



One-way random effects ANOVA with random sample sizes: An application to a Brazilian database on cancer registries

Gilberto Capistrano, Célia Nunes, Dário Ferreira, Sandra S. Ferreira, and João T. Mexia

Citation: [AIP Conference Proceedings](#) **1648**, 110009 (2015); doi: 10.1063/1.4912416

View online: <http://dx.doi.org/10.1063/1.4912416>

View Table of Contents: <http://scitation.aip.org/content/aip/proceeding/aipcp/1648?ver=pdfcov>

Published by the [AIP Publishing](#)

Articles you may be interested in

[One-way fixed effects ANOVA with missing observations](#)

AIP Conf. Proc. **1648**, 110008 (2015); 10.1063/1.4912415

[ANOVA with random sample sizes: An application to a Brazilian database on cancer registries](#)

AIP Conf. Proc. **1558**, 825 (2013); 10.1063/1.4825623

[One-way random effects ANOVA: An extension to samples with random size](#)

AIP Conf. Proc. **1479**, 1678 (2012); 10.1063/1.4756492

[Inference with Inducer Pivot Variables, an Application to the One-Way ANOVA](#)

AIP Conf. Proc. **1389**, 1631 (2011); 10.1063/1.3636921

[Randomized Sample Size F Tests for the One-Way Layout](#)

AIP Conf. Proc. **1281**, 1248 (2010); 10.1063/1.3497917

One-way Random Effects ANOVA with Random Sample Sizes: An Application to a Brazilian Database on Cancer Registries

Gilberto Capistrano*, Célia Nunes[†], Dário Ferreira[†], Sandra S. Ferreira[†] and João T. Mexia**

*University Center of Itajubá - FEPI, Brazil and Center of Mathematics, University of Beira Interior, Portugal

[†]Department of Mathematics and Center of Mathematics, University of Beira Interior, Covilhã, Portugal

**Department of Mathematics and CMA, Faculty of Science and Technology, New University of Lisbon, Portugal

Abstract. ANOVA is routinely used in many situations, namely in medical research, where the sample sizes may not be previously known. This leads us to consider the samples sizes as realizations of random variables. The aim of this paper is to extend one-way random effects ANOVA to those situations and apply our results to a Brazilian database on cancer registries.

Keywords: *F*-tests, Random effects model, Random sample sizes, Cancer registries.

AMS: 62J12, 62J10, 62J99.

INTRODUCTION

In many situations, such as in medical research, sample sizes may not be previously known. This often occurs when there is a specific time span for collecting the observations. The aim of this paper is to extend one-way random effects models to those situations and apply our results to a Brazilian database on cancer registries.

In these situations it is more correct to consider the samples sizes as realizations n_1, \dots, n_r of independent random variables N_1, \dots, N_r , see [2, 3, 4, 5, 6, 7]. We assume that these random variables will be independent Poisson distributed with parameters $\lambda_1, \dots, \lambda_r$. Therefore $n = \sum_{i=1}^r n_i$ will be a realization of the random variable $N = \sum_{i=1}^r N_i$ which is Poisson distributed with parameter $\lambda = \sum_{i=1}^r \lambda_i$, see [2, 6, 7].

When $N_i = n_i$, $i = 1, \dots, r$, the one-way random effects model can be written as, see e.g. [8],

$$Y_{i,j} = \mu + \alpha_i + e_{i,j}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, r, \quad (1)$$

where μ is a fixed unknown parameter and α_i and $e_{i,j}$, $j = 1, \dots, n_i$, $i = 1, \dots, r$, are independent normal variables with null mean values and variances, respectively, σ_α^2 and σ^2 , $\alpha_i \sim N(0, \sigma_\alpha^2)$, $e_{i,j} \sim N(0, \sigma^2)$, $j = 1, \dots, n_i$, $i = 1, \dots, r$.

