# Comparing Mean Vectors Via Generalized Inference in Multivariate Log-Normal Distributions

Kamel Abdollahnezhad[*1], Ali Akbar Jafari[2], Nasrin Tatari[1]
1. Golestan Universiy
2. Yazd University

**Extended Abstract**

## Introduction

The skewed distributions are particularly common when mean values are small, variances are large and values cannot be negative (for example lengths of latent periods of infectious diseases), and often closely fit the log-normal distribution. The log-normal distribution has been widely used in medical, biological and economic studies, where data are positive and have a right-skewed distribution. In this paper, we develop a procedure that is readily applicable for both hypothesis testing and confidence region construction for comparing the mean vectors of several multivariate log-normal distributions. The classical theory of statistical analysis is based on the assumption that the data are normally distributed. However, when the data are positive and skewed, the techniques of normality are inadequate to support inferential tests. For these data, the log-normal distribution is potentially a more suitable candidate to describe the positive and skewed random variables. Furthermore, once the data contains multiple and dependent variables, a multivariate approach instead of a univariate approach is practically necessary.

The mean of a log-normal random variable involves a linear combination of the mean and variance of the normal distribution of the log-data, and thus the inference procedures are more complicated concerning the log-normal mean. There are many applications for confidence intervals and tests regarding the log-normal means. Therefore, inference procedures concerning a single log-normal mean, or that of comparing the means of two or more independent log-normal distributions, attract more attention in the literature; see Zhou and Gao (1997), Taylor et al. (2002), Wu et al. (2002, 2003), Zhou et al. (1997), Chen and Zhou (2006), Gupta and Li (2006), Krishnamoorthy and Mathew (2003), Li (2009), Lin and Wang (2013) and Jafari and Abdollahnezhad (2015, 2017).

The applications of the multivariate log-normal distribution are similar to those of the univariate and bivariate log-normal distributions including, for example, the study of the size distribution of aerosol particles, airborne fibers, biomedical applications, etc. So far, there are no existing methods available in the literature for comparing mean vectors of multivariate log-normal distributions when several multivariate log-normal populations are not homogeneous. The multivariate analysis of variance (MANOVA) can only be applied when several log-data sets are homogeneous. Thus, the purpose of this study is to fill the gap by developing a procedure for comparing the mean vectors of several multivariate log-normal distributions based on the generalized variable approach (GVA).

## Material and methods

The GVA, including the generalized p-values and generalized confidence interval, was introduced by Tsui and Weerahandi (1989). The method has turned out to be extremely fruitful for obtaining tests and confidence intervals involving non-standard parameters. Numerous articles have revealed that this approach can satisfactorily deal with the problems which are heteroscedasticity, high dimensions or small sample sizes; see the books by Weerahandi for a detailed discussion along with numerous examples.

Let $X_1, \dots, X_n$ be a random sample from a $p$-variate normal population with mean $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)'$ and covariance matrix $\Sigma = (\sigma_{ij})$, $i = 1, \dots, p$, $j = 1, \dots, p$. Let $\bar{X}$ and $S$ denote the sample mean and sample covariance matrix, respectively. That is,

$$\bar{X} = \frac{1}{n}\sum_{h=1}^{n} X_h, \qquad S = (S_{lj}) = \sum_{h=1}^{n}(X_h - \bar{X})(X_h - \bar{X})'.$$

It is well-known that $\bar{X}$ and $S$ are independent with

$$\bar{X} \sim N_p\left(\boldsymbol{\mu}, \frac{1}{n}\Sigma\right), \quad S = W_p(n-1, \frac{1}{n-1}\Sigma).$$

where $W_p(r, \Sigma)$ denotes the $p$-dimentsional Wishart distribution with $r$ degrees of freedom and scale parameter matrix $\Sigma$.

We give a generalized pivotal variable for $\Sigma$. By Cholesky decomposition we can find a lower triangular matrix $\Lambda$ such that $\Sigma = \Lambda\Lambda'$. Let $A$ be Cholesky decomposition of $S$ such that $S = AA'$. Also, consider $B = (B_{lj}) = \Lambda^{-1}A$, then

$$BB' = \Lambda^{-1}A(\Lambda^{-1}A)' = \Lambda^{-1}S\Lambda'^{-1} \sim W_p(n-1, I_p).$$

