

Research Article

A Nonparametric Two-Sample Wald Test of Equality of Variances

David Allingham¹ and J. C. W. Rayner^{2,3}

¹ Centre for Computer-Assisted Research Computation and Its Applications, School of Mathematical and Physical Sciences, University of Newcastle, Newcastle, NSW 2308, Australia

² Centre for Statistical and Survey Methodology, School of Mathematics and Applied Statistics, University of Wollongong, Wollongong, NSW 2522, Australia

³ School of Mathematical and Physical Sciences, University of Newcastle, Newcastle, NSW 2308, Australia

Correspondence should be addressed to David Allingham, david.allingham@newcastle.edu.au

Received 31 August 2011; Accepted 17 December 2011

Academic Editor: YanXia Lin

Copyright © 2011 D. Allingham and J. C. W. Rayner. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We develop a test for equality of variances given two independent random samples of observations. The test can be expected to perform well when both sample sizes are at least moderate and the sample variances are asymptotically equivalent to the maximum likelihood estimators of the population variances. The test is motivated by and is here assessed for the case when both populations sampled are assumed to be normal. Popular choices of test would be the two-sample F test if normality can be assumed and Levene's test if this assumption is dubious. Another competitor is the Wald test for the difference in the population variances. We give a nonparametric analogue of this test and call it the R test. In an indicative empirical study when both populations are normal, we find that when both sample sizes are at least 25 the R test is nearly as robust as Levene's test and nearly as powerful as the F test.

1. Introduction: Testing Equality of Variances for Two Independent Samples

In the two-sample problem, we are given two independent random samples X_{11}, \dots, X_{1n_1} and X_{21}, \dots, X_{2n_2} . The location problem attracts most attention. Assuming that the samples are from normal populations, the pooled t -test is used to test equality of means assuming equal variances and Welch's test can be used when equality of variances is suspect but normality is not. When normality is in doubt, the Wilcoxon test is often used.

The corresponding dispersion problem is of interest to confirm the validity of, for example, the pooled t -test, and when dispersion differences are of direct interest. For example, testing for reduced variability is of interest in confirming natural selection (see, e.g., [1, Section 5.5]) and if some processes are in control. In exploratory data analysis, it is

sensible to assess if one population is more variable than another. If it is, the cause may be that one population is bimodal and the other is not; the consequences of this in both the scenario and the model can then be explored in depth.

The study here introduces a new test of equality of variances. The asymptotic null distribution of the test statistic is χ_1^2 , but, depending on the populations sampled from, this may or may not be satisfactory for small to moderate sample sizes.

To assess this, and other aspects of the proposed test, an indicative empirical study is undertaken. We give comparisons when both populations sampled are normal, when both samples have the same sample size, and for 5% level tests. First, we derive a finite sample correction to the critical values based on the asymptotic null distribution when sampling from normal populations. Different corrections would be needed when sampling from different distributions. Next, we show that in moderate samples the new test is nearly as powerful as the F test when normality may be assumed and, finally, that it is nearly as robust as Levene's test when normality is in doubt. Moore and McCabe [2, page 519] claim that the " F test and other procedures for inference about variances are so lacking in robustness as to be of little use in practice." The new test gives a counterexample to that proposition.

We acknowledge that the new test is the most effective for at least moderate sample sizes. In the normal case, each random sample should have at least 25 observations, which is what we would expect of a serious study aiming at reasonable power that cannot be hoped for with samples of size 10 or so. See Section 4.

We are aware of more expansive comparative studies such as [3, 4]. Our goal here is not to emulate these studies but to show that the new test is competitive in terms of test size, robustness, and power.

In Section 2, the new test is introduced. Sections 3, 4, and 5 give the results of an empirical investigation when the populations sampled are assumed to be normal. In Section 3, we investigate test size, showing that the asymptotic χ^2 critical values should only be used for moderate to large sample sizes. For smaller sample sizes, a finite sample correction to the asymptotic 5% χ^2 critical value is given. This results in actual test sizes between 4.6% and 5.3%.

In Section 4, we show that when normality holds, the new test is not as powerful as the Levene test for smaller sample sizes but overtakes it for moderate samples of about 25. The new test is always inferior to the optimal F test. However, the R test has power that approaches that of the F test, being at least 95% that of the F test throughout most of the parameter space in samples of at least 80.

