

AN ALTERNATIVE APPROACH FOR MULTIPLE COMPARISON PROBLEMS WHEN THERE ARE A LARGE NUMBER OF GROUPS: ANOM TECHNIQUE

M. Mendeş* and S. Yiğit

Çanakkale Onsekiz Mart University, Faculty of Agriculture, Department of Animal Science, Biometry and Genetics Units,
Çanakkale, Turkey

*Correspondence author's e-mail: mmendes@comu.edu.tr

ABSTRACT

Analysis of Means (ANOM) is a powerful tool for comparing means, variances, proportions and other location and scale measures. This procedure can also be used efficiently as a multiple comparison test especially when there are a large number of groups. A simulation study has been carried out to investigate the performance of ANOM as a multiple comparison approach. Total accuracy value was found at least 60.00% regardless of experimental conditions. ANOM and Tukey procedures have also been compared on a real data set where there are 80 different wheat varieties. It is concluded that the ANOM has a great advantage over Tukey test in terms of determining superior and worse varieties especially when there are more than 10 treatment groups.

Key words: ANOM, multiple comparisons, sensitivity, total accuracy.

INTRODUCTION

ANOVA is perhaps one of the most commonly used statistical methods for comparing independent group means. However, ANOVA only tells us if there is a significant difference among the group means. It does not tell us anything about which group mean is different from the others (Williams and Abdi, 2010). If null hypothesis (H_0) is rejected then an appropriate multiple comparison test (e.g. Tukey, SNK, Duncan) should be applied after ANOVA to determine different group means (Duncan, 1955). Multiple comparison procedures enable us to determine which treatments are better / superior and which ones are worse (Tukey, 1949). It is not difficult to determine different group means as long as the number of treatment groups is not large (e.g. if number of group is not exceeding 10). However, in some cases, especially studies related to plant breeding in crop science, researchers are commonly interested in comparing differences among a large number of treatment groups. For example, one might want to compare 140 corn genotypes or 200 different wheat genotypes. Or one might want to determine superior varieties among 180 corn varieties. There will be a large number of treatment groups in such cases. Classical multiple comparison tests like Tukey, SNK, Duncan, LSD, Bonferroni tests will not be appropriate in terms of determining different group means if the number of group means being compared is very large. However, many researchers have still been using these kinds of multiple comparison tests even though they have a large number of treatment groups. For example, Munyiri *et al.* (2013) used Least Significant Difference (LSD) pairwise comparison procedure to determine different and superior maize varieties among 25 varieties. Likewise, in their study, Flint-

Garcia *et al.* (2009) used LSD procedure after ANOVA-F test to determine different corns among 55 varieties. Wietholter *et al.* (2008) compared 41 corn populations by using ANOVA-F test and then they used Tukey pairwise comparison test to identify different populations. Sharma *et al.* (2012) compared 40 wheat varieties by using Least Significant Difference pairwise comparison test (LSD) after ANOVA. Then, they also applied Principal Component Analysis (PCA) to figure out superior wheat varieties. In their study, Karaköy *et al.* (2014) compared 178 faba bean varieties. They used ANOVA-F test to compare differences among these varieties in terms of measured morphological characteristics. After ANOVA-F test, Cluster and Principal Component Analyses performed to classify those varieties. Garcia-Lara and Bergvinson (2013) used LSD pairwise comparison procedure after ANOVA to determine different varieties among 35 corn varieties. For such cases, there are some other techniques like **Analysis of Means (ANOM)** which can be used for both identify treatments which are superior or worse than overall mean and similar to the overall mean. Therefore, the usage of ANOM technique instead of the classical pairwise comparison procedures will make possible to get more detailed and reliable results. Because the ANOM technique is not only used for comparing treatment group means and determining superior or worse treatments but also comparing variances, proportions, correlation coefficients, and other location and scale measures across several groups. That is why, especially recently it is remarkable that the ANOM procedure has become increasingly popular among the scientists and it has been implemented in many statistical package programs such as SAS (1999), MINITAB 18 (2017), JMP 12 (2016), and R (Pallmann and Hothorn, 2016). This paper has aimed to show usability of ANOM technique as a multiple

comparison test especially when there are a large number of treatment groups and to determine performances of ANOM technique under different experimental conditions via Monte Carlo Simulation Study. ANOM Technique has also been applied to a real data set to determine superior, worse and similar treatments when compared to overall mean.

