

# A TEST BASED ON THE COMPUTATIONAL APPROACH FOR EQUALITY OF MEANS UNDER THE UNEQUAL VARIANCE ASSUMPTION

Esra Yiğit Gökpınar\*† and Fikri Gökpınar\*

Received 05 : 01 : 2011 : Accepted 01 : 11 : 2011

## Abstract

The classical F-test to compare several populations means depends on the assumption of homogeneity of variances of the population and on normality. When these assumptions - especially the equality of variance - is dropped, the classical F-test fails to reject the null hypothesis even if the data actually provide strong evidence for it. This can be considered a serious problem in some applications especially when the sample sizes are not large. To deal with this problem, a number of tests are available in the literature. Recently Pal, Lim and Ling (*A computational approach to statistical inferences*, J. Appl. Probab. Stat. **2**(1), 13–35, 2007) developed a computational technique, called the Computational Approach Test (CAT), which looks similar to a parametric bootstrap for hypothesis testing. Chang and Pal (*A revisit to the Behren-Fisher Problem: Comparison of five test methods*, Communications in Statistics - Simulation and Computation **37**(6), 1064–1085, 2008) applied CAT to test the equality of two population means when the variances are unknown and arbitrary. In this study we apply a developed CAT to test the equality of  $k$  population means when the variances are unequal. Also the Brown-Forsythe, Weerahandi's Generalized F, Parametric Bootstrap and Welch tests are recalled and a simulation study performed to compare these tests according to type one errors and powers in different combinations of parameters and various sample sizes.

**Keywords:** Brown-Forsythe Test, Computational Approach Test, Generalized F test, Parametric Bootstrap Test, Classic F Test, Welch Test.

*2000 AMS Classification:* 62 F 03, 62 F 40.

---

\*Gazi University, Faculty of Science, Department of Statistics, Teknikokullar, Ankara, Turkey.  
E-mail: (E. Y. Gökpınar) [eyigit@gazi.edu.tr](mailto:eyigit@gazi.edu.tr) (F. Gökpınar) [fikri@gazi.edu.tr](mailto:fikri@gazi.edu.tr)

†Corresponding Author.

## 1. Introduction

In applied statistics an experimenter wants to compare two or more populations measured on independent samples. The classical F (*CF*) test is used under the assumption that the populations have normal distributions with the same variances. Bishop and Dudewicz [1] showed that the CF test is not robust when the population variances are unequal, especially if the sample sizes are not equal. Also Krutchkoff [8] and Lee and Ahn [9] showed that the empirical type I errors are much greater than the nominal level  $\alpha$ , especially when sample sizes are negatively related with their population variances.

Alternative methods are developed due to this problem. For some of these test statistics the distribution is not known and the  $p$ -value can be found by simulation [13, 14]. Approximate methods are used quite often with the development of computer technology. Also these tests have been applied to solve a number of problems when conventional methods are difficult to apply or fail to provide exact solutions. In practice, some exact procedures such as the CF, Welch (*W*) and Brown-Forsythe (*BF*) tests are widely used [2, 15].

Tsui and Weerahandi [11] generalized the conventional definition of the  $p$  value so that problems such as the Behrens-Fisher problem can be resolved. Weerahandi [12] defined the notion of generalized  $p$ -value for comparing the means of  $k$  populations when the variances are not equal. Krishnamoorthy *et al.* [7] proposed a Parametric Bootstrap test (*PB*); Xu and Wang [16, 17] developed a generalized F-test based on the generalized value  $p$  (*XW*).

Recently Pal *et al.* [10] developed a computational technique, called the Computational Approach Test (*CAT*) which looks similar to the parametric bootstrap for hypothesis testing. Chang and Pal [3] applied CAT to test the equality of two population means when the variances are unknown and arbitrary. Also Chang *et al.* [4] have demonstrated that for homoscedastic one-way ANOVA the CAT is as powerful as the classical F test.