In Section 5, we show that if we sample from t -distributions with varying degrees of freedom, the F test is highly nonrobust for small degrees of freedom, as is well known for fat-tailed distributions. When both populations sampled are gamma distributions, chosen to model skewness differences from normality, or t -distributions, chosen to model kurtosis differences, the new test performs far better than the F test, and is competitive with the Levene test.

That the new test gives a good compromise between power and robustness and is valid when normality does not hold are strong reasons for preferring the new test for sample sizes that are at least moderate, and normality is dubious.

2. Competitor Tests for Equality of Variance

Initially we assume that we have two independent random samples X_{i1}, \dots, X_{in_i} from normal populations, $N(\mu_i, \sigma_i^2)$ for $i = 1$ and 2. We wish to test $H : \sigma_1^2 = \sigma_2^2$ against the

alternative $K : \sigma_1^2 \neq \sigma_2^2$, with the population means being unknown nuisance parameters. If $S_i^2 = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 / (n_i - 1)$, in which $\bar{X}_i = \sum_{j=1}^{n_i} X_{ij} / n_i$, $i = 1, 2$, then the S_i^2 are the unbiased sample variances, and the so-called F test is based on their quotient, $S_2^2 / S_1^2 = F$, say. It is well known and will be confirmed yet again in Section 5 that the null distribution of F , F_{n_1-1, n_2-1} , is sensitive to departures from normality. If $F_{n_1-1, n_2-1}(x)$ is the cumulative distribution function of the F_{n_1-1, n_2-1} distribution, and if c_p is such that $F_{n_1-1, n_2-1}(c_p) = p$, then the F test rejects H at the 100 $\alpha\%$ level when $F \leq c_{\alpha/2}$ and when $F \geq c_{(1-\alpha/2)}$.

Common practice when normality is in doubt is to use Levene's test or a nonparametric test such as Mood's test. Levene's test is based on the ANOVA F test applied to the residuals. There are different versions of Levene's test using different definitions of residual. The two most common versions use residuals based on the group means, $|X_{ij} - \bar{X}_i|$, and the group medians, $|X_{ij} - \tilde{X}_i|$, in which \tilde{X}_i is the median of the i th sample. The latter is called the Brown-Forsythe test. Again it is well known that these tests are robust in that when the population variances are equal but the populations themselves are not normal, they achieve levels close to nominal. However, this happens at the expense of some power. As the empirical study in this paper is intended to be indicative rather than exhaustive, we will henceforth make comparisons only with the Levene test, based on a test statistic we denote by L .

We now construct a new test that we call the R test. For univariate parameters θ , a Wald test statistic of $H : \theta = \theta_0$ against the alternative $K : \theta \neq \theta_0$ is based on $\hat{\theta}$, the maximum likelihood estimator of θ , usually via the test statistic $(\hat{\theta} - \theta_0)^2 / \text{est var}(\hat{\theta})$, where $\text{est var}(\hat{\theta})$ is the asymptotic variance of $\hat{\theta}$ evaluated when $\theta = \hat{\theta}$. Under the null hypothesis, this test statistic has an asymptotic χ_1^2 distribution. As well as being equivalent to the likelihood ratio test, the F test is also a Wald test for testing $H : \theta = \sigma_2^2 / \sigma_1^2 = 1$ against $K : \theta \neq 1$.

Rayner [5] derived the Wald test for testing $H : \theta = \sigma_2^2 - \sigma_1^2 = 0$ against $K : \theta \neq 0$. The test statistic is

$$\frac{(S_1^2 - S_2^2)^2}{2S_1^4 / (n_1 + 1) + 2S_2^4 / (n_2 + 1)} = W, \quad \text{say.} \tag{2.1}$$

Being a Wald test, the asymptotic distribution of W is χ_1^2 , while its exact distribution is not immediately obvious. However, W is a 1-1 function of F , so the two tests are equivalent. Since the exact distribution of F is available, the F test is the obvious test to use.