MATERIALS AND METHODS

The main objective of this study is to investigate some characteristics of ANOM technique as an alternative of multiple comparison tests when there are a large number of groups. In this study, a total of 50,000 data sets along with the combinations of number of group i.e.: $k=10, 20, 30, 40, 50, 80, 100$ and 200 , sample size i.e.: $n=5, 10, 20, 30$ and 50 ; effect size i.e.: $0.3, 0.6, 0.9$ and 1.2 and percentage of different groups in reality i.e.: $d=20\%, 30\%$ and 40% have been generated from five different distributions ranging from the normal to severe deviations from normality i.e.: Normal $(0, 1)$, Beta $(10,10)$, Beta $(5,10)$, Beta $(10,5)$, and Chi-Sq with 3 df by using RNNOA, RNBET, and RNCHI functions of IMSL library of Fortran Developer Studio's (Anonymous, 1994). Fixed numbers have been added and subtracted to a certain percentage of these groups ($d=20\%, 30\%$ and 40%) to create different groups among the groups being compared. These scenarios have been simulated only under homogeneity of variances which is one of the basic assumptions of the ANOM technique. Then Specificity, Sensitivity and Total accuracy values were estimated for each experimental condition to find out how many of these differentiated and undifferentiated groups were correctly identified by the ANOM technique.

Performance criteria computed as follows:

$$\text{False Positive Rate (FPR): } \text{FPR} = \frac{B}{m_0}$$

$$\text{False Negative Rate (FNR): } \text{FNR} = \frac{C}{m_1}$$

$$\text{Sensitivity} = \frac{D}{m_1} = 1 - \text{FNR}$$

$$\text{Specificity} = \frac{A}{m_0} = 1 - \text{FPR}$$

$$\text{Total accuracy} = \frac{A + D}{m}$$

Analysis of Means (ANOM): Analysis of Means (ANOM), introduced by Ott (1983), is a technique for comparing group means. Since it presents the comparisons graphically, the researchers can easily see which treatment

mean(s) are different. This is a big advantage especially for non-statisticians. Although the ANOM is accepted as a graphical alternative to ANOVA for comparing independent group means, it has two advantages over ANOVA especially when authors are interested in studying main effects (Nelson *et al.*, 2005; Mendes and Yiğit, 2013). The first advantage of the ANOM over the ANOVA is: if any of the group mean is statistically different from the others, it enables the researchers to see exactly which ones is different easily. The second advantage of the ANOM is: since the ANOM presented the results graphically, it is easy to assess both the statistical and practical significance of observed differences simultaneously. Even though the null hypotheses of two techniques are the same, the alternative hypotheses of them are different. In its alternative hypothesis, the ANOM, states that at least one treatment mean is significantly different from the grand mean, while the alternative hypothesis of the ANOVA claims that at least one treatment mean is significantly different from the others. Although both techniques provide similar results in general, these methods might lead to different results due to differences in the alternative hypotheses of them.

The results of the ANOM technique are based on confidence interval or decision lines (upper decision line (UDL) and lower decision line (LDL)).

The computation of the upper and lower decision lines are given as below (Nelson *et al.*, 2005; Mendes and Yiğit, 2013):

Computation of UDL and LDL for equal sample size:

$$\text{UDL} = \bar{Y}_{..} + h(\alpha, k, N - k) \sqrt{\text{MSE}} \sqrt{\frac{(k-1)}{N}}$$

$$\text{LDL} = \bar{Y}_{..} - h(\alpha, k, N - k) \sqrt{\text{MSE}} \sqrt{\frac{(k-1)}{N}}$$

Computation of UDL and LDL for unequal sample size:

$$\text{UDL} = \bar{Y}_{..} + h(\alpha, k, N - k) \sqrt{\text{MSE}} \sqrt{\frac{(N - n_i)}{N n_i}}$$

$$\text{LDL} = \bar{Y}_{..} - h(\alpha, k, N - k) \sqrt{\text{MSE}} \sqrt{\frac{(N - n_i)}{N n_i}}$$

Where $\bar{Y}_{..}$, h , k , and MSE denote overall mean, critical table values for the ANOM technique based on α number of treatment group number, and mean square error (Schilling, 1973; Nelson, 1983, 1985, 1989, 1993; Nelson *et al.*, 2005).