Therefore, $B$ is a lower triangular matrix and $B_{ij}$'s are independent (Muirhead, 1982) and

$$B_{ll}^2 \sim \chi^2_{(n-l)}, \qquad B_{lj} \sim N(0,1), \quad l > j.$$

Let $s$ and $a$ be the observed matrices of $S$ and $A$, respectively. Then

$$V = (V_{lj}) = aB^{-1}B'^{-1}a' = aA^{-1}\Sigma A'^{-1}a',$$

is a generalized pivotal variable for $\Sigma$,

Now, we present a generalized pivotal variable for $\boldsymbol{\mu}$. Consider $\bar{x}$ is observed value for $\bar{X}$. Since $Z = \sqrt{n}\Lambda^{-1}(\bar{X} - \boldsymbol{\mu}) \sim N_p(\mathbf{0}, I_p)$, a generalized pivotal variable for $\boldsymbol{\mu}$ is

$$T_{\boldsymbol{\mu}} = \bar{x} - \frac{aB^{-1}}{\sqrt{n}}Z = \bar{x} - aA^{-1}(\bar{X} - \boldsymbol{\mu}).$$

Consider $T_{\boldsymbol{\sigma}} = (\sqrt{V_{11}}, \dots, \sqrt{V_{PP}})'$. Therefore, $T_{\boldsymbol{\theta}} = T_{\boldsymbol{\mu}} + \frac{1}{2}T_{\boldsymbol{\sigma}}$ is a generalized pivotal variable for $\left(\mu_1 + \frac{1}{2}\sigma_1, \dots, \mu_p + \frac{1}{2}\sigma_p\right)$.

## Results and discussion

Let $Y_{i1}, \dots, Y_{in_i}$ be a random sample from multivariate log-normal distribution with parameters $\boldsymbol{\mu}_i$ and $\Sigma_i$ ($i = 1, \dots, k$) where

$$\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{ip})', \qquad \Sigma_i = [\sigma_{i,st}], \quad i = 1, \dots, k.$$

Our problem is testing equality of means of this $k$ populations which equivalent to $H_0 : H\boldsymbol{\theta} = \mathbf{0}$, vs. $H_1 : H\boldsymbol{\theta} \neq \mathbf{0}$ where $\boldsymbol{\theta} = (\boldsymbol{\theta}_1', \dots, \boldsymbol{\theta}_k')'$, $\boldsymbol{\theta}_i = \boldsymbol{\mu}_i + \frac{1}{2}\boldsymbol{\sigma}_i$, $\boldsymbol{\sigma}_i = (\sigma_{i,11}, \dots, \sigma_{i,pp})'$, $H = Q \otimes I_p$, $Q = [-\mathbf{1}_{k-1} : I_{k-1}]$, $\mathbf{1}_{k-1}$ is a column vector of 1's

A generalized pivotal variable for the mean of $i$th population is $T_i = T_{\boldsymbol{\mu}_i} + \frac{1}{2}T_{\boldsymbol{\sigma}_i}$. Therefore, $T^* = H(T_1', \dots, T_k')'$, is a generalized pivotal variable for $H\boldsymbol{\theta}$.

Let $\widetilde{T}^*$ denote the standard expression of $T^*$ with $\widetilde{T}^* = \Sigma_{T^*}^{-\frac{1}{2}}(T^* - \mu_{T^*})$, where $\mu_{T^*}$ and $\Sigma_{T^*}$ are the conditional expectation and conditional variance-covariance matrix of $T^*$. The observed value of $T^*$ is $\tilde{t}^* = \Sigma_{T^*}^{-\frac{1}{2}}(H\theta - \mu_{T^*})$. Therefore, the generalized $p$-value for the test can be given by

$$p = P\big((T^* - \mu_{T^*})'\Sigma_{T^*}^{-1}(T^* - \mu_{T^*}) > \mu_{T^*}'\Sigma_{T^*}^{-1}\mu_{T^*}\big).$$

and $H_0$ will be rejected when $p$ is less than the level $\alpha$.

We performed a simulation study to compare the proposed generalized p-value and multivariate analysis of variance (MANOVA) to test the equality of means of several multivariate log-normal populations. The simulations show that the actual sizes of the both approaches are very smaller than the nominal level and there is not different between them. But the power of the proposed approach is larger than the power of MANOVA in small sample sizes.

## Conclusion

We consider the problem of means in several multivariate log-normal distributions and propose a useful method called as generalized variable method. Simulation studies show that suggested method has an appropriate size and power regardless sample size. To evaluation this method, we compare this method with traditional MANOVA such that the actual sizes of the two methods are close but the power of test and coverage probability of proposed methods are better than MANOVA in most cases specially when the sample sizes are small. Therefore, we can use this method when the variance-covariance matrices are not equal and there is not a suitable method.

**Keywords:** Generalized Variable, Multivariate log-normal distribution, Size and power of test, Coverage probability, Cholesky decomposition.

[*]Corresponding Author:    kamel_abdollahnezhad@yahoo.com