$$

The Hessian-vector product requires the derivative of  $\partial_{\mathbf{X}} \text{tr} [\mathbf{X}^+ \mathbf{M}]$ :

$$\partial_{\mathbf{X}} \langle \partial_\Sigma \mathcal{D}, \mathbf{M} \rangle = \partial_{\mathbf{X}} \text{tr} \left[ \left( \mathbf{\Lambda} \mathbf{X} - \mathbf{X}^{+\top} \right)^\top \mathbf{M} \right] = \partial_{\mathbf{X}} \text{tr} [\mathbf{\Lambda} \mathbf{X} \mathbf{M}] - \partial_{\mathbf{X}} \text{tr} [\mathbf{X}^+ \mathbf{M}]. \quad (7)$$

Goulob and Pereya (1972) Eq. 4.12 gives the derivative of a fixed-rank pseudoinverse:

$$\partial \mathbf{X}^+ = -\mathbf{X}^+ (\partial \mathbf{X}) \mathbf{X}^+ + \mathbf{X}^+ \mathbf{X}^{+\top} (\partial \mathbf{X})^\top (1 - \mathbf{X} \mathbf{X}^+) + (1 - \mathbf{X}^+ \mathbf{X}) (\partial \mathbf{X})^\top \mathbf{X}^{+\top} \mathbf{X}^+ \quad (8)$$

Since  $\mathbf{X}$  is  $N \times K$  with rank  $K$ ,  $\mathbf{X}^+ \mathbf{X}$  is full-rank. Therefore  $\mathbf{X}^+ \mathbf{X} = \mathbf{I}_K$  and the final term in (8) vanishes. The derivative of the pseudoinverse can now be written as:

$$\partial \mathbf{X}^+ = -\mathbf{X}^+ (\partial \mathbf{X}) \mathbf{X}^+ + \mathbf{X}^+ \mathbf{X}^{+\top} (\partial \mathbf{X})^\top (\mathbf{I}_n - \mathbf{X} \mathbf{X}^+) \quad (9)$$

Since the derivative of a trace of a matrix-valued function is just the (transpose) of the scalar derivative,

$$\begin{aligned} \partial_{\mathbf{X}} \text{tr} [\mathbf{X}^+ \mathbf{M}] &= \{ -\mathbf{X}^+ \mathbf{M} \mathbf{X}^+ + \mathbf{X}^+ \mathbf{X}^{+\top} \mathbf{M}^\top (\mathbf{I}_n - \mathbf{X} \mathbf{X}^+) \}^\top \\ &= -\mathbf{X}^{+\top} \mathbf{M}^\top \mathbf{X}^{+\top} + (\mathbf{I} - \mathbf{X}^{+\top} \mathbf{X}^\top) \mathbf{M} \mathbf{X}^+ \mathbf{X}^{+\top}. \end{aligned} \quad (10)$$

Overall, we obtain the following Hessian-vector product:

$$\partial_{\mathbf{X}} \langle \partial_\Sigma \mathcal{D}, \mathbf{M} \rangle = \mathbf{\Lambda} \mathbf{M}^\top + \mathbf{X}^{+\top} \mathbf{M}^\top \mathbf{X}^{+\top} - (\mathbf{I} - \mathbf{X}^{+\top} \mathbf{X}^\top) \mathbf{M} \mathbf{X}^+ \mathbf{X}^{+\top} \quad (11)$$

### 0.2.1 $\Sigma = \mathbf{X}\mathbf{X}^\top$ when $\mathbf{X}$ is full-rank

Equations (6) and (11) are also valid if  $\mathbf{X}$  is a rank- $L$  triangular (Choleskey) factorization of  $\Sigma$ . In this case the pseudoinverse can be replaced by the full inverse, and various terms simplify:

$$\begin{aligned}\partial_{\mathbf{X}}\mathcal{D} &= \Lambda\mathbf{X} - \mathbf{X}^{-\top} \\ \partial_{\mathbf{X}}\langle\partial_{\mathbf{X}}\mathcal{D}, \mathbf{M}\rangle &= \Lambda\mathbf{M}^\top + \mathbf{X}^{-\top}\mathbf{M}^\top\mathbf{X}^{-\top}\end{aligned}\quad (12)$$