The model (1) can be written in matrix notation as, see e.g. [1],

$$\mathbf{Y} = \mu + D(\mathbf{1}_{n_1}, \dots, \mathbf{1}_{n_r})\boldsymbol{\alpha} + \mathbf{e}, \quad (2)$$

where $\mathbf{Y} = (Y_{1,1}, \dots, Y_{1,n_1}, Y_{2,1}, \dots, Y_{2,n_2}, \dots, Y_{r,1}, \dots, Y_{r,n_r})'$, $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_r)'$ and $D(\mathbf{1}_{n_1}, \dots, \mathbf{1}_{n_r})$ denotes a block diagonal matrix with $\mathbf{1}_{n_1}, \dots, \mathbf{1}_{n_r}$ along the blocks and $\mathbf{1}_m$ is the vector with all m components equal to 1. Therefore, $\boldsymbol{\alpha}$ and \mathbf{e} are normal distributed, with null mean vectors and variance-covariance matrices, respectively, $\sigma_\alpha^2 \mathbf{I}_r$ and $\sigma^2 \mathbf{I}_n$, $\boldsymbol{\alpha} \sim N(\mathbf{0}, \sigma_\alpha^2 \mathbf{I}_r)$, $\mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, with \mathbf{I}_n the $n \times n$ identity matrix.

The vector \mathbf{Y} , when $N = n$, will be conditionally normal with mean vector $\mu = \mu \mathbf{1}_r$ and variance-covariance matrix given by, see [1] and [9], $\boldsymbol{\Sigma} = \sigma_\alpha^2 D(\mathbf{J}_{n_1}, \dots, \mathbf{J}_{n_r}) + \sigma^2 \mathbf{I}_n$, where $\mathbf{J}_m = \mathbf{1}_m \mathbf{1}_m'$. We put

$$\mathbf{Y} \underset{(N=n)}{\sim} N(\mu, \sigma_\alpha^2 D(\mathbf{J}_{n_1}, \dots, \mathbf{J}_{n_r}) + \sigma^2 \mathbf{I}_n).$$

We are interested in test the hypotheses

$$H_0 : \sigma_\alpha^2 = 0 \text{ vs } H_1 : \sigma_\alpha^2 > 0. \quad (3)$$

In what follows we present the test statistic for testing these hypotheses and their conditional and unconditional distributions. We will consider the expression of the test statistic obtained in [5]. Then we present an application based on real medical data, particularly on cancer registries from São Paulo, Brazil. Finally we conclude with some final remarks.

STATISTIC AND THEIR DISTRIBUTIONS

In this section we will obtain the common conditional distribution of the test statistic and also its unconditional distribution, under the assumption that we have a global minimum dimension for the samples, which avoids highly unbalanced cases.

When $N_i = n_i, i = 1, \dots, r$, we have the averages $Y_{i,\bullet}, i = 1, \dots, r$, for the samples $Y_{i,1}, \dots, Y_{i,n_i}, i = 1, \dots, r$. The sum of squares for the error will be given by, see e.g. [9],

$$S = \sum_{i=1}^r \sum_{j=1}^{n_i} (Y_{i,j} - Y_{i,\bullet})^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} Y_{i,j}^2 - \sum_{i=1}^r \frac{T_i^2}{n_i}, \quad (4)$$

with $T_i = \sum_{j=1}^{n_i} Y_{i,j}, i = 1, \dots, r$ and S will be the product by σ^2 of a central chi-square with $g(n) = n - r$ degrees of freedom, $S \sim \sigma^2 \chi_{g(n)}^2$. As we can see, for instance in [9], S can be written using the matrix formulation as

$$S = \mathbf{Y}' \left(\mathbf{I}_n - D \begin{pmatrix} \mathbf{J}_{n_1} \\ \mathbf{J}_{n_2} \\ \vdots \\ \mathbf{J}_{n_r} \end{pmatrix} \right) \mathbf{Y}. \quad (5)$$

When $N_i = n_i$, the sample means

$$Y_{i,\bullet} = \mu + \alpha_i + e_{i,\bullet}, \quad i = 1, \dots, r,$$

have mean values μ and variances $\sigma_\alpha^2 + \frac{\sigma^2}{n_i}, i = 1, \dots, r$. So, when the hypothesis H_0 holds,

$$Y_{i,\bullet} \sim N\left(\mu, \frac{\sigma^2}{n_i}\right), \quad i = 1, \dots, r, \text{ and } Y_{\bullet,\bullet} = \frac{1}{n} \sum_{i=1}^r n_i Y_{i,\bullet} \sim N\left(\mu, \frac{\sigma^2}{n}\right).$$

