

# STATISTICAL HYPOTHESES TESTING IN VARIANCE ANALYSIS IN CASE OF CLASSICAL ASSUMPTIONS FAILURE<sup>1</sup>

V.M. PONOMARENKO, B.YU. LEMESHKO  
*Novosibirsk State Technical University*  
*Novosibirsk, Russia*  
e-mail: ponomarenkov@mail.ru

## Abstract

The distributions of likelihood ratio statistic used in variance analysis for testing hypotheses on model parameters have been investigated by statistical simulation methods. The case of non-normal error distributions has been considered. It has been shown that statistic distributions depend on the error distribution law and the method of model parameter estimation used. The approximations of the limiting statistic distributions for certain observation error laws have been constructed.

## 1 Introduction

In variance analysis while analysing the rate of investigated factor influence on the response under observation it is necessary to solve the tasks of testing hypotheses on significance of the differences in factors. As a rule, the likelihood ratio test statistic is used. The test statistic is given with [1]

$$Q = \frac{\left(\mathbf{K}^T \hat{\theta} - \mathbf{b}\right)^T \left(\mathbf{K}^T \mathbf{G} \mathbf{K}\right)^{-1} \left(\mathbf{K}^T \hat{\theta} - \mathbf{b}\right)}{\left(\mathbf{Y} - \mathbf{X} \hat{\theta}\right)^T \left(\mathbf{Y} - \mathbf{X} \hat{\theta}\right)} \frac{n-r}{k}, \quad (1)$$

where  $\mathbf{K}^T$  is a given  $(k \times m)$  matrix, determining the hypothesis under test

$$H : \mathbf{K}^T \hat{\theta} = \mathbf{b}, \quad (2)$$

$rg(\mathbf{K}) = k \leq m$ ;  $\mathbf{b}$  is a given  $(k \times 1)$  vector,  $\hat{\theta}$  is an estimate of parameter  $\theta$  vector of order  $m$  at the model

$$\mathbf{Y} = \mathbf{X} \theta + \mathbf{e}, \quad (3)$$

where  $\mathbf{Y}$  is the vector of observations of order  $n$ ,  $\mathbf{X}$  is the design matrix with  $(n \times m)$  dimension,  $r = rg(\mathbf{X})$ ,  $\mathbf{e}$  is the  $(n \times 1)$  vector of error terms;  $\mathbf{G}$  is a generalized inverse of a matrix  $\mathbf{X}^T \mathbf{X}$ . The components  $(e_1, \dots, e_n)$  of the  $\mathbf{e}$  vector are implied to be independent, identically distributed random variables with the null mathematical expectation and a certain variance  $\sigma^2$ .

---

<sup>1</sup>The research is supported by the Ministry of Education of the Russian Federation (project No. T02-3.3-3356)

In the classical variance analysis error terms  $e_i$  are assumed to have a normal distribution. Then the statistic (1) at the limit submits the  $F$ -distribution with  $k$  and  $n - r$  degrees of freedom. In case of using of the model (3)  $Q$  statistic submits the  $F$ -distribution with  $rg(\mathbf{K})$  and  $n - rg(\mathbf{X})$  degrees of freedom. However the assumption of error terms normality frequently doesn't hold in practice. So the purposes of the paper are: a) to investigate distributions of the  $Q$  statistic (1) by statistical simulation methods for non-normal error terms; b) to analyse the influence of the model parameter estimation method and the noise level on statistic distributions obtained.

## 2 The investigation of statistic distribution

$Q$  statistic distributions have been investigated for the exponential distribution family (De) with different values of the form parameter  $\lambda$  considered as an error distribution in the model (3):

$$De(\lambda) = f(x) = \frac{\theta_2}{2\theta_1\Gamma(1/\theta_2)} \exp \left\{ - \left( \frac{|x - \theta_0|}{\theta_1} \right)^\lambda \right\}. \quad (4)$$

The form parameter  $\lambda$  has been taken equal to 0.3, 0.5, 1, 2, 3, 5, 10.