### 0.3 $\Sigma = \mathbf{A}^\top \text{diag}[\mathbf{v}]\mathbf{A}$

Let  $\Sigma = \mathbf{A}^\top \text{diag}[\mathbf{v}]\mathbf{A}$ , where  $\mathbf{A}$  is fixed and  $\mathbf{v} \in \mathbb{R}^L$  are free parameters. Define  $\text{diag}[\cdot]$  as an operator that constructs a diagonal matrix from a vector, or extracts the main diagonal from a matrix if its argument is a matrix. The gradient of  $\mathcal{D}$  in  $\mathbf{v}$  is:

$$\begin{aligned}\partial_{\mathbf{X}}\mathcal{D} &= \partial_{\mathbf{X}}\frac{1}{2} \{ \text{tr} [\Lambda\mathbf{A}^\top \text{diag}[\mathbf{v}]\mathbf{A}] - \ln |\mathbf{A}^\top \text{diag}[\mathbf{v}]\mathbf{A}| \} \\ &= \frac{1}{2} \{ \text{diag}[\Lambda\mathbf{A}\mathbf{A}^\top] - \frac{1}{\mathbf{v}} \}\end{aligned}\quad (13)$$

The hessian in  $\mathbf{v}$  is a matrix in this case:

$$\mathbf{H}_{\mathbf{v}}\mathcal{D} = \frac{1}{2} \text{diag} \left[ \frac{1}{v^2} \right]. \quad (14)$$

This parameterization is useful for spatiotemporal inference problems, where the matrix  $\mathbf{A}$  represents a fixed convolution which can be evaluated using the Fast Fourier Transform (FFT).

### 0.4 Inverse-diagonal approximation

Let  $\Sigma^{-1} = \Lambda + \text{diag}[\mathbf{p}]$ . To obtain the gradient in  $\mathbf{p}$ , combine the derivatives  $\partial_{\Sigma}\mathcal{D}$  (Eq. (2)) and  $\partial_{\mathbf{p}}\Sigma$  using the chain rule. If  $\{(\Sigma)\}$  is a function of  $\Sigma$ , and  $\Sigma(\theta_i)$  is a function of a parameter  $\theta_i$ , then the chain rule is (The Matrix Cookbook; Eq. 136):

$$\partial_{\theta_i}\{ = \langle\partial_{\Sigma}\{, \partial_{\theta_i}\Sigma\rangle = \sum_{kl} (\partial_{\Sigma_{kl}}\{) (\partial_{\theta_i}\Sigma_{kl}) \quad (15)$$

From (2) we have  $\partial_{\Sigma}\mathcal{D} = \frac{1}{2} (\Lambda - \Sigma^{-1})$ ; Since  $\Sigma^{-1} = \Lambda + \text{diag}[\mathbf{p}]$ , this simplifies to:

$$\begin{aligned}\partial_{\Sigma}\mathcal{D} &= \frac{1}{2} (\Lambda - \Sigma^{-1}) \\ &= \frac{1}{2} (\Lambda - \Lambda - \text{diag}[\mathbf{p}]) \\ &= -\frac{1}{2} \text{diag}[\mathbf{p}]\end{aligned}\quad (16)$$

We also need  $\partial_{\mathbf{p}_i}\Sigma$ . Let  $\mathbf{Y} = \Sigma^{-1}$ . The derivative  $\partial\mathbf{Y}^{-1}$  is given as identity (59) in The Matrix Cookbook as  $\partial\mathbf{Y}^{-1} = -\mathbf{Y}^{-1}(\partial\mathbf{Y})\mathbf{Y}^{-1}$ . Using this, we can obtain  $\partial_{\mathbf{p}_i}\Sigma$ :

