

Derivation of fisher information of a multivariate Gaussian density

Kevin W Li

March 7, 2017

We are interested to find the fisher information of this density with respect to parameter θ upon which μ and Σ depends. This notes assumes fluency in matrix calculus (denominator layout) in the matrix cookbook or wikipedia.

Preliminaries

For scalar x , scalar function f , vectors a, b, c , and matrices A, B, C :

$$\begin{aligned}d \frac{a^\top Ab}{da} &= Ab \\d \frac{a^\top Ab}{dA} &= d \frac{\text{Tr}(a^\top Ab)}{dA} = d \frac{\text{Tr}(Aa^\top b)}{dA} = ab^\top \\d \frac{f(b(x))}{dx} &= \left[\frac{df(b)}{db} \right]^\top \nabla_x b \\d \frac{f(A(x))}{dx} &= \text{Tr} \left[\left(\frac{dx}{dA} \right)^\top \nabla_x A \right] \\d \frac{a(c)^\top b}{dc} &= \frac{da^\top}{dc} b \\d \frac{\log A}{dx} &= A^{-T} \nabla_x A \\dA^{-1} &= -A^{-1} dA A^{-1} \\d \frac{A^{-1}}{dx} &= -A^{-1} \nabla_x A A^{-1} \\d \frac{\text{Tr}(AB^{-1}C)}{dB} &= -B^{-T} C A B^{-T}\end{aligned}$$

In addition, the first derivative of a vector w.r.t a scalar is a vector, first derivative of a matrix w.r.t a scalar is a matrix.

Derivation

The multivariate Gaussian density is written as:

$$p(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu) \right) = \frac{|\Lambda|^{1/2}}{(2\pi)^{D/2}} \exp \left(-\frac{1}{2}(x - \mu)^\top \Lambda(x - \mu) \right)$$

where Λ is the precision matrix and is symmetric. Define

$$\begin{aligned}\nabla\mu &= \frac{d\mu}{ds} \\ \nabla\Sigma &= \frac{d\Sigma}{ds} \\ \nabla\Lambda &= \frac{d\Lambda}{ds} = \frac{d\Sigma^{-1}}{ds} = -\Sigma^{-1}\nabla\Sigma\Sigma^{-1}\end{aligned}$$

In the following, we will use the fact that $\Lambda = \Lambda^\top$

And write the log density as:

$$\log p(x) = \frac{1}{2} \log |\Lambda| - \frac{1}{2} (x - \mu)^\top \Lambda (x - \mu) + \text{constant}$$

There two ways ahead. One can write out the second derivative of the log density using the chain rule. It will be composed of four terms in which two of them will involve a score term $\frac{d\log p(x)}{d\mu}$ and $\frac{d\log p(x)}{d\Lambda}$. The expectation of these two terms will be zeros, so we ended up having only two terms. However, $\frac{d^2 \log p(x)}{d\mu^2}$ is a matrix and $\frac{d^2 \log p(x)}{d\Lambda^2}$ is a four-dimensional tensor which we should avoid... It is still possible to work, but less clear.

The second way is to compute the first derivative first in terms of the parameters, and then take derivative again.

$$\begin{aligned}\frac{d\log p(x)}{ds} &= \left[\frac{d\log p(s)}{d\mu} \right]^\top \nabla\mu + \text{Tr} \left[\frac{d\log p(s)}{d\Lambda} \nabla\Lambda \right] \\ &= \underbrace{(x - \mu)^\top \Lambda \nabla\mu}_{M_1} + \frac{1}{2} \underbrace{\text{Tr} \left[[\Lambda^{-1} - (x - \mu)(x - \mu)^\top] \nabla\Lambda \right]}_{M_2}\end{aligned}$$

Taking the square yields something hard to proceed. Let's continue with the other equivalent definition of FI that uses second derivatives.

Take the derivative of M_1 and then evaluate the negative expectation

$$\begin{aligned}\frac{dM_1}{ds} &= \left[\frac{dM_1}{d\mu} \right]^\top \nabla\mu + \text{Tr} \left[\left[\frac{dM_1}{d\Lambda} \right]^\top \nabla\Lambda \right] \\ &= [-\Lambda \nabla\mu + \nabla^2 \mu \Lambda (x - \mu)]^\top \nabla\mu + \text{Tr} [(x - \mu) \nabla\mu^\top]^\top \nabla\Lambda \\ - \left\langle \frac{dM_1}{ds} \right\rangle &= \nabla\mu^\top \Lambda \nabla\mu = \nabla\mu^\top \Sigma^{-1} \nabla\mu\end{aligned}$$

Do the same for M_2

$$\begin{aligned}\frac{dM_2}{ds} &= \left[\frac{dM_2}{d\mu} \right]^\top \nabla\mu + \text{Tr} \left[\left[\frac{dM_2}{d\Lambda} \right]^\top \nabla\Lambda \right] \\ &= - \left[\frac{d}{d\mu} \text{Tr} [(x - \mu)(x - \mu)^\top \nabla\Lambda] \right]^\top \nabla\mu + \text{Tr} \left[\left[\frac{d}{d\Lambda} \text{Tr} [\Lambda^{-1} \nabla\Lambda] - \frac{d}{d\Lambda} \text{Tr} [(x - \mu)(x - \mu)^\top \nabla\Lambda] \right]^\top \nabla\Lambda \right] \\ &= -2(x - \mu)^\top \nabla\Lambda \nabla\mu + \text{Tr} \left[[-\Lambda^{-1} \nabla\Lambda \Lambda^{-1} + \Lambda^{-1} \nabla^2 \Lambda - (x - \mu)(x - \mu)^\top \nabla^2 \Lambda]^\top \nabla\Lambda \right] \\ &= -2(x - \mu)^\top \nabla\Lambda \nabla\mu + \text{Tr} [-\Lambda^{-1} \nabla\Lambda \Lambda^{-1} \nabla\Lambda] + \text{Tr} [\Lambda^{-1} - (x - \mu)(x - \mu)^\top] \nabla^2 \Lambda \nabla\Lambda \\ - \left\langle \frac{dM_2}{ds} \right\rangle &= 0 + \text{Tr} [-\Lambda^{-1} \nabla\Lambda \Lambda^{-1} \nabla\Lambda] + \text{Tr} [\Lambda^{-1} \nabla^2 \Lambda - \Sigma \nabla^2 \Lambda] \nabla\Lambda \\ - \left\langle \frac{dM_2}{ds} \right\rangle &= \text{Tr} [\Lambda^{-1} \nabla\Lambda \Lambda^{-1} \nabla\Lambda] = \text{Tr} [\Lambda^{-1} \Sigma^{-1} \nabla\Sigma \Sigma^{-1} \Lambda^{-1} \Sigma^{-1} \nabla\Sigma \Sigma^{-1}] = \text{Tr} [\nabla\Sigma \Sigma^{-1} \nabla\Sigma \Sigma^{-1}]\end{aligned}$$

The last equality is because the first contains first moment, and the third term is zero even.

So the fisher information is then

$$J(s) = \nabla \mu^\top \Sigma^{-1} \nabla \mu + \frac{1}{2} \text{Tr} [\nabla \Sigma \Sigma^{-1} \nabla \Sigma \Sigma^{-1}]$$