In all considered cases the full factorial experiment and balanced design of observations in the model (3) have been implied. No more restrictions have been added in the model (3). Under such conditions the hypothesis of factor confidence can be formulated on the basis of estimable functions (EF) system [1]:

$$\begin{cases} \psi_1 = \alpha_2 - \alpha_1, \\ \psi_2 = \alpha_3 - \alpha_1, \\ \dots \\ \psi_{p-1} = \alpha_p - \alpha_1, \end{cases} \quad (5)$$

where  $\alpha_1, \dots, \alpha_p$  are the model parameters concerning the factor under test. The parameters are the elements of the true parameter vector  $\theta$ .

The error vector  $\mathbf{e}$  with components having some given distribution has been simulated with the given design matrix  $\mathbf{X}$ , vector of model parameters  $\theta$ , the noise level (10% of a significant signal, if something else isn't stipulated). Model parameters have been estimated by the least squares method (LSM) or by maximum likelihood method (MLM). Numerical investigations of statistic distributions have been carried out by means of the technique perfectly proved in statistical regularity investigations [2, 3].

For normal error terms the empirical distributions of the  $Q$  statistic, as it should be, fit to the limiting  $F(k, n - r)$  distribution with a high significance level achieved [2, 3]. Models of different dimensions (with different values of  $n, m, r, k$ ) have been considered.

It has been shown that the distributions of (1) with errors distributed by (4) turn out to be more robust to changes in error distributions using LSM. For example, figure 1 illustrates how the distributions of statistic (1) change with different error distributions ( $e \sim De(\lambda)$ ) for similar models for  $n = 36, m = 8, r = 6, k = 2$ . The classical limiting

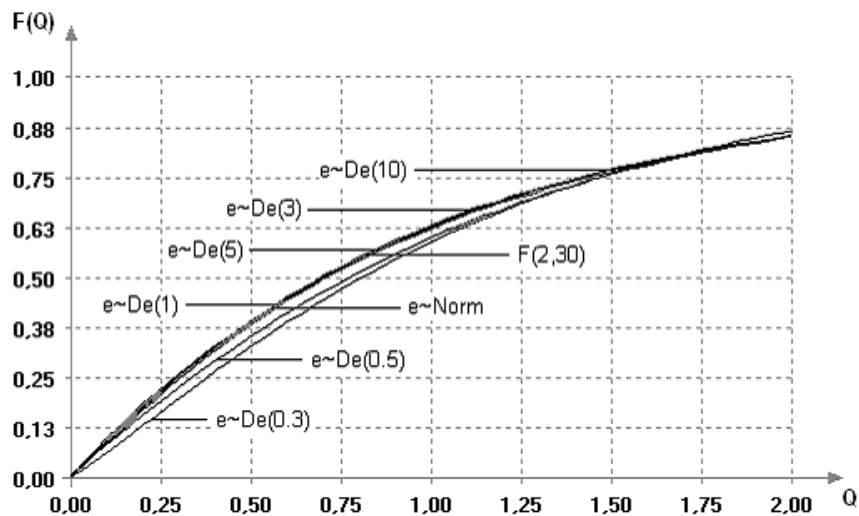


Figure 1:  $Q$  statistic distributions in dependence on the observed law, using LSM for estimation of model parameters

$F$ -distribution has been given for comparison. The visible deviation of the observed  $Q$  statistic distribution from the Fisher distribution takes place only in case of heavy-tailed error distributions ( $De(0.3)$  and  $De(0.5)$ ) using LSM for model parameter estimation.  $Q$  statistic distributions essentially depend on the error distribution law when MLM is used. In figure 2 the example represented in figure 1 is shown in case of using MLM. The distributions of statistic (1) here considerably differ from the Fisher  $F$ -distribution, but they are perfectly approximated with beta-distribution of the second kind. It has been shown that the revealed pattern of changes in  $Q$  statistic distributions also remains when noise level increases.