In W , the variances $\text{var}(S_j^2)$ are estimated optimally using the Rao-Blackwell theorem. This depends very strongly on the assumption of normality. If normality is in doubt, then we can estimate $\text{var}(S_1^2 - S_2^2)$ using results given, for example, in [6]. For a random sample, Y_1, \dots, Y_n with population and sample central moments μ_r and $m_r = \sum_{j=1}^n (Y_j - \bar{Y})^r / n$, $r = 2, 3, \dots$, [6] gives that

$$E[m_r] = \mu_r + O(n^{-1}), \quad \text{var}(m_2) = \frac{(\mu_4 - \mu_2^2)}{n + O(n^{-2})}. \tag{2.2}$$

Applying [6, 10.5] to the numerator of $\text{var}(m_2)$, μ_2^2 may be estimated to $O(n^{-1})$ by m_2^2 , or, equivalently, by $nm_2^2 / (n - 1) = S^4$, where S^2 is the unbiased sample variance. It follows

that, to order $O(n^{-2})$, $\text{var}(m_2)$ may be estimated by $(m_4 - m_2^2)/n$. We thus propose a robust alternative to W , given by

$$\frac{(S_1^2 - S_2^2)^2}{(m_{14} - S_1^4)/n_1 + (m_{24} - S_2^4)/n_2} = R, \quad \text{say,} \quad (2.3)$$

in which m_{i4} are the fourth central sample moments for the i th sample, $i = 1, 2$. We call the test based on R the R test. As the sample sizes increase, the distributions of the standardised sample variances approach standard normality, the denominator in R will approximate $\text{var}(S_1^2 - S_2^2)$, and R will have asymptotic distribution χ_1^2 . Thus, if c_α is the point for which the χ_1^2 distribution has weight α in the right hand tail, then the R test rejects H at approximately the 100 α % level when $R \geq c_\alpha$.

We emphasise that although the motivation for the derivation of R is under the assumption of sampling from normal populations, it is a valid test statistic for testing equality of variances no matter what the populations sampled.

If the sample variances are equal to or asymptotically equivalent to the maximum likelihood estimators of the population variances, as is the case when sampling from normal populations, then the R test is a Wald test for equality of variances in the sense described above. Since it does not depend on any distributional assumptions about the data, it can be thought of as a nonparametric Wald test. As such, it can be expected to have good properties in large samples.

We note that all the above test statistics are invariant under transformations $a(X_{ij} - b_i)$, for constants a , b_1 , and b_2 and for $j = 1, \dots, n_i$ and $i = 1, 2$.

The next three sections report an empirical study when the distributions sampled are assumed to be normal. As this is an indicative study, we fix the samples sizes to be equal, $n_1 = n_2 = n$, say, and the significance level to be 5% throughout.

3. Test Size under Normality

Under the null hypothesis, the distribution of F is known exactly, the distribution of L is known approximately, and, as mentioned above, the distribution of R is known asymptotically. In analysing data, these distributions are used to determine P values and critical values. We now investigate their use in determining test size, the probability of rejecting the null hypothesis when it is true.

Two empirical assessments of test size will now be undertaken. Since the test statistics are scale invariant, it is sufficient under the null hypothesis to take both population variances to be one.

In the first assessment, we assume normality. For various values of the common sample size n , we estimate the 5% critical points for each test by generating 100,000 pairs of random samples of size n , calculating the test statistics, ordering them, and hence identifying the 0.95th percentile. The estimated critical points of R approach the χ_1^2 5% critical point 3.841. These estimated critical points will subsequently be used in the power study to give tests with test the size of exactly 5%.

To see the extent of the error caused by using the asymptotic critical point 3.841, Figure 1 gives the proportion of rejections in 100,000 pairs of random samples for sample sizes up to 100. For $n = 10$, the proportion of rejections is nearly 20% and although for $n = 40$ this has dropped to nearly 7%, most users would hope for observed test sizes closer to 5%.

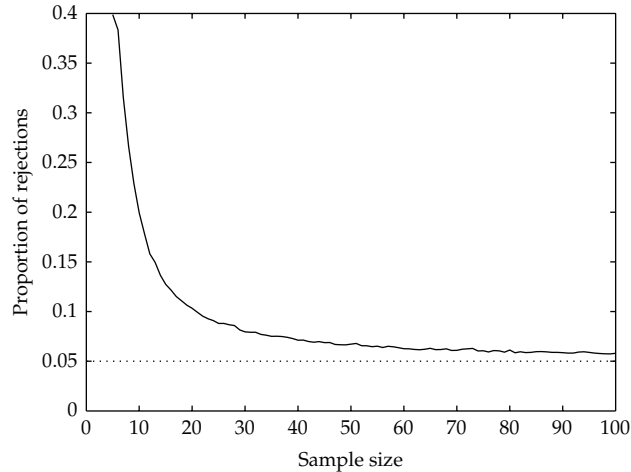


Figure 1: Proportion of rejections of the R test using the χ_1^2 5% critical point 3.841 for sample sizes up to 100.

