

Simple Review of Gaussian Variational Approximate Inference

Daeyoung Lim*
Department of Statistics
Korea University

January 24, 2016

1 Figure 2

This section reviews the figure plotted in p.13 of the associated paper.

Variational Coefficients					
β_0	β_1	β_2	β_3	β_4	β_5
-1.325	0.883	-0.933	0.481	-0.160	0.339

GVA Standard Errors					
β_0	β_1	β_2	β_3	β_4	β_5
1.179	0.131	0.400	0.346	0.055	0.203

- Coefficient estimates via GVA seem to align well with the plot given in the paper.
- The caption under the *confidence interval* plots should be a misnomer since interval estimation is called the *credible interval* in Bayesian frameworks. Requires revision.
- The credible intervals plotted in box and whisker plots are consistent with the numerical results calculated with the results above.
- *Adaptive Gauss-Hermite quadrature* and *penalized quasi-likelihood* methods were unavailable.

2 Figure 3

This section reviews the code for figure 3 in p.14.

- The code returned an error that the number of *adaptive Gauss-Hermite quadrature* greater than 1 is only available for models with a single, scalar random-effects term.
- The line of code `glmer(formula=y ~ Base*trt+Age+Visit+(Visit|subject), data=epil2, family = poisson, nAGQ = 20)` is the source of this error.

*Prof. Taeryon Choi

3 Figure 4

This section reviews the code for figure 4 in p.15.

Variational Coefficients			
β_0	β_1	β_2	β_3
-1.434	-0.134	-0.377	-0.130

GVA Standard Errors			
β_0	β_1	β_2	β_3
0.386	0.526	0.043	0.065

- R package {glmmAK} is deprecated under R version 3.2.3 (2015-12-10) – “Wooden Christmas-Tree”.
- The required dataset ‘toenail’ in R package {glmmAK} was downloaded from <https://github.com/cran/glmmAK>
- The results all match what is written in the paper.

4 Model review

Model Specifications	
Likelihood	$\mathbf{y}_i \mathbf{u}_i \sim \exp \left\{ \mathbf{y}_i^T (X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) + \mathbf{1}_i^T c(\mathbf{y}_i) \right\}$
Prior	$\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$

The log-marginal likelihood $p(\mathbf{y})$ is, with the abuse of notation, $\int \ell(\mathbf{u}, \boldsymbol{\beta}, \Sigma) d\mathbf{u}$ where

$$\begin{aligned} \ell(\mathbf{u}, \boldsymbol{\beta}, \Sigma) &= \log \left(\prod_{i=1}^m \exp \left\{ \mathbf{y}_i^T (X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) + \mathbf{1}_i^T c(\mathbf{y}_i) \right\} \cdot \frac{1}{\sqrt{|2\pi\Sigma|}} e^{-\frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i} \right) \\ &= \sum_{i=1}^m \left[\mathbf{y}_i^T (X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) + \mathbf{1}_i^T c(\mathbf{y}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i - \frac{1}{2} \log |2\pi\Sigma| \right] \end{aligned}$$

Therefore,

$$\begin{aligned} \ell(\boldsymbol{\beta}, \Sigma) &= \sum_{i=1}^m \left[\mathbf{y}_i^T X_i \boldsymbol{\beta} + \mathbf{1}_i^T c(\mathbf{y}_i) - \frac{1}{2} \log |2\pi\Sigma| \right] \\ &\quad + \sum_{i=1}^m \log \int_{\mathbb{R}^k} \exp \left\{ \mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i \right\} d\mathbf{u} \end{aligned}$$

The integral does not change even if we multiply the integrand by 1:

$$\int_{\mathbb{R}^k} \exp \left\{ \mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i \right\} \frac{\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)}{\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)} d\mathbf{u} \quad (1)$$

where $\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)$ is the multivariate Gaussian p.d.f. of a random vector \mathbf{u} with a mean vector μ_i and a covariance matrix Λ_i . Then (1) is rewritten in terms of the expectation with respect to \mathbf{u} :

$$(1) = \mathbb{E}_{\mathbf{u}} \left[\frac{\exp \left\{ \mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i \right\}}{\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)} \right]$$