Let us consider

$$\mathbf{Z} = \mathbf{Y}_\bullet - \mathbf{Y}_{\bullet,\bullet} = \mathbf{B}\mathbf{Y}_\bullet,$$

where $\mathbf{Y}_\bullet = (Y_{1,\bullet}, \dots, Y_{r,\bullet})', \mathbf{Y}_{\bullet,\bullet} = (Y_{\bullet,\bullet}, \dots, Y_{\bullet,\bullet})'$ and

$$\mathbf{B} = \mathbf{I}_r - \begin{bmatrix} \frac{n_1}{n} & \dots & \frac{n_r}{n} \\ \vdots & \ddots & \vdots \\ \frac{n_1}{n} & \dots & \frac{n_r}{n} \end{bmatrix} = \begin{bmatrix} \frac{n-n_1}{n} & \dots & \frac{-n_r}{n} \\ \vdots & \ddots & \vdots \\ \frac{-n_1}{n} & \dots & \frac{n-n_r}{n} \end{bmatrix} = \mathbf{I}_r - \frac{1}{n} \mathbf{1}_r \mathbf{n}'.$$

When $N_i = n_i, i = 1, \dots, r$, \mathbf{Z} will be normal distributed with null mean vector and variance-covariance matrix $\sigma^2 \mathbf{V}$, $\mathbf{Z} \sim N(\mathbf{0}, \sigma^2 \mathbf{V})$, with

$$\mathbf{V} = \mathbf{B} D \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \mathbf{B}' = D \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \frac{1}{n} \mathbf{J}_r.$$

The sum of squares for the main effects of the factor will be given by $S_{num} = \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z}$, where \mathbf{V}^{-1} denote a generalized inverse of matrix \mathbf{V} . When H_0 holds, we have

$$S_{num} = \mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z} \sim \sigma^2 \chi_g^2,$$

with $g = \text{rank}(\mathbf{V}) = r - 1$, see [5].

Therefore, when $N = n$ and H_0 holds, the common conditional distribution of the test statistic

$$\mathfrak{S} = \frac{S_{num}}{S} = \frac{\mathbf{Z}' \mathbf{V}^{-1} \mathbf{Z}}{S} \quad (6)$$

will be $\bar{F}(\cdot | g, g(n))$, which is the distribution of the quotient of independent central chi-squares with g and $g(n)$ degrees of freedom.

For carrying out the inference we assume that $N \geq n^\bullet$, which means that we assume that we have a global minimum dimension for the samples, see [2, 5, 6, 7]. So we take

$$p_{n^\bullet}(n) = pr(N = n | N \geq n^\bullet) = \frac{pr(N = n)}{pr(N \geq n^\bullet)} = \frac{e^{-\lambda} \lambda^n / n!}{1 - \sum_{u=0}^{n^\bullet-1} e^{-\lambda} \lambda^u / u!} = \frac{\lambda^n}{n! \left(e^\lambda - \sum_{u=0}^{n^\bullet-1} \lambda^u / u! \right)}, \quad n = n^\bullet, \dots \quad (7)$$

and we obtain the unconditional distribution given by

$$\bar{F}(z) = \sum_{n=n^\bullet}^{\infty} pr(N = n | N \geq n^\bullet) pr(\mathfrak{S} \leq z | N = n) = \sum_{n=n^\bullet}^{\infty} p_{n^\bullet}(n) \bar{F}(z | g, g(n)). \quad (8)$$

AN APPLICATION

The data used in this application were provided by the National Cancer Institute (INCA) and are from city of São Paulo, Brazil, 2010. We selected $r = 6$ different kinds of cancer, using simple random sampling. The Table 1 illustrates the kinds of cancer which have been selected and the number of patients.

TABLE 1. Kinds of cancers and number of patients

Kinds of Cancer	Number of patients
Body of Stomach	91
Spinal cord and other parties S.N.C.	42
Melanoma on the trunk	107
Encephalon	93
Ascending Colon	201
Upper lobe, bronchus or lung	155

As we saw, given $N = n$, the conditional distribution of \mathfrak{S} is, when H_0 holds, a central \bar{F} distribution with $g = 5$ and $g(n) = n - 6$ degrees of freedom, $\bar{F}(\cdot|5, n - 6)$.