The application of the test is improved by the use a Bartlett-type correction. This was found by plotting the estimated 5% critical points against n and using standard curve fitting techniques to find that the exact critical points were well approximated between $n = 10$ and 100 by $c(n, 0.05) = 3.84146(1.339 - 4.953/\sqrt{n} + 24.171/n)$. For larger n , it is sufficient to use the asymptotic 5% value 3.84146, the error being at most 0.8%. We checked the exact probabilities of rejection under the null hypothesis when applying the test with critical value $c(n, 0.05)$, and all were between 4.6% and 5.3%.

For levels other than 5%, and when sample sizes are unequal, further empirical work needs to be done to find critical values. However, as this study was intended to be indicative, we leave extensive tabulation of critical values for another time.

4. Power under Normality

For the F , Levene, and R tests, we estimate the power as the proportion of rejections from 100,000 pairs of random samples of size n when the first sample is from an $N(0, 1)$ population and the second is from an $N(0, \sigma^2)$ population with $\sigma^2 \geq 1$. To compare like with like, we use estimated critical values that give exact 5% level tests. It is apparent that for sample sizes less than about 20 the Levene test is more powerful than the R test and that for a sample size between approximately 20 and 30 the R test takes over from the Levene test; thereafter the R test is always more powerful than the Levene test. This is shown in Figure 2.

Both the Levene and R tests are always less powerful than the F test. This is explored in Figure 3 that compares the Levene test to the F test in the left hand panel and the R test to the F test in the right hand panel. What is given is a contour plot of the regions in which the ratios of the power of the stated test to the F test are less than 95%, between 95% and 99.99%, and greater than 99.99%. Generally, for any given n and σ^2 , it is clear that the power of the Levene test is at most that of the R test. For example, it appears that for $n_1 = n_2 = 60$ approximately the power of the R test is always at least 95% of that of the F test, whereas there is a considerable region where the power of the Levene test are less than 95% of that of the F test.

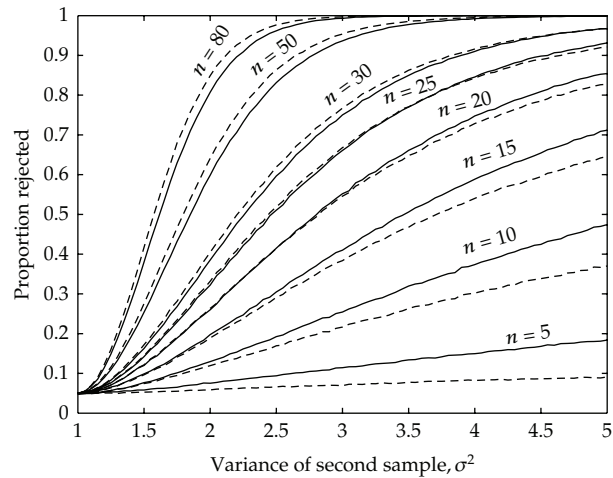


Figure 2: Power of the 5% level L test (solid line) and R test (dashed line) for various sample sizes.