According to *Jensen's inequality*, $\log(\mathbb{E}(X)) \geq \mathbb{E}(\log X)$, which thus yields

$$\begin{aligned} & \log \left(\mathbb{E}_{\mathbf{u}} \left[\frac{\exp \left\{ \mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i \right\}}{\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)} \right] \right) \\ & \geq \mathbb{E}_{\mathbf{u}} \left[\mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i - \log(\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)) \right] \end{aligned}$$

The inequality gives us the lower bound of the variational approximation for the generalized linear mixed model.

$$\begin{aligned} \underline{\ell}(\boldsymbol{\beta}, \Sigma, \boldsymbol{\mu}, \boldsymbol{\Lambda}) &= \sum_{i=1}^m \left[\mathbf{y}_i^T X_i \boldsymbol{\beta} + \mathbf{1}_i^T c(\mathbf{y}_i) - \frac{1}{2} \log |2\pi \Sigma| \right] \\ &+ \sum_{i=1}^m \mathbb{E}_{\mathbf{u}} \left[\mathbf{y}_i^T Z_i \mathbf{u}_i - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \frac{1}{2} \mathbf{u}_i^T \Sigma^{-1} \mathbf{u}_i - \log(\varphi_{\Lambda_i}(\mathbf{u} - \mu_i)) \right] \end{aligned}$$

Computing the expectation is not hard except for the unspecified function $b(\cdot)$ which can be easily approximated via *Gauss-Hermite quadrature*. Therefore, we have a scalar function as the objective function and are only left with optimizing the function over variational parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Lambda}$ along with the original parameters $\boldsymbol{\beta}$ and Σ . The computational formula for *Newton-Raphson* algorithm is given in the supplementary document. However, the real problem is making sense of why the formula works and how the original authors of the paper computed the necessary elements such as the Hessian matrix and gradient vector for each parameter.

Stuff yet to be resolved

- How to compute the abscissa of Hermite polynomials that are used in *Gauss-Hermite quadrature*.
- How to determine the number of nodes of *Gauss-Hermite quadrature*.
- How to compute the Hessian of $\text{vech}(\Sigma)$ or other terms with $\text{vech}(\cdot)$.
- Why the *Newton-Raphson* formula is written as what is given in the supplementary paper.

5 Examples

We will from now on follow the examples and try to calculate the gradient vectors and Hessian matrices whose general forms are given in the supplementary material. To achieve this, we first need to compute $\mathcal{B}^{(1)}(\boldsymbol{\beta}, \boldsymbol{\mu}_i, \boldsymbol{\Lambda}_i)$.

5.1 Poisson Mixed

Specifying the parametric family of the GLM adds clarity to how we should compute the derivatives necessary for Newton's method. As a gentle reminder, the general model was

$$\mathbf{y}_i | \mathbf{u}_i \sim \exp \left\{ \mathbf{y}_i^T (X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) - \mathbf{1}_i^T b(X_i \boldsymbol{\beta} + Z_i \mathbf{u}_i) + \mathbf{1}_i^T c(\mathbf{y}_i) \right\}.$$

In Poisson mixed model, $b(x) = e^x$ and $c(x) = -\log(x!)$.

$$B(\mu, \sigma^2) = \exp \left(\mu + \frac{1}{2} \sigma^2 \right)$$

$$\mathcal{B}(\boldsymbol{\beta}, \boldsymbol{\mu}_i, \boldsymbol{\Lambda}_i) = B(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\mu}_i, \text{dg}(\mathbf{Z}_i \boldsymbol{\Lambda}_i \mathbf{Z}_i^T))$$

$\frac{\partial \mathcal{B}}{\partial \boldsymbol{\beta}}$ is our goal. For notational convenience, let's replace \mathcal{B} with simple f .