In this paper we apply CAT to test the equality of  $k$  population means when the variances are unknown and arbitrary. In the following section, we describe the *W*, *BF*, generalized F (*GF*) due to Weerahandi [13], *PB* and *CAT* tests. Monte Carlo comparison studies by Gamage and Weerahandi [5], Gerami and Zahiden [6] and Yiğit and Gökpinar [18] showed that, out of these and other tests, only the *W*, *GF* and *PB* tests emerged satisfactory provided the sample sizes are moderate or large. So we chose these tests for the simulation study. These methods are compared according to type I errors and powers in different combinations of parameters and various sample sizes.

## 2. Tests for one way ANOVA

Let  $X_{i1}, \dots, X_{in_i}$  be a random sample from  $N(\mu_i, \sigma_i^2)$ ,  $i = 1, \dots, k$ . The problem of interest involves testing

$$(1) \quad H_0 : \mu_1 = \mu_2 = \dots = \mu_k \text{ against } H_A : \text{Not all } \mu_i \text{ s are equal, } i = 1, \dots, k.$$

If the  $\sigma_i^2$ 's are unequal, then the testing procedure given in equation (1) defines the standardized between-group sum of squares in equation (2) and the standardized error sum of squares in equation (3).

$$(2) \quad \tilde{S}_b = \tilde{S}_b(\sigma_1^2, \dots, \sigma_k^2) = \sum_{i=1}^k \frac{n_i \bar{X}_i^2}{\sigma_i^2} - \frac{\left(\sum_{i=1}^k \frac{n_i \bar{X}_i}{\sigma_i^2}\right)^2}{\sum_{i=1}^k \frac{n_i}{\sigma_i^2}},$$

$$(3) \quad \tilde{S}_e = \sum_{i=1}^k \frac{(n_i - 1) S_i^2}{\sigma_i^2},$$

where  $S_i^2$  is the unbiased estimator of  $\sigma_i^2$  and  $\bar{X}_i$  is the estimator of  $\mu_i$ . Most of the test statistics to test the equality of means under heteroscedasticity are based on the standardized between-group sum of squares and standardized error sum of squares. In the rest of this section some test statistics are briefly introduced.

*The Welch Test*

If  $w_i = \frac{n_i}{S_i^2}$ , Welch [15] gives the following test statistics [14]:

$$(4) \quad W = \frac{\tilde{S}_b (S_1^2, \dots, S_k^2) / (k - 1)}{1 + \frac{2(k-2)}{k^2-1} \sum_{i=1}^k \frac{1}{n_i-1} \left(1 - \frac{w_i}{\sum w_j}\right)^2} = \frac{\sum_{i=1}^k w_i [(\bar{X}_i - \bar{X})^2 / (k-1)]}{1 + \frac{2(k-2)}{k^2-1} \sum_{i=1}^k \frac{1}{n_i-1} \left(1 - \frac{w_i}{\sum w_j}\right)^2},$$

where  $\bar{X}$  is the estimator of the overall mean. If  $H_0$  is true, then the distribution of  $W$  is  $F_{k-1, f}$ , where

$$f = \frac{1}{\frac{3}{k^2-1} \sum_{i=1}^k \frac{1}{n_i-1} \left(1 - \frac{w_i}{\sum w_j}\right)^2}.$$

For a given level  $\alpha$ , and an observed value  $W_h$  of  $W$ , this test rejects the  $H_0$  in equation (1) whenever the  $p$ -value is given as  $P(F_{k-1, f} > W_h) < \alpha$ .

*The Brown-Forsythe Test*

Brown and Forsythe [2] give the following test statistics.

$$B = \frac{\sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2}{\sum_{i=1}^k \left(1 - \frac{n_i}{n}\right) S_i^2}.$$

If  $H_0$  is true, then the distribution of  $B$  is  $F_{k-1, v}$ , where

$$v = \frac{\left[\sum_{i=1}^k \left(1 - \frac{n_i}{n}\right) S_i^2\right]^2}{\sum_{i=1}^k \frac{\left(1 - \frac{n_i}{n}\right)^2 S_i^4}{n_i - 1}}.$$

For a given level  $\alpha$  and an observed value  $B_h$  of  $B$ , this test rejects the  $H_0$  in equation (1) whenever the  $p$ -value is given as  $P(F_{k-1, v} > B_h) < \alpha$ .