At the end of hypothesis testing, if all means fall between the UDL and LDL, then the null hypothesis is accepted and it is concluded that all means are equal. Any of the group mean, on the other hand, falls outside the decision lines, then the null hypothesis is rejected and it is concluded that at least one group mean is significantly different from the overall or grand mean.

RESULTS AND DISCUSSION

Results of Simulation Study: Sensitivity, specificity, and total accuracy estimates of the ANOM have been presented in Figure 1-20, respectively. Specificity values have been found very close to 100% regardless of number of groups being compared (k), sample size (n), effect size (δ), and percentage of different groups (d). They changed between 97.04 % and 99.98 % in general. These findings are one of the indicators that the ANOM is very specific technique in terms of identifying treatment means which are similar to the grand mean. In other words, the ANOM technique can correctly identify almost all groups that are between upper and lower decision lines (or not significantly deviate from overall mean). This is a very big and an important advantage of the ANOM over the Classical Pairwise Multiple Comparison Procedures like LSD, Tukey, Duncan, Dunnett, SNK etc. especially when there are a large number of treatment groups.

As it is seen in Figure 1-20, numbers of groups (k), sample size (n) and effect size (δ) have significant impact on sensitivity estimates while distribution shape and percentage of different groups (d) have not. The sensitivity values of the ANOM have varied between 0.12% and 100.00% in general. The lower sensitivity values have been obtained when $n=5$, $\delta=0.30$ and $k=200$. However, as it is expected as the n and δ were increased, the sensitivity values increased as well. Although three criteria namely sensitivity, specificity and total accuracy have been taken into account, interpretations have been made on the total accuracy values since total accuracy is computed based on the sensitivity and the specificity values. Therefore, the total accuracy can be considered an index that considers both specificity and sensitivity and thus the interpretation of the results through the total accuracy allows the real situation to be explained in a better way. For this reason, while the results are interpreted, it has been focused on the total accuracy values.

As in the sensitivity values, the total accuracy values have also been affected by the number of group, sample size and effect size. Total accuracy values have estimated between 60.02% and 99.98% in general. Slightly decreases have been observed in total accuracy estimates when the number of group being compared has been increased. For example, when samples have been taken from Normal distributed population under the sample size of 30 and effect size of 0.60, the total accuracy values have been estimated between 83.01% and 91.40% for $k=30$. When k increased to 80, the total accuracy values have been estimated between 78.23% and 89.06% and it has been estimated between 74.31% and 87.15% when k increased to 200. However, the total accuracy of the ANOM has not fallen below 60.02% regardless of experimental conditions.

As the effect size increased, the total accuracy values increased as well. For instance, when samples have

been taken from Normal distributed population and $k=200$, for $\delta=0.30$ the total accuracy values have varied between 60.04% and 81.24%. When δ increased to 0.60 and 0.90 the total accuracy values have changed between 60.40% and 94.47%, and 61.98% and 99.92% respectively. For $\delta=1.2$, the total accuracy values have changed between 66.58% and 99.98% under the same experimental conditions.

Percentage of different groups (d) has also an important impact on the total accuracy estimates. Increases in the percentage of different groups have caused slightly decreases in the total accuracy estimates. For example, when samples are taken from Normal distributed populations for $k=200$ and $n=5$, the total accuracy values for $d=20\%$, $d=30\%$ and $d=40\%$ have been varied between 80.01% and 83.27 %, 70.02% and 74.93 %, and 60.04 % and 66.58%. When sample size is increased to 30, the total accuracy values have been estimated between 80.42% and 99.95% for $d=20\%$, 70.65 % and 99.93% for $d=30\%$, 60.87% and 99.92% for $d=40\%$ under the same conditions. Total accuracy estimates have found between 60.02% and 99.98% in general and these estimates have not fell below the level of 60.00 % under none of the experimental conditions. For example, the total accuracy of the ANOM has been found as 60.02% even if sample size is 5, effect size is $=0.30$ and number of group is 200. That means, the ANOM technique can reveal the real situation with at least 60.00% accuracy regardless of experimental conditions. The total accuracy or correct classification rate of the ANOM has gotten close to 100.00% when sample size increased. Therefore, it is possible to conclude that the correct classification percentage of ANOM will be at least 60.00 % no matter which experimental conditions are considered. This is also an indication that the ANOM technique can accurately figure out a fairly large proportion of groups that are actually different / not different from the grand mean or exceed / not exceed upper or lower decision lines.