$$df = d \left(\exp(\mathbf{X}_i \boldsymbol{\beta}) \circ \exp \left\{ \mathbf{Z}_i \boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i \boldsymbol{\Lambda}_i \mathbf{Z}_i^T) \right\} \right)$$

$$= d \exp(\mathbf{X}_i \boldsymbol{\beta}) \circ \exp \left\{ \mathbf{Z}_i \boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i \boldsymbol{\Lambda}_i \mathbf{Z}_i^T) \right\}$$

where \circ indicates *Hadamard product*(or *Schur product*) which is simply the elementwise product between two matrices(or vectors) of the same size. We differentiate the notation between **Diag** and **dg**. The former is applied to **vectors** and converts it into a diagonal matrix whose diagonal elements correspond to the elements of the vector. On the other hand, **dg** is applied to **matrices** and collects the diagonal entries and creates a column vector with them.

5.1.1 $D_{\boldsymbol{\beta}} \underline{\ell}$

Let

$$s = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i \boldsymbol{\Lambda}_i \mathbf{Z}_i^T)$$

$$e = \exp(s)$$

$$f = \mathbf{1}^T e$$

Then,

$$ds = \mathbf{X}_i d\boldsymbol{\beta} + \mathbf{Z}_i d\boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i d\boldsymbol{\Lambda}_i \mathbf{Z}_i^T).$$

But since we are only interested in $d\boldsymbol{\beta}$, we set $d\boldsymbol{\mu}_i = 0$, $d\boldsymbol{\Lambda}_i = 0$.

$$de = d \exp(s)$$

$$= e \circ ds$$

$$= e \circ \mathbf{X}_i d\boldsymbol{\beta}$$

Therefore, $df = \mathbf{1}^T de = e^T \mathbf{X}_i d\boldsymbol{\beta}$.

$$\frac{\partial \underline{\ell}}{\partial \boldsymbol{\beta}} = \sum_{i=1}^m (\mathbf{y}_i^T - e^T) \mathbf{X}_i$$

5.1.2 $D_{\text{vech}}(\Sigma) \underline{\ell}$

Note that $A : B$ where A and B are matrices of the same size, is the *Frobenius product*. Then we can rewrite $d\underline{\ell}$ as follows:

$$d\underline{\ell} = \frac{1}{2} \sum_{i=1}^m [d(\log \det(\Sigma^{-1} \Lambda_i)) - d(\mu_i \mu_i^T : \Sigma^{-1}) - d(\Lambda_i : \Sigma^{-1})] \quad (2)$$

$$= \frac{1}{2} \sum_{i=1}^m [d(\text{Tr} \log(\Sigma^{-1} \Lambda_i)) - d(\mu_i \mu_i^T : \Sigma^{-1}) - d(\Lambda_i : \Sigma^{-1})] \quad (3)$$

$$= \frac{1}{2} \sum_{i=1}^m [(\Sigma^{-1} \Lambda_i)^{-T} : (d\Sigma^{-1} \Lambda_i) - \mu_i \mu_i^T : d\Sigma^{-1} - \Lambda_i : d\Sigma^{-1}] \quad (4)$$

$$= \frac{1}{2} \sum_{i=1}^m [(\Sigma^T \Lambda_i^{-T}) \Lambda_i^T - \mu_i \mu_i^T - \Lambda_i] : d\Sigma^{-1} \quad (5)$$

$$= \frac{1}{2} \sum_{i=1}^m [-\Sigma^T + \mu_i \mu_i^T + \Lambda_i] : \Sigma^{-1} d\Sigma \Sigma^{-1} \quad (6)$$

$$= \frac{1}{2} \sum_{i=1}^m [\Sigma^{-T} (\mu_i \mu_i^T + \Lambda_i) \Sigma^{-T} - \Sigma^{-T}] : d\Sigma \quad (7)$$

- From equations (1) to (2), $\log \det(A) = \text{Tr} \log(A)$ where $\log(A)$ is NOT the elementwise logarithm but rather the matrix logarithm defined in terms of the Taylor series of e^A .
- Let $f(X) = \log \det(X)$ where X is a positive definite matrix.

$$\frac{\partial}{\partial X} f(X) = X^{-1}.$$

This should not be surprising since $(\log x)' = 1/x$. Positive-definiteness is important since there should be a guarantee that there is no eigenvalue whose value is 0. (The determinant should not be 0 since $\log(0)$ is not defined.)