Statistical properties of the parametric functions  $\hat{\psi}_{LSM}$  and  $\hat{\psi}_{MLM}$ , considered as LSM and MLM estimates of EF (5), have been investigated.  $\hat{\psi}_{LSM}$  and  $\hat{\psi}_{MLM}$  functions obtained by substitution of  $\theta$  vector for its estimates  $\hat{\theta}_{LSM}$  and  $\hat{\theta}_{MLM}$  at the EF considered. The noise level has been taken for 100% during investigations of EF estimates properties. In case of normal error terms the empirical distributions of EF estimates simulated using MLM and LSM have coincided as it should be (with different values of  $n, m, r, k$ ). In case of non-normal error terms the distributions of MLM estimates of EF are essentially different from the distributions of LSM estimates. And MLM estimates of EF have better statistical properties here. MLM estimates of EF, as a rule, have a lesser variance under similar conditions. However, the distributions of MLM estimates of EF depend on error distribution to a greater extent. And at the same time MLM estimates of EF are not robust, especially for heavy-tailed error distributions. The distributions of LSM estimates of EF are on the contrary more robust to deviations of error distribution from the normal one. The investigation of EF estimate properties has explained differences in  $Q$  statistic distributions obtained using MLM and LSM for non-normal error distributions.

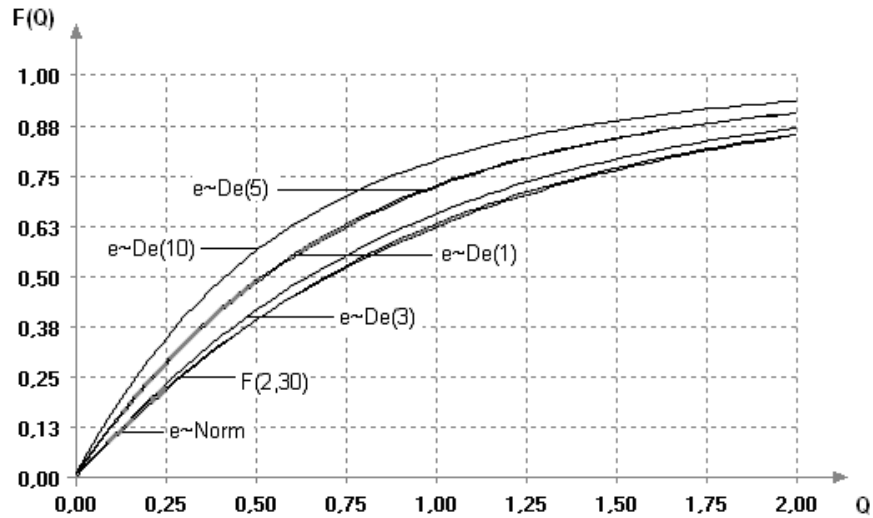


Figure 2:  $Q$  statistic distributions in dependence on the observed law, using MLM for estimation of model parameters

### 3 Conclusions

The investigations have shown that the level of influence of error non-normality on the limiting  $Q$  statistic distribution strongly depends on the method of model parameter estimation. The influence on the limiting  $Q$  statistic distribution is not large when using LSM. Considerable deviations have been observed only for heavy-tailed error distributions. When using MLM, the limiting  $Q$  statistic distribution essentially depends on the error distribution.

### References

- [1] Markova E.V., Denisov V.I., Poletaeva I.A., Ponomarev V.V. (1982). Variance analysis and synthesis of designs performing on computer. Nauka, Moscow. (in Russian)
- [2] R 50.1.033-2001. (2002). Recommendations for standardization. Applied statistics. Rules of check of experimental and theoretical distribution of the consent. Part I. Goodness-of-fit tests of a type chi-square. Publishing house of the standards, Moscow. (in Russian)
- [3] R 50.1.037-2002. (2002). Recommendations for standardization. Applied statistics. Rules of check of experimental and theoretical distribution of the consent. Part II. Nonparametric goodnessoffit test. Publishing house of the standards, Moscow. (in Russian)