*The Weerahandi's Generalized F test*

The sample variances (MLEs) of the  $k$  populations are denoted by  $S_i^2$ , where  $S_i^2 = \frac{1}{n_i} \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2$ . Define

$$B_j = \frac{\left(\sum_{i=1}^j \frac{n_i S_i^2}{\sigma_i^2}\right)}{\left(\sum_{i=1}^{j+1} \frac{n_i S_i^2}{\sigma_i^2}\right)} \sim \text{Beta} \left( \sum_{i=1}^j \frac{(n_i - 1)}{2}, \frac{(n_{j+1} - 1)}{2} \right), \quad j = 1, \dots, k - 1$$

where  $\tilde{S}_e$  and  $B_j$  are all independent random variables. Note that the random variables  $\frac{n_i S_i^2}{\sigma_i^2}$  can be expressed as

$$\begin{aligned} \frac{n_1 S_1^2}{\sigma_1^2} &= \tilde{S}_e B_1 B_2 \dots B_{k-1}, \\ \frac{n_i S_i^2}{\sigma_i^2} &= \tilde{S}_e (1 - B_{i-1}) B_i \dots B_{k-1} \text{ for } i = 2, \dots, k - 1, \\ \frac{n_k S_k^2}{\sigma_k^2} &= \tilde{S}_e (1 - B_{k-1}). \end{aligned}$$

Therefore, the generalized  $p$  value can be expressed as

$$(5) \quad p = 1 - E \left( H_{k-1, n-k} \left\{ \frac{n-k}{k-1} \tilde{s}_b \left[ \frac{n_1 s_1^2}{B_1 B_2 \dots B_{k-1}}, \frac{n_2 s_2^2}{(1-B_1) B_2 \dots B_{k-1}}, \dots, \frac{n_k s_k^2}{(1-B_{k-1})} \right] \right\} \right),$$

where  $H_{k-1, n-k}$  is the cumulative distribution function of the  $F$ -distribution with  $k-1$  and  $N-k$  degrees of freedom. This test rejects the  $H_0$  in equation (1) whenever  $p < \alpha$  [12].

#### The parametric bootstrap test

The PB approach is defined as follows. In the case where the population variances  $\sigma_i^2$ 's are unknown; a test statistic can be obtained by replacing  $\sigma_i^2$  in equation (2) by  $S_i^2$  and is given by

$$(6) \quad \tilde{S}_b(S_1^2, \dots, S_k^2) = \sum_{i=1}^k \frac{n_i \bar{X}_i^2}{S_i^2} - \frac{\left( \sum_{i=1}^k \frac{n_i \bar{X}_i}{S_i^2} \right)^2}{\sum_{i=1}^k \frac{n_i}{S_i^2}}.$$

As the test statistic in equation (6) is location invariant, without loss of generality, we can take the common mean to be zero [7].

Let  $\bar{X}_{Bi} \sim Z_i \left( \frac{S_i}{\sqrt{n_i}} \right)$ , where  $Z_i$  is a standard normal random variable and  $S_{Bi}^2 \sim S_i^2 \chi_{n_i-1}^2 / (n_i - 1)$ . Then the PB pivot variable can be obtained by replacing  $\bar{X}$ ,  $S_i^2$  in equation (6) by  $\bar{X}_{Bi}$ ,  $S_{Bi}^2$ , and is given by

$$S_{bB} = \sum_{i=1}^k \frac{n_i}{S_{Bi}^2} \bar{X}_{Bi}^2 - \frac{\left[ \sum_{i=1}^k \frac{n_i \bar{X}_{Bi}}{S_{Bi}^2} \right]^2}{\sum_{i=1}^k \frac{n_i}{S_{Bi}^2}} = \sum_{i=1}^k \frac{Z_i^2 (n_i - 1)}{\chi_{n_i-1}^2} - \frac{\left[ \sum_{i=1}^k \frac{\sqrt{n_i} Z_i (n_i - 1)}{S_i \chi_{n_i-1}^2} \right]^2}{\sum_{i=1}^k \frac{n_i (n_i - 1)}{S_i^2 \chi_{n_i-1}^2}}.$$

$H_0$  is rejected if  $P\{\tilde{S}_{bB}(Z_i, \chi_{n_i-1}^2; s_i^2) > \tilde{s}_b\} < \alpha$ .