- In relation to the previous point,

$$\frac{\partial}{\partial X} \text{Tr} f(X) = \left(\frac{d}{dx} f(x) \Big|_{x=X} \right)^T$$

We can rewrite $\log \det(X) = \text{Tr} \log(X)$ whose derivatives are consistent with each other provided that X is positive definite.

- It is a well-proven fact that the differential of the inverse of a matrix is as follows:

$$d\Sigma^{-1} = -\Sigma^{-1} d\Sigma \Sigma^{-1}.$$

- The *Frobenius product* $A : B$ can be described differently (in reference to the transition from eqn (4) to (5)):

$$A : B = \text{Tr}(B^T A).$$

Since trace operator is not affected by transposition,

$$A : BC = \text{Tr}(C^T B^T A) = \text{Tr}(A^T BC) = B^T A : C.$$

- Another variant of the *Frobenius product* is

$$A : B = \text{vec}(A)^T \text{vec}(B).$$

Using $\Sigma^T = \Sigma$ and $\text{vec}(X) = D_p \text{vech}(X)$, we can recast the last line as follows:

$$\begin{aligned} d\ell &= \frac{1}{2} \sum_{i=1}^m \text{vec} \left[\Sigma^{-1} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \Lambda_i) \Sigma^{-1} - \Sigma^{-1} \right]^T (D_K \text{vech}(d\Sigma)) \\ \therefore \frac{\partial \ell}{\partial \text{vech}(\Sigma)} &= \frac{1}{2} \sum_{i=1}^m \text{vec} \left\{ \Sigma^{-1} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \Lambda_i) \Sigma^{-1} - \Sigma^{-1} \right\}^T D_K \end{aligned}$$

5.1.3 $D_{\boldsymbol{\mu}_i} \underline{\ell}$

This is similar to section 5.1.1 $D_{\boldsymbol{\beta}} \underline{\ell}$.

$$\begin{aligned} \frac{\partial f}{\partial \boldsymbol{\mu}_i} &= e^T \mathbf{Z}_i d\boldsymbol{\mu}_i \\ \therefore \frac{\partial \underline{\ell}}{\partial \boldsymbol{\mu}_i} &= (y_i - e)^T \mathbf{Z}_i - \boldsymbol{\mu}_i^T \Sigma^{-1} \end{aligned}$$

5.1.4 $D_{\text{vech}(\Lambda_i)} \underline{\ell}$

Collecting just the terms that contain Λ_i ,

$$d\underline{\ell} = -d \left(\mathbf{1}_i^T \exp \left(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i \Lambda_i \mathbf{Z}_i^T) \right) \right) + \frac{1}{2} \left\{ d \log \det(\Sigma^{-1} \Lambda_i) - d(\Sigma^{-1} : \Lambda_i) \right\}$$

As per the first term, for convenience, let's write

$$\begin{aligned} s &= \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i \Lambda_i \mathbf{Z}_i) \\ e &= \exp(s) \\ f &= \mathbf{1}^T e \end{aligned}$$

Then,

$$\begin{aligned} ds &= \mathbf{X}_i d\boldsymbol{\beta} + \mathbf{Z}_i d\boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i d\Lambda_i \mathbf{Z}_i^T) \\ de &= e \circ ds \\ df &= \mathbf{1}^T de \\ &= \mathbf{1}^T (e \circ ds) \\ &= e^T ds \\ &= e^T \left(\mathbf{X}_i d\boldsymbol{\beta} + \mathbf{Z}_i d\boldsymbol{\mu}_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i d\Lambda_i \mathbf{Z}_i^T) \right) \end{aligned}$$

Since we are working on the gradient with respect to $\mathbf{\Lambda}_i$, set $d\boldsymbol{\beta} = 0$ and $d\boldsymbol{\mu}_i = 0$.