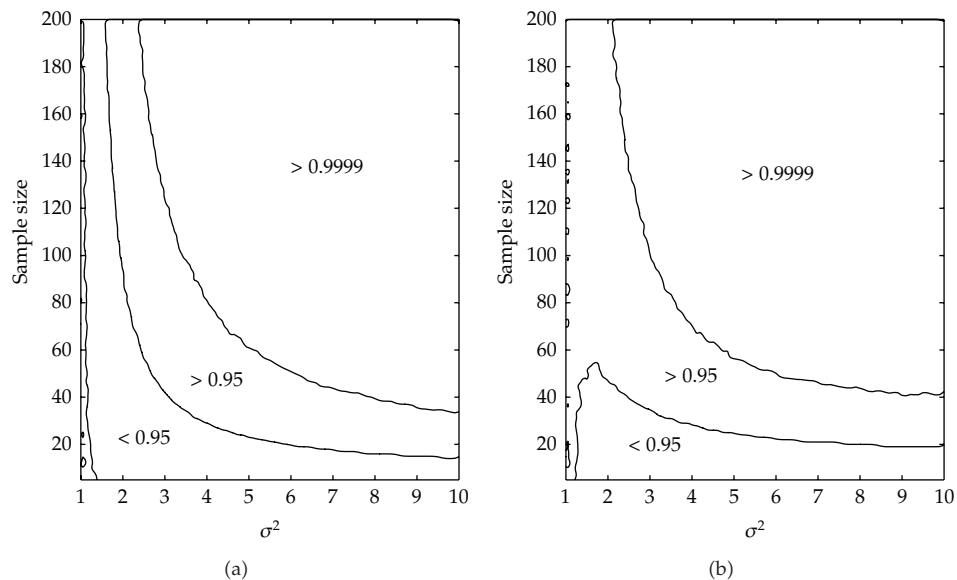


Figure 3: Contour plots of the L test (a) and R test (b) relative to the F test showing regions in which the power ratios are less than 95%, between 95% and 99.99%, and greater than 99.99%.

5. Robustness

Even if the R test has good power, the test is of little value unless it is robust in the sense that when the distributions from which we are sampling are not from the nominated population (here, the normal distribution), the P values are reasonably accurate. It is thus of interest to estimate the proportion of rejections when the null hypothesis is true and both the populations from which we sample are not normal. We have looked at variable kurtosis and skewness. First, variable kurtosis was considered via t -distributions with varying degrees of freedom ν , say. Second, variable skewness was considered through gamma distributions

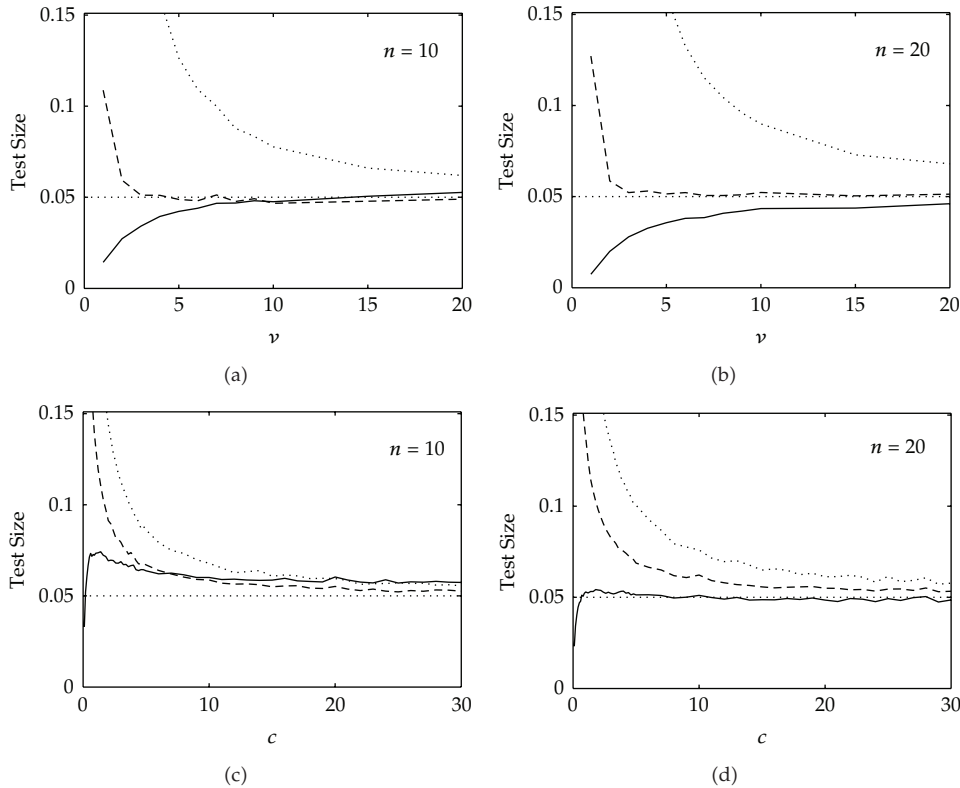


Figure 4: Test sizes for the F (dots), L (dashes), and R (solid line) tests for sample sizes of 10 and 20 for (a), (b) t -distributions with varying degrees of freedom, ν , from 1 to 20, and (c), (d) gamma distributions with scale parameter 1 and varying shape parameter, c , from 0.1 to 30.

