



## Estimation de la variance généralisée

Thu Pham-Gia

► **To cite this version:**

Thu Pham-Gia. Estimation de la variance généralisée. 42èmes Journées de Statistique, 2010, Marseille, France, France. 2010. <inria-00494686>

**HAL Id: inria-00494686**

**<https://hal.inria.fr/inria-00494686>**

Submitted on 24 Jun 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## ESTIMATION OF THE GENERALIZED VARIANCE IN MULTIVARIATE ANALYSIS

by Thu PHAM-GIA, (phamgit@umoncton.ca)

*Département de math. et de stat., Université de Moncton, Canada E1A 3E9.*

**ABSTRACT:** Expressing the density of the determinant of a sample covariance matrix in terms of Meijer's G-function, we provide a confidence interval for the determinant if the population covariance matrix. The Bayesian approach is used to provide a set of bounds for the posterior generalized variance of the mean vector and other cases. Numerical applications of the results are also presented.

**RESUME:** Exprimant la densité du déterminant de la matrice de covariance d'échantillon sous forme de fonction de Meijer, nous dérivons un intervalle de confiance de la variance généralisée de la population. Ceci permet, à travers l'approche bayésienne, d'obtenir des bornes pour la variance généralisée postérieure du vecteur moyen. Des applications numériques sont aussi présentées.

**MOTS-CLES :** Statistique mathématique, Méthodes bayésiennes.

### Main results

We consider  $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and a random sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$  of  $\mathbf{X}$ . We define

$$\mathbf{A} = \sum_{\alpha=1}^N (\mathbf{X}_\alpha - \bar{\mathbf{X}})(\mathbf{X}_\alpha - \bar{\mathbf{X}})^T, \quad (1)$$

with  $\bar{\mathbf{X}} = \sum_{\alpha=1}^N \mathbf{X}_\alpha / N$ , the sample sum of squares and products matrix. The sample covariance matrix is:

$\mathbf{S} = \mathbf{A} / n$ . For  $n > p$  we have  $\mathbf{A} \sim W_p(n, \boldsymbol{\Sigma})$ , where  $W_p(n, \boldsymbol{\Sigma})$  is the central Wishart distribution in  $R^p$ , with  $n$  degrees of freedom and covariance matrix  $\boldsymbol{\Sigma}$ , which is a multivariate generalization of  $\chi_n^2$ . Hence,  $\mathbf{S} \sim W_p(n, \boldsymbol{\Sigma} / n)$ . Also, it can be proven that :

$$|\mathbf{A}| / |\boldsymbol{\Sigma}| \cong \chi_n^2 \cdots \chi_{n-p+1}^2, \quad (2)$$

where the  $\chi_j^2$  variables, with  $j$  degrees of freedom,  $n - p + 1 \leq j \leq n$ , are independent.

This is also the distribution of  $n^p \frac{|\mathbf{S}|}{|\boldsymbol{\Sigma}|}$ , with  $n = N - 1$ .

On the other hand, Meijer's function  $\mathbf{G}(x)$  is defined as follows:

$$\mathbf{G} \begin{matrix} m & r \\ p & q \end{matrix} \left[ x \left| \begin{matrix} a_1, \dots, a_p \\ b_1, \dots, b_p \end{matrix} \right. \right] = \frac{1}{2\pi i} \int_L \frac{\prod_{j=1}^m \Gamma(b_j - s) \prod_{j=1}^r \Gamma(1 - a_j + s)}{\prod_{j=m+1}^q \Gamma(1 - b_j + s) \prod_{j=r+1}^p \Gamma(a_j - s)} x^s ds, \text{ i.e. it is the integral along the}$$