To carry out the calculations when we may assume that $\sum_{n=0}^{n^* - 1} p_{n^*}(n) \simeq 0$, which means that, with high probability, we have $N \geq n^*$, the unconditional distribution of the statistic will be given by

$$\bar{\bar{F}}(z) = \sum_{n=n^*}^{\infty} p_{n^*}(n) \bar{F}(z|5, n - 6).$$

Moreover, to the monotony property of the \bar{F} distribution, see [2], we have $\bar{F}(z|g, n - 6) < \bar{F}(z|g, n^o - 6)$, with $n < n^o$, so

$$\bar{F}(z|5, n^* - 6) \leq \bar{\bar{F}}(z) \leq 1,$$

which gives us a lower bound for $\bar{\bar{F}}(z)$. Thus, from $\bar{F}(z|5, n^* - 6)$, we can obtain upper bounds for the quantiles of the unconditional distribution $\bar{\bar{F}}(z)$. If we use these upper bounds as critical values we will have tests with sizes that do not exceed the theoretical values.

If the statistic's value exceeds the upper bounds, also exceeds the real critical value (obtaining considering random sample sizes) and in this case we reject the hypothesis. When the statistic's value is lower than the upper bound we must compute the real critical values or calculate the minimum value of n^* that leads to reject the null hypothesis.

In this case we obtain

$$S = \sum_{j=1}^{91} (y_{1,j} - y_{1,\bullet})^2 + \sum_{j=1}^{42} (y_{2,j} - y_{2,\bullet})^2 + \dots + \sum_{j=1}^{155} (y_{6,j} - y_{6,\bullet})^2 = 209006, \quad (9)$$

with the sample means $y_{1,\bullet} = 62.05; y_{2,\bullet} = 46.17; y_{3,\bullet} = 67.89; y_{4,\bullet} = 49.90; y_{5,\bullet} = 71.10; y_{6,\bullet} = 66.26$. The general mean obtained is equal to $y_{\bullet,\bullet} = 64.02$.

The numerator of the \mathfrak{S} statistic, when the hypothesis holds, was defined by $S_{num} = \mathbf{Z}'\mathbf{V}^{-1}\mathbf{Z} \sim \sigma^2\chi_5^2$. In this case we obtain

$$\mathbf{Z} = \mathbf{B}\mathbf{Y}_{\bullet} = \begin{bmatrix} -1.965375 \\ -17.853653 \\ 3.867530 \\ -14.988060 \\ 7.084158 \\ 2.237745 \end{bmatrix}; \mathbf{V}^{-1} = \begin{bmatrix} 88.6758568 & 0.9076503 & -4.3801736 & -2.5542037 & -26.409124 & -13.504328 \\ 0.9076503 & 43.3329148 & 0.3069035 & 0.8449923 & -7.814687 & -2.859674 \\ -4.3801736 & 0.3069035 & 99.9126454 & -4.6868906 & -34.691226 & -18.684709 \\ -2.5542037 & 0.8449923 & -4.6868906 & 90.2067458 & -27.384874 & -14.105984 \\ -26.4091237 & -7.8146868 & -34.6912261 & -27.3848744 & 95.674453 & -66.066861 \\ -13.5043283 & -2.8596742 & -18.6847092 & -14.1059835 & -66.066861 & 115.739094 \end{bmatrix}$$

and consequently $S_{num} = \mathbf{Z}'\mathbf{V}^{-1}\mathbf{Z} = 47067.08$. Therefore, the statistic's value, \mathfrak{S}_{Obs} , is given by

$$\mathfrak{S}_{Obs} = \frac{47067.08}{209006} = 0.2252.$$

If we use the common conditional distribution of \mathfrak{S} , which corresponds to $\bar{F}(z|5, 683)$, we will obtain the quantiles, $z_{1-\alpha}$, given in Table 2. So, since $\mathfrak{S}_{Obs} > z_{1-\alpha}$, we can conclude that we reject H_0 for the usual levels of significance.

Let us now assume that $n^* = 59$, which means that we have about 10 observations per treatment. The Table 2 shows the upper bounds for the quantiles, $z_{1-\alpha}^u$, for probability $1 - \alpha$ of the unconditional distribution $\bar{\bar{F}}(z)$.