with probability density functions $b^{c-1} \exp(-x)/\Gamma(c)$ for $x > 0$. Thus, the scale parameter is set at one and the scale parameter c was varied. As ν and c increase, the distributions become increasingly more normal, the distribution sampled will be close enough to normal that we could expect the proportion of rejections to be close to the nominal.

In Figure 4, we plot the proportion of rejections for the Levene, F , and R tests when sampling from t_ν distributions, for $\nu = 1, \dots, 20$, and from gamma distributions with scale parameter 1 and shape parameter c , for $c = 0.1, \dots, 30$. We show curves for each test with common sample sizes $n = 10$ and 20 . The critical values used for the R test are the $c(n, 0.05)$ from Section 3; the critical values used for the Levene test are those estimated by simulation to give exactly the nominal level when normality holds. Thus, this comparison favours the Levene test.

For these distributions, samples are increasingly nonnormal as ν and c decrease. It is apparent that the F test performs increasingly poorly as these parameters decrease.

When sampling from the t -distribution, the Levene test generally has exact level closer to the nominal level than the R test except for very nonnormal samples (small ν). However, the level of the R test is almost always reasonable, and while for very small ν , the level is not as close to the exact level as perhaps we may prefer, the same is the case for the Levene test.

When sampling from the gamma distribution when $n = 10$, the R test outperforms the Levene test at small values of c and is slightly inferior for larger values of c , although both

tests have exact test sizes acceptably close to the nominal test size. When $n = 20$, the R test is uniformly closer to the nominal test size than the Levene test. Although we only display results for $n = 10$ and $n = 20$, sample sizes of more than 20 yield very similar conclusions.

6. Conclusion

First, we reflect on testing for equality of variances when it is assumed that the populations sampled are normal. The F test is both the likelihood ratio test and a Wald test, and is the appropriate test to apply. When normality does not hold, the F test is no longer an asymptotically optimal test, and its well-known nonrobustness means that tests such as the Levene are more appropriate for small to moderate sample sizes. However, for sample sizes of about 25 or more, the R test is more powerful than the Levene, and with the small sample corrected critical values, it holds its nominal significance level well. For these sample sizes, it can be preferred to the Levene test.

Second, consider testing for equality of variances when both samples are drawn from the same population. If that population is nominated, then the R test may be applied after determining critical values or using P values calculated by Monte Carlo methods. When the sample variances are asymptotically equivalent to the maximum likelihood estimators of the population variances (as, e.g., is the case when sampling from normal populations but not Poisson populations), the R test is a nonparametric Wald test and hence will have good power in sufficiently large samples. If the population is not specified, the R test can be confidently applied when the sample sizes are large, using the asymptotic χ^2_1 null distribution to calculate P values or critical values.

References

- [1] B. F. J. Manly, *The Statistics of Natural Selection on Animal Populations*, Chapman and Hall, London, UK, 1985.
- [2] D. S. Moore and G. P. McCabe, *Introduction to the Practice of Statistics*, W. H. Freeman, New York, NY, USA, 5th edition, 2006.
- [3] W. J. Conover, M. E. Johnson, and M. M. Johnson, "A comparative study of tests of homogeneity of variances, with applications to the outer continental shelf bidding data," *Technometrics*, vol. 23, no. 4, pp. 351–361, 1981.
- [4] D. D. Boos and C. Brownie, "Bootstrap methods for testing homogeneity of variances," *Technometrics*, vol. 31, no. 1, pp. 69–82, 1989.
- [5] J. C. Rayner, "The asymptotically optimal tests," *Journal of the Royal Statistical Society Series D*, vol. 46, no. 3, pp. 337–346, 1997.
- [6] A. Stuart and J. K. Ord, *Kendall's Advanced Theory of Statistics*, vol. 1, Edward Arnold, London, UK, 6th edition, 1994.