#### The test based on CAT

In this section, we give a test procedure based on CAT for one-way ANOVA when the population variances are unequal. The main point in using the CAT approach in this problem is to obtain the Restricted Maximum Likelihood (RML) estimators of the parameters. The RML estimators can be obtained as follows.

The RML method gives the ML estimator of  $\mu$  and  $\sigma_i^2$  based on the assumption that  $H_0$  ( $\mu_1 = \mu_2 = \dots = \mu_k = \mu$ ). The likelihood function of the sample  $(X_{11}, \dots, X_{kn_k})$  can be given as below.

$$(7) \quad L = \frac{1}{(2\pi)^{\sum_{i=1}^k n_i/2} \prod_{i=1}^k \sigma_i^{n_i}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(X_{ij} - \mu)^2}{\sigma_i^2} \right\},$$

$$L^* = \ln(L) = -\sum_{i=1}^k \frac{n_i}{2} \ln(2\pi) - \sum_{i=1}^k n_i \ln(\sigma_i) - \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(X_{ij} - \mu)^2}{\sigma_i^2}$$

$$= -\sum_{i=1}^k \frac{n_i}{2} \ln(2\pi) - \sum_{i=1}^k \frac{n_i}{2} \ln(\sigma_i^2) - \frac{1}{2} \frac{\sum_{i=1}^k (n_i S_i^2 + n_i (\bar{X}_i - \mu)^2)}{\sigma_i^2}.$$

By differentiating Equation (7) with respect to  $\mu$  and  $\sigma_i^2$ , the following equations can be obtained.

$$\hat{\mu} = \frac{\sum_{i=1}^k \frac{n_i}{\sigma_i^2} \bar{X}_i}{\sum_{i=1}^k \frac{n_i}{\sigma_i^2}},$$

$$\hat{\sigma}_i^2 = S_i^2 + \left( \frac{\sum_{i=1}^k \frac{n_i}{\sigma_i^2} D_{ij}}{\sum_{i=1}^k \frac{n_i}{\sigma_i^2}} \right)^2,$$

where  $D_{ij} = \bar{X}_i - \bar{X}_j$  and  $S_i^2$  is the sample variance (MLEs) of the  $i^{\text{th}}$  population.

To apply the developed CAT we first express  $H_0$  in terms of a suitable scalar  $\eta$ . Define  $\eta$  as

$$\eta = \eta(\mu_1, \dots, \mu_k) = \sum_{i=1}^k n_i (\mu_i - \bar{\mu})^2, \quad \bar{\mu} = \sum_{i=1}^k n_i \mu_i / n, \quad n = \sum_{i=1}^k n_i.$$

It is seen that testing  $H_0$  against  $H_A$  in equation (1) is equivalent to testing  $H_0^* : \eta = 0$  against  $H_A^* = \eta > 0$ . To test the hypothesis in equation (1) we use  $\eta$  as a test statistic. If  $H_0$  is true then  $\eta = 0$  otherwise  $\eta$  becomes larger than 0.