Results of Real Data Set

ANOM vs Tukey: Results of Tukey pairwise multiple comparison test have been presented in Table 2. Table 2 shows the results of an experiment where there are 80 wheat varieties. As can be seen from the results of Tukey pairwise comparison test, it is not easy (almost impossible) to determine superior or worse varieties due to a large number of varieties. At the same time, it is possible to say that the reliability of these results is highly dubious and questionable because of a large number of pairwise comparisons. For example, since there are 80 wheat varieties in this study, there will be 3160 pairwise comparisons ($k(k-1)/2=(80)(79)/2=3160$). And it will not be easy to determine superior and worse varieties among of 80 wheat varieties due to 3160 pairwise comparisons. That is why, Tukey (1992) recommended researchers to prefer Multiple Comparison with Mean Procedure (MCMP) which is proposed by Halperin *et al.* (1955) to All-Pairwise

Multiple Comparison Procedures (MCAP) for large number of group cases (k). It is because, the results of k comparisons in Multiple Comparison with Mean Procedures would be easier to comprehend than the result of $k(k-1)/2$ pairwise comparisons in All-pairwise Comparison Procedure when k is large. In other word, for this study, evaluating the results of 80 wheat varieties instead of evaluating the results of 3160 pairwise comparisons ($(80)(79)/2=3160$ pairwise comparisons) definitely will be easier and reliable. As can be seen in table 2, each variety has been represented by a different capital letter and due to a large number of varieties there are many capital letters have been used to determine different varieties. Some varieties have shared one or more capital letters to show non-statistical significant differences among those varieties. Although classical pairwise comparison procedures are not appropriate for the cases where there are a large number of treatment groups or a large number of pairwise comparisons, on the other hand, the ANOM technique can be used efficiently in determining superior and worse varieties or similar varieties easily regardless of number of treatment means. It is because ANOM technique is a graphical technique and it is very easy to see different varieties in terms of interested variable (Nelson *et al.*, 2005; Wu and Liao, 2004; Mendes and Yiğit, 2013; Pallmann and Hothorn, 2016). Homa (2007) reported that the ANOM graphs are very useful in order to see if any group mean is significantly different from the overall mean. ANOM chart which has been generated for this data set has been given in figure 21. As it is seen from figure 21, mean of some varieties fall outside the decision lines while some of them not. If the lines of the varieties fall outside the lower and upper decision lines then it is concluded that there is a statistically significant difference between varieties means and overall mean. Therefore, from the figure 21 it is very easy and clear to see different varieties. And it is also possible to evaluate practical significance of the observed differences along with statistical significance as well. As it is seen in figure 21, some means have exceeded upper and lower decision lines. That shows both different and superior or worse varieties. In other way, average of some varieties (i.e. number 11, 48, 3, and 29) are significantly greater / smaller than that of the overall or population mean. For example, the varieties number 11, 12, 41, 48, 54, 58, 61, 63, 71, 92, 93, 94, 95, 96 and 97 have exceeded upper decision line. Therefore, these fifteen varieties are superior because the highest grain yields have obtained in varieties 11, 12, 41, 48, 54, 58, 61, 63, 71, 92, 93, 94, 95, 96 and 97 when compared to the rest. And these varieties might be selected for the further studies (i.e. breeding studies) (CIMMYT, 1966-2016). Likewise, the varieties number 3, 21, 22, 25, 28, 29, 31, 32, 44, 73, 76, 80, 82, 83, 90 and 98 have exceeded lower decision line. That means the averages of these varieties are significantly lower than that of the population average and it will not be appropriate to choose

these varieties for further studies. Since the lowest grain yields have obtained in varieties 3, 21, 22, 25, 28, 29, 31, 32, 44, 73, 76, 80, 82, 83, 90 and 98, these varieties are the worst among 80 varieties in terms of grain yield and these varieties will not select for the further studies. Thus, the varieties number 11, 12, 41, 48, 54, 58, 61, 63, 71, 92, 93, 94, 95, 96 and 97 have higher grain yield than that of the varieties 3, 21, 22, 25, 28, 29, 31, 32, 44, 73, 76, 80, 82, 83, 90 and 98.