$$\begin{aligned}
df &= \frac{1}{2} e^T \text{dg} (\mathbf{Z}_i d\mathbf{\Lambda}_i \mathbf{Z}_i^T) \\
&= \frac{1}{2} \text{Diag} (e^T) : \mathbf{Z}_i d\mathbf{\Lambda}_i \mathbf{Z}_i^T \\
&= \frac{1}{2} E : \mathbf{Z}_i d\mathbf{\Lambda}_i \mathbf{Z}_i^T \\
&= \frac{1}{2} \mathbf{Z}_i^T E \mathbf{Z}_i : d\mathbf{\Lambda}_i
\end{aligned}$$

The second term is rather easy.

$$\begin{aligned}
d \log \det (\boldsymbol{\Sigma}^{-1} \mathbf{\Lambda}_i) &= d \text{Tr} \log (\boldsymbol{\Sigma}^{-1} \mathbf{\Lambda}_i) \\
&= (\boldsymbol{\Sigma}^{-1} \mathbf{\Lambda}_i)^{-T} : d\boldsymbol{\Sigma}^{-1} \mathbf{\Lambda}_i \\
&= \boldsymbol{\Sigma}^{-T} (\boldsymbol{\Sigma}^{-1} \mathbf{\Lambda}_i)^{-T} : d\mathbf{\Lambda}_i \\
&= \mathbf{\Lambda}_i^{-T} : d\mathbf{\Lambda}_i
\end{aligned}$$

Combining all three terms,

$$\begin{aligned}
d\ell &= -\frac{1}{2} \mathbf{Z}_i^T E \mathbf{Z}_i : d\mathbf{\Lambda}_i + \frac{1}{2} \left[\mathbf{\Lambda}_i^{-T} : d\mathbf{\Lambda}_i - \boldsymbol{\Sigma}^{-1} : d\mathbf{\Lambda}_i \right] \\
&= \frac{1}{2} (\mathbf{\Lambda}_i^{-1} - \boldsymbol{\Sigma}^{-1} - \mathbf{Z}_i^T E \mathbf{Z}_i) : d\mathbf{\Lambda}_i
\end{aligned}$$

Using the vec operator, the differential form translates to

$$\begin{aligned}
d\ell &= \frac{1}{2} \text{vec} (\mathbf{\Lambda}_i^{-1} - \boldsymbol{\Sigma}^{-1} - \mathbf{Z}_i^T E \mathbf{Z}_i)^T D_K \text{vech} (d\mathbf{\Lambda}_i) \\
\therefore \frac{\partial \ell}{\partial \text{vech}(\mathbf{\Lambda}_i)} &= \frac{1}{2} \text{vec} (\mathbf{\Lambda}_i^{-1} - \boldsymbol{\Sigma}^{-1} - \mathbf{Z}_i^T E \mathbf{Z}_i)^T D_K
\end{aligned}$$

5.1.5 $\mathbf{H}_{\boldsymbol{\beta}\boldsymbol{\beta}} \ell$

Recall that

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = \sum_{i=1}^m (\mathbf{y} - e)^T \mathbf{X}_i.$$

We set $f = \frac{\partial \ell}{\partial \boldsymbol{\beta}}$.

$$\begin{aligned}
df &= - \sum_{i=1}^m (de)^T \mathbf{X}_i \\
&= - \sum_{i=1}^m (e \circ \mathbf{X}_i d\boldsymbol{\beta})^T \mathbf{X}_i
\end{aligned}$$

Here we should note that between two vectors x, y ,

$$x \circ y = \text{Diag} (x) y = \text{Diag} (y) x.$$

Therefore,

$$\begin{aligned} df &= - \sum_{i=1}^m (\text{Diag}(e) \mathbf{X}_i d\boldsymbol{\beta})^T \mathbf{X}_i \\ &= - \sum_{i=1}^m (d\boldsymbol{\beta})^T \mathbf{X}_i^T \text{Diag}(e) \mathbf{X}_i \end{aligned}$$