$$\begin{aligned}\partial_{\mathbf{p}_i}\Sigma &= \partial_{\mathbf{p}_i}\mathbf{Y}^{-1} = -\mathbf{Y}^{-1} (\partial_{\mathbf{p}_i}\mathbf{Y}) \mathbf{Y}^{-1} = -\Sigma (\partial_{\mathbf{p}_i}\Sigma^{-1}) \Sigma \\ &= -\Sigma\partial_{\mathbf{p}_i} [\Lambda + \text{diag}[\mathbf{p}_i]] \Sigma = -\Sigma\mathbf{J}_{ii}\Sigma \\ &= -\sigma_i\sigma_i^\top\end{aligned}\quad (17)$$

where  $\sigma_i$  is the  $i^{\text{th}}$  row of  $\Sigma$  and  $\mathbf{J}_{ii}$  is a matrix which is zero everywhere, except for at index  $(i, i)$ , where it is 1.

Applying (15) to (16) and (17) for a particular element  $\mathbf{p}_i$  gives:

$$\begin{aligned}\partial_{\mathbf{p}_i} \mathcal{D} &= \sum_{kl} [\partial_{\Sigma_{kl}} \mathcal{D}] [\partial_{\mathbf{p}_i} \Sigma_{kl}] = \sum_{kl} \left\{ -\frac{1}{2} \text{diag} [\mathbf{p}] \right\}_{kl} \left\{ -\sigma_i \sigma_i^\top \right\}_{kl} \\ &= \frac{1}{2} \sum_{kl} \delta_{k=l} \mathbf{p}_k \sigma_{ik} \sigma_{il} = \frac{1}{2} \sum_k \mathbf{p}_k \sigma_{ik} \sigma_{ik} = \frac{1}{2} \sum_k \mathbf{p}_k \sigma_{ik}^2 \\ &= \frac{1}{2} \mathbf{p} \sigma_i^{\circ 2}\end{aligned}\tag{18}$$

where  $(\cdot)^{\circ 2}$  denotes the element-wise square of a vector or matrix. In matrix notation, this is:

$$\partial_{\mathbf{p}} \mathcal{D} = \frac{1}{2} \mathbf{p} \Sigma^{\circ 2} = \frac{1}{2} \text{diag} [\Sigma \text{diag} [\mathbf{p}] \Sigma],\tag{19}$$

The Hessian-vector product is cumbersome, since each term in the expression  $\Sigma (\text{diag} [\mathbf{p}]) \Sigma$  depends on  $\mathbf{p}$ . In the case of the log-linear Poisson GLM, the gradient (??) simplifies further and optimization becomes tractable. We will explore this further in later notes.

This parameterization resembles the closed-form covariance update for a linear, Gaussian model, where  $1/\mathbf{p}$  is a vector of measurement noise variances. It is also a useful parameterization for variational Bayesian solutions for non-conjugate Generalized Linear Models (GLMs), where  $\mathbf{p}$  becomes a free parameter to be estimated.

## 0.5 $\Sigma = \mathbf{F}^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{F}$

Let  $\Sigma = \mathbf{F}^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{F}$ , where  $\mathbf{Q} \in \mathbb{R}^{K \times K}$ ;  $K < L$  is the free parameter and  $\mathbf{F} \in \mathbb{R}^{K \times L}$  is a fixed transformation. If  $\mathbf{Q}$  is a lower-triangular matrix, then this approximation involves optimizing  $K(K+1)/2$  parameters.

Since the trace is invariant under cyclic permutation,  $\text{tr} [\Lambda \mathbf{F}^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{F}] = \text{tr} [\mathbf{F} \Lambda \mathbf{F}^\top \mathbf{Q} \mathbf{Q}^\top]$ . The derivatives have the same form as (12) with  $\tilde{\Lambda} = \mathbf{F} \Lambda \mathbf{F}^\top$ :

$$\begin{aligned}\partial_{\mathbf{Q}} \mathcal{D} &= \tilde{\Lambda} \mathbf{Q} - \mathbf{Q}^{-\top} \\ &= \mathbf{F} \Lambda \mathbf{F}^\top \mathbf{Q} - \mathbf{Q}^{-\top} \\ \partial_{\mathbf{Q}} \langle \partial_{\mathbf{Q}} \mathcal{D}, \mathbf{M} \rangle &= \tilde{\Lambda} \mathbf{M}^\top + \mathbf{Q}^{-\top} \mathbf{M}^\top \mathbf{Q}^{-\top} \\ &= \mathbf{F} \Lambda \mathbf{F}^\top \mathbf{M}^\top + \mathbf{Q}^{-\top} \mathbf{M}^\top \mathbf{Q}^{-\top}\end{aligned}\tag{20}$$