The test procedure is as follows:

- 1) The ML estimators of the parameter are  $\hat{\mu}_{i(ML)} = \bar{X}_i$ ,  $\hat{\sigma}_{i(ML)}^2 = S_i^2$ . Also the test statistic is  $\hat{\eta}_{ML} = \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2$
- 2) If  $H_0$  is true,  $X_{ij} \sim N(\mu, \sigma_i^2)$  ( $1 \leq j \leq n_i, 1 \leq i \leq k$ ),

$$\hat{\mu}_{RML} = \frac{\sum_{i=1}^k (n_i \bar{X}_i / \hat{\sigma}_{i(RML)}^2)}{\sum_{i=1}^k (n_i / \hat{\sigma}_{i(RML)}^2)},$$

$$\hat{\sigma}_{i(RML)}^2 = S_i^2 + \left( \frac{\sum_{i=1}^k \frac{n_i}{\hat{\sigma}_{i(RML)}^2} D_{ij}}{\sum_{i=1}^k \frac{n_i}{\hat{\sigma}_{i(RML)}^2}} \right)^2, \quad i = 1, \dots, k.$$

As seen from these equations, the RML estimates of the  $\mu$  and  $\sigma_i^2$  parameters have no closed forms. Therefore, the RML estimates of these parameters can be obtained using iteration.

- 3) Generate  $\bar{X}_i$   $M$  times such that

$$(8) \quad \bar{X}_i \sim N(\hat{\mu}_{RML}, \hat{\sigma}_{i(RML)}^2 / n_i).$$

- 4) For each replication of  $\bar{X}_i$ , calculate  $\hat{\eta}_{ML}^* = \sum_{i=1}^k n_i (\hat{\mu}_{i(ML)} - \hat{\mu}_{ML})^2$ .
- 5) Calculate  $p = \frac{\#(\hat{\eta}_{ML} > \hat{\eta}_{ML}^*)}{M}$ ; when  $p < \alpha$ ,  $H_0$  is rejected.

### 3. Simulation study

In this section we compare the CF, BF, W, GF, PB and CAT tests according to type I errors and powers in different combinations of parameters and sample sizes. We consider the balanced and unbalanced cases from smaller to larger sample sizes where  $k = 3$  and  $k = 5$  for comparing the tests. The values for the variances vary over a large range so that  $\sigma_1^2 < \dots < \sigma_k^2$  and  $\sigma_1^2 > \dots > \sigma_k^2$ .

For each combination of  $n_i$  and  $\sigma_i^2$  the rejection rate of each testing procedure is calculated and compared with the nominal level 0.05 when the means are all equal. To estimate the type I error rates of the CF, W and BF tests, we used simulation consisting

of 5000 runs for each of the sample sizes and parameter configurations. The CF, W and BF test statistics are calculated from these generated data and type I errors are estimated by the proportion of test statistics that exceed the critical values calculated from the distributions. To estimate the type I error rates of the GF, PB and CAT tests, we use a two-step simulation. For estimating the type I error rates of the CAT test we generate 5000 observed vectors in equation (8) and use 5000 runs for each observed vector to estimate the  $p$  value.

Finally the type I error rates of the CAT test are estimated by the proportion of the 5000  $p$ -values that are less than the nominal level  $\alpha$ . The type I error rates of the PB and GF tests are similarly estimated. In both cases of equal and unequal variances for  $k = 3$  and  $k = 5$  simulated type I error rates are given in Tables 1 and 2 respectively.

We observed the following from the numerical results in Tables 1 and 2.

As seen from Table 1 the new test based on CAT has empirical type I error rates close to 0.05 when  $k = 3$ . Also the BF, GF, PB, W tests are close to 0.05 in this situation. But the W, PB and GF tests appear to be more powerful than the CAT test when  $k = 3$  and the sample sizes are small ( $(n_1, n_2, n_3) = (7, 9, 11)$ ).

When the  $n_i$ 's and  $\sigma_i^2$ 's are in reverse order the CF test's type I error rate is getting far away from its nominal level. When the differences between the means of the groups increase, the power of the new CAT test is superior to the other test for almost every combination of sample sizes and population variances. When the group size is 5 and the  $n_i$ 's and  $\sigma_i^2$ 's are in reverse order the type I error rate of CF is approximately around 0.10. This is not a acceptable level for a type I error. Also the type I error rate of GF is greater than its nominal level, especially when the sample sizes are small. The type I error of the new test based on CAT is smaller than its nominal level especially if the  $n_i$ 's and  $\sigma_i^2$ 's are in reverse order. When we investigate the power of the tests we can easily see that the new test is superior to the other tests.