Conclusion: Many studies in practice aim at determining superior and worse treatments among a large number of treatment groups. To achieve this, Pairwise Comparison Procedures like LSD, Duncan, Tukey etc are commonly used. However, classical pairwise comparison procedures are not appropriate for the cases where there are a large number of treatment groups or a large number of pairwise comparisons. The large number of group being compared and multiplicity problems are two big challenges for ANOVA and Pairwise Comparison Procedures. On the other hand, ANOM technique can be used efficiently to determine superior, worse and similar varieties when number of treatment groups are very high. Nelson and Dudewicz (2002) describes it as a “**graphical analogue of ANOVA**”. Likewise, Pallmann and Hothorn (2016) reported that the ANOM is a graphical method for presenting multiple group comparisons with an overall or grand mean. However, ANOM technique is not only used for comparing group means but also used for comparing variances, proportions, correlation coefficients, and other location and scale measures across several groups (Wludyka and Nelson, 1997; Kumar and Rao, 1998; Rao and Kumar, 2002; Nelson *et al.*, 2005; Rao and Deva Raaj, 2006).

ANOM technique can also be used efficiently as a multiple comparison test especially for the cases where number of treatment groups is large (Fritsch and Hsu, 1997). Since the ANOM is a graphical technique, it presents the results visually that provides a quick way for researchers and readers to evaluate both practical and statistical significant differences between the treatment groups and the overall mean. The ANOM graphs also enable the researchers to figure out superior and worse treatment groups easily when compared to the overall or population mean. This is a very big advantage of ANOM over ANOVA-F and Classical Pairwise Multiple Comparison Procedures (CPMP) especially when there are a large number of treatment groups. Recently, due to these kind of advantageous of ANOM technique, it has become increasingly popular among the scientists and it has been implemented in many statistical package programs such as SAS (1999), MINITAB 18 (2017), JMP 12 (2016), and R (Pallmann and Hothorn, 2016).

With the ANOM, it is possible to identify different or superior treatment groups among the many treatment groups easily. The results of Monte Carlo Simulation study

and illustrative example to compare the ANOM and Tukey procedures in terms of their correctly identify superior and worse treatments ability demonstrated that the ANOM procedure had a great advantage over the Tukey Procedure especially when there are a large number of treatment groups. As a result, it is highly possible to conclude that the ANOM procedure can be used efficiently as a multiple comparison procedure especially to determine superior and worse treatment group means when compared to population mean along with comparing independent means, variances, rates, and other local and scales measures.

Acknowledgements: The authors thank Dr. Cem Ömer Egesel and Dr. Fatih Kahrıman for providing the data set to illustrative the ANOM.