complex contour  $L$  of a ratio of products of gamma functions.  $\mathbf{G}$  has a relationship with the generalized hypergeometric function:  ${}_pF_{q-1}(\cdot)$ . Many common densities defined on  $(0, \infty)$ , can be expressed as a  $\mathbf{G}$ -function. For example, the gamma variable in two parameters, defined on  $(0, \infty)$ , with density  $f(x) = x^{\alpha-1} \exp(-x/\beta) / [\beta^\alpha \Gamma(\alpha)]$ ,  $0 \leq x < \infty$ ,  $\alpha, \beta > 0$  is a  $\mathbf{G}$ -function

$$\text{random variable since: } f(x) = \frac{1}{\beta \Gamma(\alpha)} \mathbf{G} \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} \left[ \frac{x}{\beta} \mid (\alpha - 1) \right], x > 0.$$

**THEOREM:** The ratio  $Y = \frac{|\mathbf{A}|}{|\boldsymbol{\Sigma}|} = n^p \frac{|\mathbf{S}|}{|\boldsymbol{\Sigma}|}$ , related to a random sample of size  $N = n+1$ , has its

density given by:

$$h(y) = \frac{1}{2^p} \left( \prod_{j=0}^{p-1} \frac{1}{\Gamma(\frac{n-j}{2})} \right) \mathbf{G} \begin{matrix} p & 0 \\ 0 & p \end{matrix} \left[ \frac{y}{2^p} \mid \frac{(n-2)}{2}, \frac{(n-3)}{2}, \dots, \frac{(n-(p+1))}{2} \right], y \geq 0 \quad (3)$$

Let  $h_{\alpha/2}$  and  $h_{1-(\alpha/2)}$  be the lower and upper  $\alpha/2$  percentiles of this density, that can be numerically determined with precision, using MAPLE. We then have:

$$h_\alpha \leq \frac{n^p |\mathbf{S}|}{|\boldsymbol{\Sigma}|} \leq h_{1-\alpha/2} \rightarrow \frac{n^p |\hat{\mathbf{S}}|}{h_{1-\alpha}/2} \leq |\boldsymbol{\Sigma}| \leq \frac{n^p |\hat{\mathbf{S}}|}{h_\alpha}, \text{ which provides the exact } (1-\alpha)100\%$$

confidence interval for the population generalized variance.

- a) Numerical example: Let us generate nine observations from a normal population (on Fisher's iris data), with covariance matrix  $\boldsymbol{\Sigma}$  taken from one of our previous publications. In  $R^4$  this distribution  $N_4(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , has known parameters,  $\boldsymbol{\mu}^T = (6.35, 6.20, 5.55, 5.23)$  and

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1.80 & 0.140 & 0.023 & 0.010 \\ 0.140 & 0.900 & 0.070 & 0.030 \\ 0.023 & 0.070 & 0.700 & 0.030 \\ 0.010 & 0.030 & 0.030 & 0.600 \end{pmatrix}. \text{ Hence } |\boldsymbol{\Sigma}| = 0.6645.$$

The nine observations generated lead to the estimated sample  $\mathbf{S}$ , with mean vector

$$\hat{\boldsymbol{\mu}}^T = \{6.485, 6.479, 5.102, 5.014\} \text{ and the estimated sample covariance matrix :}$$

$$\hat{\mathbf{S}} = \begin{pmatrix} 2.465 & 0.694 & -0.162 & 0.017 \\ 0.994 & 0.578 & -0.097 & -0.004 \\ -0.162 & -0.097 & 0.857 & 0.506 \\ 0.017 & -0.004 & 0.506 & 0.521 \end{pmatrix}, \text{ with determinant value: } |\hat{\mathbf{S}}| = 0.113.$$

a) With  $n = 8$ , the density of  $8^4 \times |\mathbf{S}| / |\boldsymbol{\Sigma}|$ , where  $\mathbf{S}$  is a random sample covariance matrix from the above distribution is given by (3). The exact 95% confidence interval for  $n^p |\mathbf{S}| / |\boldsymbol{\Sigma}|$ , is found numerically to be (65.35, 7934), as given by Fig.1. Hence, by (3), the corresponding 95% confidence interval for  $|\boldsymbol{\Sigma}|$ , based on the observed sample is: (0.0583, 7.0826), since  $8^4 \times |\hat{\mathbf{S}}| = 462.84$ , and contains the value of  $|\boldsymbol{\Sigma}|$ .