#### 4. Conclusion

In this simulation study for a range of choices of sample sizes and parameter configurations we compared the performance of the above tests for testing the equality of means of one-way ANOVA models under heteroscedasticity. The CF test is not an appropriate test for heteroscedasticity because its type I error rates exceed the nominal level. The type I errors of the new test based on CAT are close to the nominal level. The CAT test appears to be more powerful than the other tests under almost every situation except when  $k = 3$  and the sample sizes are small ( $(n_1, n_2, n_3) = (7, 9, 11)$ ).







## References

- [1] Bishop, T. A. and Dudewicz, E. J. *Heteroscedastic ANOVA*, Sankhya **43**B, 40–57, 1981.
- [2] Brown, M. B. and Forsythe, A. B. *The small sample behavior of some statistics which test the equality of several means*, Technometrics **16**, 129–132, 1974.
- [3] Chang, C. H. and Pal, N. *A revisit to the Behren-Fisher Problem: Comparison of five test methods*, Communications in Statistics-Simulation and Computation **37** (6), 1064–1085, 2008.
- [4] Chang, C. H., Pal, N. and Lim, W. K. *Comparing several population means: a parametric bootstrap method, and its comparison with usual ANOVA F test as well as ANOM*, Computational Statistics **25** (1), 71–95, 2010.
- [5] Gamage, J. ve Weerahandi, S. *Size performance of some tests in one-way ANOVA*, Communications in Statistics Simulations **27** (3), 625–640, 1998.
- [6] Gerami, A. and Zahedian, A. *Comparing the means of normal populations with unequal variances* (Proceedings of the 53rd Session of International Statistical Institute, Seoul, Korea, 2001).
- [7] Krishnamoorthy, K., Lu, F. and Thomas, M. *A parametric bootstrap approach for ANOVA with unequal variances: fixed and random models*, Computational Statistics and Data Analysis, **51**, 5731–5742, 2006.
- [8] Krutchkoff, R. G. *One-way fixed effects analysis of variance when the error variances may be unequal*, J. Statist. Comput. Simulation **30**, 259–271, 1988.
- [9] Lee, S. and Ahn, C. H. *Modified ANOVA for unequal variances*, Communications in Statistics Simulations **32**, 987–1004, 2003.
- [10] Pal, N., Lim, W. K. and Ling, C. H. *A computational approach to statistical inferences*, J. Appl. Probab. Stat. **2** (1), 13–35, 2007.
- [11] Tsui, K. and Weerahandi, S. *Generalized p-values in significance testing of hypotheses in the presence of nuisance parametres*, Journal of the American Statistical Association **84**, 602–607, 1989.
- [12] Weerahandi, S. *ANOVA under unequal error variances*, Biometrika **38**, 330–336, 1995.
- [13] Weerahandi, S. *Exact Statistical Method for Data Analysis* (Springer-Verlag, New York, 1995).
- [14] Weerahandi, S. *Generalized Inference in Repeated Measures: Exact Methods in MANOVA and Mixed Models* (Wiley, New York, 2004).
- [15] Welch, B. L. *On the comparison of several mean values: An alternative approach*, Biometrika **38**, 330–336, 1951.
- [16] Xu, L. and Wang, S. *A new generalized p-value for ANOVA under heteroscedasticity*, Statistics and Probability Letters **78**, 963–969, 2007.
- [17] Xu, L. and Wang, S. *A new generalized p-value and its upper bound for ANOVA under unequal errors variances*, Communications in Statistics Theory and Methods **37**, 1002–1010, 2007.
- [18] Yiğit, E. and Gökpınar, F. *A simulation study on tests for one-way ANOVA under the unequal variance assumption*, Commun. Fac. Sci. Univ. Ank. Series A **12**, 15–34, 2010.