Finally, we can recast the Hessian as

$$\begin{aligned} \frac{\partial^2 \underline{\ell}}{\partial \boldsymbol{\beta}^2} &= - \sum_{i=1}^m \mathbf{X}_i^T \text{Diag}(e) \mathbf{X}_i \\ &= - \sum_{i=1}^m \mathbf{X}_i^T E \mathbf{X}_i. \end{aligned}$$

5.1.6 $H_{\text{vech}(\boldsymbol{\Sigma}) \text{vech}(\boldsymbol{\Sigma}) \underline{\ell}}$

Recall that

$$\frac{\partial \underline{\ell}}{\partial \text{vech}(\boldsymbol{\Sigma})} = \frac{1}{2} \sum_{i=1}^m \text{vec} \{ \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Lambda}_i) \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \}^T D_K.$$

Let $f = \frac{\partial \underline{\ell}}{\partial \text{vech}(\boldsymbol{\Sigma})}$ this time. Before we jump in, let this be a gentle reminder of what we will be using:

- $\text{vec}(\boldsymbol{\Sigma}) = D_K \text{vech}(\boldsymbol{\Sigma})$
- $\partial X^{-1} = -X^{-1} \partial X X^{-1}$
- $\text{vec}(A \pm B) = \text{vec}(A) \pm \text{vec}(B)$
- $\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X)$
- $(A \otimes B)^T = A^T \otimes B^T$
- Lastly, let's write

$$M_i := \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Lambda}_i) \boldsymbol{\Sigma}^{-1} \quad (M_i^T = M_i).$$

Then,

$$\begin{aligned}
df &= \frac{1}{2} \sum_{i=1}^m [\text{vec}(-\Sigma^{-1} d\Sigma M_i - M_i d\Sigma \Sigma^{-1}) + \text{vec}(\Sigma^{-1} d\Sigma \Sigma^{-1})]^T D_K \\
&= \frac{1}{2} \sum_{i=1}^m [-\text{vec}(\Sigma^{-1} d\Sigma M_i) - \text{vec}(M_i d\Sigma \Sigma^{-1}) + \text{vec}(\Sigma^{-1} d\Sigma \Sigma^{-1})]^T D_K \\
&= \frac{1}{2} \sum_{i=1}^m [-(M_i \otimes \Sigma^{-1}) \text{vec}(d\Sigma) - (\Sigma^{-1} \otimes M_i) \text{vec}(d\Sigma) + (\Sigma^{-1} \otimes \Sigma^{-1}) \text{vec}(d\Sigma)]^T D_K \\
&= \frac{1}{2} \sum_{i=1}^m \left\{ -\text{vec}(d\Sigma)^T (M_i \otimes \Sigma^{-1}) - \text{vec}(d\Sigma)^T (\Sigma^{-1} \otimes M_i) + \text{vec}(d\Sigma)^T (\Sigma^{-1} \otimes \Sigma^{-1}) \right\} D_K \\
&= \frac{1}{2} \text{vec}(d\Sigma)^T \sum_{i=1}^m \{ -(M_i \otimes \Sigma^{-1} - \Sigma^{-1} \otimes M_i) + \Sigma^{-1} \otimes \Sigma^{-1} \} \\
&= \frac{1}{2} \text{vech}(d\Sigma)^T D_K^T \sum_{i=1}^m \{ -(M_i \otimes \Sigma^{-1} + \Sigma^{-1} \otimes M_i) + \Sigma^{-1} \otimes \Sigma^{-1} \} D_K \\
&= \text{vech}(d\Sigma)^T \frac{1}{2} D_K^T \left\{ -\sum_{i=1}^m (M_i \otimes \Sigma^{-1} + \Sigma^{-1} \otimes M_i) + m (\Sigma^{-1} \otimes \Sigma^{-1}) \right\} D_K
\end{aligned}$$

Therefore,

$$\frac{\partial^2 \underline{\ell}}{\partial \text{vech}(d\Sigma) \partial \text{vech}(d\Sigma)} = \frac{1}{2} D_K^T \left\{ -\sum_{i=1}^m (M_i \otimes \Sigma^{-1} + \Sigma^{-1} \otimes M_i) + m (\Sigma^{-1} \otimes \Sigma^{-1}) \right\} D_K$$