This form is convenient for spatiotemporal inference problems that are sparse in frequency space. In this application,  $\mathbf{F}$  corresponds a (unitary) Fourier transform with all by  $K$  of the resulting frequency components discarded. The product of  $\mathbf{F}$  with a vector  $\mathbf{v}$  can be computed in  $O[L \log(L)]$  time using the Fast Fourier Transform (FFT). Alternatively, if  $K \leq O(\log(L))$ , it is faster to simply multiply  $\mathbf{F} \mathbf{v}$  directly. Furthermore, if  $\mathbf{F}$  is semi-orthogonal ( $\mathbf{F} \mathbf{F}^\top = \mathbf{I}$ ), then calculation of  $\mathbf{F}^\top \mathbf{Q}$  can be re-used (for example  $\text{diag}[\Sigma] = [(\mathbf{F}^\top \mathbf{Q})^{\circ 2}]^\top \mathbf{1}$ ).

## 0.6 Conclusion

These notes provide the gradients and Hessian-vector products for four simplified parameterizations of the posterior covariance matrix for variational Gaussian process inference. If combined with the gradients

and Hessian-vector products for the expected log-likelihood, these expressions can be used with Krylov-subspace solvers to compute the Newton-Raphson update to optimize  $\Sigma$ .

We evaluated the following parameterizations for  $\Sigma$ : 1.  $\Sigma$ :

$$\begin{aligned}\partial &= \frac{1}{2} (\Lambda - \Sigma^{-1}) \\ \partial \langle \partial, \mathbf{M} \rangle &= \frac{1}{2} \Sigma^{-1} \mathbf{M}^\top \Sigma^{-1}\end{aligned}\tag{21}$$

2.  $\Sigma \approx \mathbf{X}\mathbf{X}^\top$ :

$$\begin{aligned}\partial &= \Lambda \mathbf{X} - \mathbf{X}^{+\top} \\ \partial \langle \partial, \mathbf{M} \rangle &= \Lambda \mathbf{M}^\top + \mathbf{X}^{+\top} \mathbf{M}^\top \mathbf{X}^{+\top} - (\mathbf{I} - \mathbf{X}^{+\top} \mathbf{X}^\top) \mathbf{M} \mathbf{X}^{+\top}\end{aligned}\tag{22}$$

3.  $\Sigma \approx \mathbf{A}^\top \text{diag}[\mathbf{v}] \mathbf{A}$ :

$$\begin{aligned}\partial &= \frac{1}{2} \left\{ \text{diag}[\mathbf{A} \Lambda \mathbf{A}^\top] - \frac{1}{\mathbf{v}} \right\} \\ \partial \langle \partial, \mathbf{u} \rangle &= \frac{1}{2} \left[ \frac{1}{\mathbf{v}^2} \right]^\top \mathbf{u}\end{aligned}\tag{23}$$

4.  $\Sigma \approx [\Lambda + \text{diag}[\mathbf{p}]]^{-1}$ :

$$\partial = \frac{1}{2} \mathbf{p} \Sigma^{\circ 2} = \frac{1}{2} \text{diag} [\Sigma \text{diag} [\mathbf{p}] \Sigma],\tag{24}$$

5.  $\mathbf{F}^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{F}$ :

$$\begin{aligned}\partial &= \mathbf{F} \Lambda \mathbf{F}^\top \mathbf{Q} - \mathbf{Q}^{-\top} \\ \partial \langle \partial, \mathbf{M} \rangle &= \mathbf{F} \Lambda \mathbf{F}^\top \mathbf{M}^\top + \mathbf{Q}^{-\top} \mathbf{M}^\top \mathbf{Q}^{-\top}\end{aligned}\tag{25}$$

In future notes, we will consider the full derivatives required for variational latent Gaussian-process inference for the Poisson and probit generalized linear models.