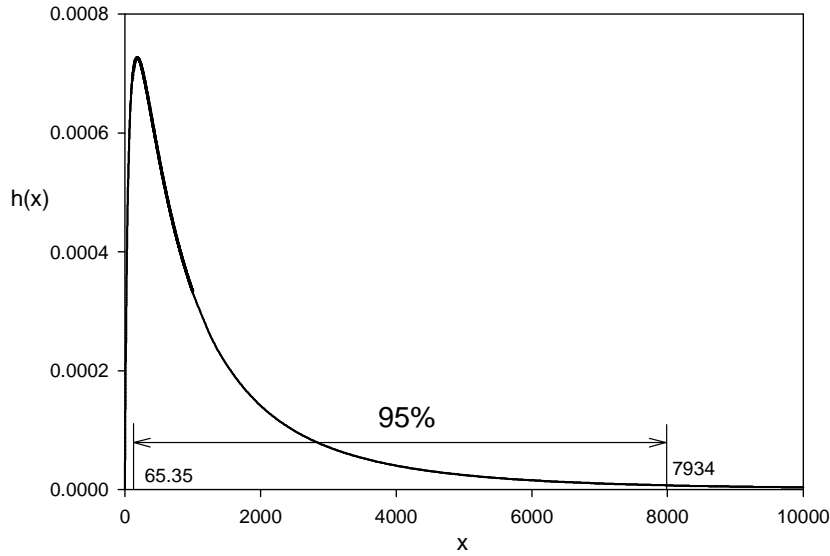


Fig. 1

Density of  $n^p |\mathbf{S}| / |\boldsymbol{\Sigma}|, n = 8, p = 4$

i) Normal approximation : Anderson (1984, p.266), based on the result that, asymptotically,  $\sqrt{n} \left( \frac{|\mathbf{S}|}{|\boldsymbol{\Sigma}|} - 1 \right) \sim N(0, 2p)$ . This would lead to the approximate confidence interval :

$$|\boldsymbol{\Sigma}| \in \left( \frac{|\hat{\mathbf{S}}|}{1 + z_{\alpha/2} \sqrt{\frac{2p}{n}}}, \frac{|\hat{\mathbf{S}}|}{1 - z_{\alpha/2} \sqrt{\frac{2p}{n}}} \right), \quad (4)$$

where  $P(Z \leq z_{\alpha/2}) = \alpha/2$ , and  $\mathbf{S}$  is defined in the introduction.

ii) Other researchers: T. Kubokawa and Srivastava (2003) provided a detailed survey of the different approaches proposed to date to estimate  $\Sigma$  and its determinant.

Iliopoulos and Kourouklis (1998) also used the product of chi-square variables, but not G-functions.

Hao and Krishnamoorthy (2001) gave some results based on chi-square approximation and missing data approaches.

Using the very rich properties of the G-function, we can derive:

$$a) \text{ Moments of } Y : E(|Y|^k) = 2^{pk} \prod_{j=1}^p \frac{\Gamma\left(\frac{n-(j-1)}{2} + k\right)}{\Gamma\left(\frac{n-(j-1)}{2}\right)}, k \geq 1, \text{ from which the moments of } \mathbf{S},$$

in function of  $\Sigma$ , can be obtained.

b) Bayesian approach: Since the sample covariance matrix  $\mathbf{S}$  is present in several results in Multivariate Bayesian analysis, we can derive the following result:

1) Let  $\mu$  be unknown in  $\mathbf{X} \sim N_p(\mu, \Sigma)$ , with a normal prior  $N(\mu_0, \tau_0)$ , where  $\tau_0$  being the matrix of precision. If  $\Sigma$  is known, we know that, with  $\bar{\mathbf{x}}$  being the mean of the sample  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , the posterior distribution of  $\mu$  is  $N_p(\mu^*, \tau_0 + N\mathbf{R})$ , where  $\mathbf{R} = \Sigma^{-1}$  and  $\mu^* = (\tau_0 \mu_0 + N\mathbf{R}\bar{\mathbf{x}})(\tau_0 + N\mathbf{R})^{-1}$ . With the 95% confidence bounds for  $\Sigma$  above, we can derive the bounds for the difference between prior and posterior precisions (for the mean) is:  $\left\{ \frac{N}{l_u}, \frac{N}{l_l} \right\}$ . Other similar results concern the cases:

2) Let  $\Sigma$  be unknown in  $\mathbf{X} \sim N_p(\mu_0, \Sigma)$ , with a Wishart prior,  $\Sigma \sim Wi_p(\alpha, \tau)$ , where  $\tau$  being the matrix of precision. We know that, with  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  being the observed sample, the posterior distribution of  $\Sigma$  is Wishart, i.e.

$$\Sigma | (\mathbf{x}_1, \dots, \mathbf{x}_N) \sim Wi_p(\alpha^* = \alpha + n, \tau^* = (\Sigma^*)^{-1} = \tau + \sum_{i=1}^N (\mathbf{x}_i - \mu_0)(\mathbf{x}_i - \mu_0)^T).$$

3) Both mean and precision matrix are unknown in  $\mathbf{X} \sim N_p(\mathbf{M}, \mathbf{R})$ . The prior of  $(\mathbf{M}, \mathbf{R}) \sim NWi_p(\alpha, \tau, \mu, \nu)$  is the Normal-Wishart distribution where :

$$\mathbf{R} \sim Wi_p(\alpha, \tau), \alpha \geq p, \mathbf{M} | \mathbf{R} = \mathbf{r} \sim N_p(\mu, \nu \mathbf{r}), \nu$$

The posterior distribution of  $(\boldsymbol{\mu}, \mathbf{R}) \sim NWi_p(\boldsymbol{\alpha}^*, \boldsymbol{\tau}^*, \boldsymbol{\mu}^*, \nu^*)$ , where

$$\boldsymbol{\alpha}^* = \boldsymbol{\alpha} + N, \quad \nu^* = \nu + N, \quad \boldsymbol{\mu}^* = \frac{\nu \boldsymbol{\mu} + N \bar{\mathbf{x}}}{\nu + N},$$

$$\text{and } \boldsymbol{\tau}^* = \boldsymbol{\tau} + \mathbf{A} + \frac{\nu N}{\nu + N} (\boldsymbol{\mu} - \bar{\mathbf{x}})(\boldsymbol{\mu} - \bar{\mathbf{x}})^T.$$

Some bounds for  $\boldsymbol{\tau}^* - \boldsymbol{\tau}$  can then be given.

## REFERENCES

- [1] Anderson, T.W., *Introduction to Multivariate Statistical Analysis*, 2<sup>nd</sup> Ed., John Wiley, New York, 1984.
- [2] Iliopoulos, G. and Kourouklis, S. (1998), On Improved Interval Estimation for the Generalized Variance, *Journ. of Stat. Plann. and Inf.*, 66, 305-420.
- [3] Hao, J. and Krishnamoorthy, K.. (2001), Inferences on a Normal Covariance Matrix and Generalized Variance with Monotone Missing Data, *J. of Mult. Anal.*, 78, 62-82.
- [4] Kubokawa, T. and Srivastava, M.S. (2003), Estimating the Covariance Matrix: a New Approach, Discussion Paper CIRJE-F-52, the University of Toronto.
- [5] Pham-Gia, T. and Turkkan, N., Density of the sample generalized Variance and applications, *Stat. Papers*, on line on the journal's site, 2009.