5.1.7 $H_{\beta \mu_i} \underline{\ell}$

Recall that

$$\frac{\partial \underline{\ell}}{\partial \beta} = \sum_{i=1}^m (\mathbf{y}_i^T - e^T) \mathbf{X}_i.$$

We need to go through the same thing again. Differentiate $\frac{\partial \underline{\ell}}{\partial \beta}$ with respect to μ_i . However, for the time being, we will only consider the i^{th} element of μ , which rules out the need for summation.

$$df = -(de)^T \mathbf{X}_i.$$

$$de = e \circ ds$$

$$ds = \mathbf{X}_i d\beta + \mathbf{Z}_i d\mu_i + \frac{1}{2} \text{dg}(\mathbf{Z}_i d\Lambda_i \mathbf{Z}_i^T)$$

Setting $d\beta = 0$, $d\Lambda_i = 0$,

$$\begin{aligned}
de &= e \circ ds \\
&= e \circ \mathbf{Z}_i d\mu_i \\
&= E \mathbf{Z}_i d\mu_i
\end{aligned}$$

where again $E = \text{Diag}(e)$. (Readers who are not familiar with the notations are referred to section 5.1.1 and 5.1.4.) Therefore,

$$df = -(d\mu_i)^T \mathbf{Z}_i^T E \mathbf{X}_i$$

which is then yields

$$\frac{\partial^2 \ell}{\partial \beta \partial \mu_i} = -\mathbf{Z}_i^T E \mathbf{X}_i$$

5.1.8 $H_{\text{vech}(\mathbf{\Lambda}_i)\mu_i}$

We start from

$$\frac{\partial \ell}{\partial \text{vech}(\mathbf{\Lambda}_i)} = \frac{1}{2} \left(\text{vec} \left(\mathbf{\Lambda}_i^{-1} - \Sigma^{-1} - \mathbf{Z}_i^T E \mathbf{Z}_i \right) \right)^T D_K$$

Now setting $f = \frac{\partial \ell}{\partial \text{vech}(\mathbf{\Lambda}_i)}$,

$$\begin{aligned} df &= -\frac{1}{2} \left(\text{vec} \left(\mathbf{Z}_i^T E \mathbf{Z}_i \right) \right)^T D_K \\ dE &= \text{Diag}(de). \end{aligned}$$

We now use the following theorem:

$$\text{vec} \left(\mathbf{Z}^T \text{Diag}(de) \mathbf{Z} \right) = \mathcal{Q}(\mathbf{Z})^T de$$

where $\mathcal{Q}(\mathbf{Z}) = (\mathbf{Z} \otimes \mathbf{1}^T) \circ (\mathbf{1}^T \otimes \mathbf{Z})$. $\mathbf{1}$ is a column vector filled with ones of the same dimension as the columns of \mathbf{Z} . Then,

$$\begin{aligned} df &= -\frac{1}{2} \left(\text{vec} \left(\mathbf{Z}_i^T \text{Diag}(de) \mathbf{Z}_i \right) \right)^T D_K \\ &= -\frac{1}{2} (de)^T \mathcal{Q}(\mathbf{Z}_i) D_K \\ &= -\frac{1}{2} (e \circ ds)^T \mathcal{Q}(\mathbf{Z}_i) D_K \\ &= -\frac{1}{2} (\text{Diag}(e) \mathbf{Z}_i d\mu_i)^T \mathcal{Q}(\mathbf{Z}_i) D_K \\ &= -\frac{1}{2} (d\mu_i)^T \mathbf{Z}_i^T E \mathcal{Q}(\mathbf{Z}_i) D_K \end{aligned}$$

Thus,

$$\frac{\partial^2 \ell}{\partial \text{vech}(\mathbf{\Lambda}_i) \partial \mu_i} = -\frac{1}{2} \mathbf{Z}_i^T E \mathcal{Q}(\mathbf{Z}_i) D_K.$$