REFERENCES

- Anonymous, (1994). FORTRAN Subroutines for Mathematical Applications. IMSL Math/Library. 1-2. Visual Numerics, Inc., Houston, USA.
- Duncan, D. B. (1955). Multiple range and multiple F tests. *Biometrics* 11(1): 1-42.
- Flint-Garcia, S. A., Bodnar, A. L. and M. P. Scott (2009). Wide variability in kernel composition, seed characteristics, and zein profiles among diverse maize inbreds, landraces, and teosinte. *Theor. Appl. Genet.* 119(6): 1129-1142.
- Fritsch, K. S. and J. C. Hsu (1997). Multiple comparisons with the mean. In N. Balakrishnan (Ed.), *Advances in statistical decision theory and applications*. Birkhäuser; Boston (USA). pp. 189-204.
- García-Lara, S. and D. J. Bergvinson (2013). Identification of maize landraces with high level of resistance to storage pests *Sitophilus zeamais* motschulsky and *Prostephanus truncatus* horn in Latin America. *Rev. Fitotec. Mex.* 36(3A): 347-356.
- Halperin, M., S. W. Greenhouse, J. Cornfield and J. Zalkar (1955). Tables of Percentage Points for the Studentized Maximum Absolute Deviate in Normal Samples. *J. Am. Stat. Assoc.* 50(269): 185-195.
- Homa, K. (2007). Analysis of means used to compare providers' referral patterns. *Qual. Manag. Health Care* 16(3): 256-264.
- International Maize and Wheat Improvement Center (CIMMYT), 1966-2016, El Batan, Mexico.
- JMP®, Version 12. (2016). SAS Institute Inc., Cary, NC.
- Karaköy T., F. S. Baloch, F. Toklu and H. Özkan (2014). Variation for selected morphological and quality-related traits among 178 faba bean landraces collected from Turkey. *Plant Genet. Resour.-C*, 12(1): 5-13
- Kumar, M.P. and Rao, C.V. (1998). ANOM-type graphical method for testing the equality of several variances. *Commun. Stat. Simul. Comput.* 27(2): 459-468.
- Mendes, M. and Yiğit, S. (2013). Comparison of ANOVA-F and ANOM tests with regard to type I error rate and test power. *J. Stat. Comput. Simul.* 83(11): 2093-2104.
- MINITAB 18, Statistical Software (2017). State College, PA: Minitab, Inc.
- Munyiri, S. W., Mugo, S. N., Otim, M., Tefera, T., Beyene, Y., Mwololo, J. K. and Okori, P. (2013). Responses of tropical maize landraces to damage by *Chilopartellus* stem borer. *Afr. J. Biotechnol.* 12(11): 1229-1235.
- Nelson, P.P. (1983). The analysis of means for balanced experimental designs. *J. Qual. Technol.* 15(1): 45-54.
- Nelson, P. R. (1985). Power Curves for the Analysis of Means. *Technometrics.* 27(1): 65-73.
- Nelson, P.P. (1989). Multiple comparisons of means using simultaneous confidence intervals. *J. Qual. Technol.* 21(4): 232-241.
- Nelson, P. R. (1993). Additional Uses for the Analysis of Means and Extended Tables of Critical Values. *Technometrics.* 35(1): 61-71.
- Nelson, P.P. and E. J. Dudewicz (2002). Exact analysis of means with unequal variances. *Technometrics.* 44(2): 152-160.
- Nelson, P.P., P. S. Wludyka, P.S. and K. A. F. Copeland (2005). *The analysis of means: A graphical method for comparing means, rates and proportions*. 1st Ed. SIAM; Philadelphia (USA). 239 p.
- Ott, E. R. (1983). Analysis of means- A graphical procedure. *J. Qual. Technol.* 15(1): 10-18.
- Pallmann, P. and L. A. Hothorn (2016). Analysis of means: a generalized approach using R. *J. Appl. Stat.* 43(8): 1541-1560.
- Rao, C. V. and V. J. Deva Raaj (2006). ANOM for testing the equality of several intercepts in a k-sample regression model. *Statistical Methods.* 8(1), 60-72.
- Rao, C.V. and M. P. Kumar (2002). ANOM-type graphical methods for testing the equality of several correlation coefficients. *Gujarat Stat. Rev.* 29: 47-56.
- SAS, (1999). Institute Inc., SAS OnlineDoc®, Version 8, Cary, NC.
- Sharma, R. C., S. Rajaram, S. Alikulov, Z. Ziyaev, S. Hazratkulova, M. Khodarahami, S. M. Nazeri, S. Belen, Z. Khalikulov, M. Mosaad, Y. Kaya, M. Keser, Z. Eshonova, A. Kokhmetova, M. G. Ahmedov, M. R. Jalalkamali, and A. I. Morgounov (2013). Improved winter wheat

- genotypes for Central and West Asia. *Euphytica*. 190(1): 19-31.
- Schilling, E. G. (1973). A Systematic Approach to the Analysis of Means. *J. Qual. Technol.* 5(4): 92-108.
- Wietholter, P., M. C. de MeloSerenio, T. de Freitas Terra, S. D. dos Anjos-e-Silva and J. F. Barbosa Neto (2008). Genetic variability in corn landraces from Southern Brazil. *Maydica*. 53(2): 151-159.
- Williams, L. J. and H. Abdi (2010). Fisher's least significance difference (LSD) test. In N. J. Salkind (Ed.), *Encyclopedia of Research Design*. SAGE; London (UK). pp. 491-494.
- Wludyka, P. S. and P. R. Nelson (1997). An analysis-of-means-type test for variances from normal populations. *Technometrics* 39(3): 274-285.
- Wu, S. F. and B.X. Liao (2004). A simulation study of multiple comparisons with the average under heteroscedasticity. *Commun. Stat. Simul. Comput.* 33(3): 639-659.
- Tukey, J. W. (1949). Comparing individual means in the analysis of variance. *Biometrics* 5(2): 99-114.
- Tukey, J. W. (1992). Where should multiple comparisons go next?. In F. M. Hoppe Ed., *Multiple Comparisons, Selection, and Applications in Biometry: A Festschrift in Honor of Charles W. Dunnett*. Marcel Dekker; New York (USA). pp. 187-208.