The results in this table may lead us to take a contrary decision that we had taken using the conditional distribution of the statistics for $\alpha = 0.05$ and $\alpha = 0.01$. Therefore the increase of the critical values points towards the possibility of non-rejection when random sample sizes are considered.

TABLE 2. The quantiles of the conditional distribution of the statistics and upper bounds for the quantiles

Values de α	0.1	0.05	0.01
$z_{1-\alpha}$	0.0136	0.0163	0.0223
$z_{1-\alpha}^u$	0.1848	0.2254	0.3193

TABLE 3. Minimum value of n^\bullet that leads to reject the null hypothesis

Values of α	0.1	0.05	0.01
n^\bullet	51	60	79

Assuming the values of the test statistic remain unchanged, for ensuring rejection, we should have the total sample sizes presented in Table 3. Since for higher values of n^\bullet we would get lower values for the quantiles, we have $\mathfrak{S}_{Obs} > z_{1-\alpha}^u$ for all $n^\bullet \geq 79$, which means that, in this case, we reject H_0 considering the usual levels of significance. Thus the kind of cancer has a significant random effect, so the ages of disease detection may differ significantly with it's kind.

FINAL REMARKS

When we cannot previously know the sample sizes it is much more correct to consider them as realizations of random variables. Through the application on cancer registries we can prove the relevance of the unconditional approach in order to possibly avoid false rejections. We also can conclude that when the samples dimensions increase, the conditional and unconditional approach converge to the same decision.

During our treatment we worked with \bar{F} distributions since they are more treatable and statistically equivalent. This equivalence enabled us to consider our tests as F -tests.

ACKNOWLEDGMENTS

This work was partially supported by Center of Mathematics, University of Beira Interior, through the project PEst-OE/MAT/UI0212/2014 and by CMA, Faculty of Science and Technology, New University of Lisbon, through the project PEst- OE/MAT/UI0297/2014.

REFERENCES

1. A.I. Khuri, T. Mathew, B.K. and B.K. Sinha, *Statistical Tests for Mixed Linear Models*. John Wiley & Sons, Inc., New York, (1998).
2. J.T. Mexia, C. Nunes, D. Ferreira, S.S. Ferreira and E. Moreira, "Orthogonal fixed effects ANOVA with random sample sizes", in *Proceedings of the 5th International Conference on Applied Mathematics, Simulation, Modelling (ASM'11)*, Corfu, Greece, 2011, pp. 84-90.
3. C. Nunes, D. Ferreira, S.S. Ferreira and J.T. Mexia, "F Tests with Random Sample Sizes", in *8th International Conference on Numerical Analysis and Applied Mathematics, AIP Conf. Proc.*, 2010, 1281(II), pp. 1241-1244.
4. C. Nunes, D. Ferreira, S.S. Ferreira and J.T. Mexia, F -tests with a rare pathology. *Journal of Applied Statistics* **39**(3), 551–561 (2012).
5. C. Nunes, D. Ferreira, S.S. Ferreira, M.M. Oliveira and J.T. Mexia, "One-way random effects ANOVA: An extension to samples with random size", in *10th International Conference on Numerical Analysis and Applied Mathematics, AIP Conf. Proc.*, 2012, 1479, pp. 1678-1681.
6. C. Nunes, G. Capistrano, D. Ferreira and S. S. Ferreira, "ANOVA with Random Sample Sizes: An Application to a Brazilian Database on Cancer Registries", in *11th International Conference on Numerical Analysis and Applied Mathematics, AIP Conf. Proc.*, 2013, 1558, 825-828.
7. C. Nunes, D. Ferreira, S. S. Ferreira and J.T. Mexia. Fixed effects ANOVA: an extension to samples with random size. *Journal of Statistical Computation and Simulation*, **84**(11), 2316-2328, (2014).
8. H. Scheffé, *The analysis of variance*. Wiley series in Probability and Statistics, John Wiley & Sons, New York, (1959).
9. S.R. Searle, G. Casella and C.E. McCulloch, *Variance Components*. Wiley series in Probability and statistics. John Wiley & Sons, New York, (1992).