

# A Powerful and Robust Nonparametric Statistic for Joint Mean-Variance Quality Control

J.D. Opdyke

For statistical process control, a number of single charts that jointly monitor both process mean and variability recently have been developed. For quality control-related hypothesis testing, however, there has been little analogous development of joint mean-variance tests: only one two-sample statistic that is not computationally intensive has been designed specifically for the one-sided test of  $H_0: \mu_2 \leq \mu_1$  and  $\sigma_2 \leq \sigma_1$  vs.  $H_a: \mu_2 > \mu_1$  OR  $\sigma_2 > \sigma_1$  (see Opdyke (2005)). For these joint hypotheses, under many conditions tests of stochastic dominance (e.g. one-sided Kolmogorov-Smirnov) can severely violate the nominal level, and exceedance tests (e.g. Rosenbaum (1954)) and tests of distributional equality (e.g. most permutation tests) can have virtually no power. This paper further develops the maximum test proposed in Opdyke (2005) and demonstrates via thorough simulation that under typical quality control conditions a) it always maintains good level control; b) it has good power under symmetry and modest power under asymmetry; and c) it often has dramatically more power *and* much better level control than the only widely endorsed competitor. The statistic – OBMax2 – is not computationally intensive, and although initially designed for quality control testing in regulatory telecommunications, its range of application is as broad as the number of quality control settings requiring a one-sided test of the first two moments.

KEY WORDS: Maximum test; Location-Scale; Statistical process control; Six sigma; Telecommunications; CLEC.

## 1. INTRODUCTION

The statistical process control literature recently has seen the development of a number of single control charts that jointly monitor both process mean and variability (see Gan et al. (2004); Costa & Rahim (2004); Wu et al. (2005); Hawkins & Zamba (2005); and Reynolds & Stoumbos (2005)). Quality control-related hypothesis testing, however, has seen few analogous developments. Only two statistics known to this author (one computationally intensive – see Pesarin (2000), p.325) have been developed specifically to test the one-sided, joint mean-variance hypotheses of:

$$H_0: \mu_2 \leq \mu_1 \text{ and } \sigma_2 \leq \sigma_1 \text{ v. } H_a: \mu_2 > \mu_1 \text{ OR } \sigma_2 > \sigma_1 \quad (1)$$

Tests of stochastic dominance test fundamentally different hypotheses and, as shown below, can be entirely inappropriate for the joint hypotheses listed in (1). And exceedance tests, as

well as tests of distributional equality (e.g. most permutation tests), lack power to detect differences in both of the distributional characteristics that really matter from a quality control perspective. It is both the location and the spread (as conveniently, and often most appropriately, measured by the first two moments – the mean and the variance) of a quality metric’s distribution that typically are the quality control investigator’s primary or exclusive concern when comparing two data samples: higher moments often are irrelevant. For example, if two groups of customers are mandated to receive equal quality service, a difference in the kurtosis between the two groups’ time-to-service – *all else equal* – arguably has little or no effect on the perceived, and even actual, “quality” of service they receive. Of primary concern from a quality perspective would be a slower time-to-service *on average* for one group compared to the other, and/or a larger variability in the time-to-service for one group compared to the other. And this is exactly what is tested by the one-sided, two-sample statistic developed below: whether the mean and/OR variance of a quality metric of one population or process are larger than those of another.

## 2. PREVIOUS AND RELATED WORK

A number of statistics have been developed for the two-sided, location-scale hypotheses of

$$H_0: F(x) = G(x) \text{ v. } H_a: F(x) = G\left(\frac{x-\mu}{\sigma}\right) \quad (2)$$

with  $\sigma > 0$ , and  $\mu \neq 0$  and/or  $\sigma \neq 1$

(see O’Brien (1988); Podgor & Gastwirth (1994); Buning & Thadewald (2000); and Manly & Francis (2002)). But from a quality perspective we are more concerned with testing the one-sided hypotheses presented in (1) because the focus is on whether the quality of one population or process is *worse than* (better than) that of the other, not just different from that of the other. One statistic has received widespread attention as a test of (1) in the regulatory telecommunications arena. Seven years’ worth of expert testimony, as well as multiple Rulings, Opinions, and Orders handed down by various state and federal regulatory bodies, have supported use of the ‘modified’ *t* statistic (3) (Brownie et al. (1990)) to compare the quality of service provided to two groups of telecommunications customers – competing local exchange carrier (CLEC) customers and incumbent local exchange carrier (ILEC) customers. The point is to ensure that the quality received by CLEC customers is “at least equal” to that received by ILEC customers (see Telecommunications Act of 1996, Pub. LA. No. 104-104, 110 Stat. 56 (1996), at S251 (c) (2) (C)), thus ensuring that a formerly regulated industry can effectively

$$t_{\text{mod}} = \frac{(\bar{X}_2 - \bar{X}_1) - (\mu_2 - \mu_1)}{s_1 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \text{ with } df = n_1 - 1 \quad (3)$$

J.D. Opdyke is President, DataMineIt, 46 Tioga Way, Commerce Center #310, Marblehead, MA 01945 (E-mail: JDOpdyke@DataMineIt.com). The author expresses sincere appreciation to Geri S. Costanza, M.S., for numerous and valuable insightful discussions.

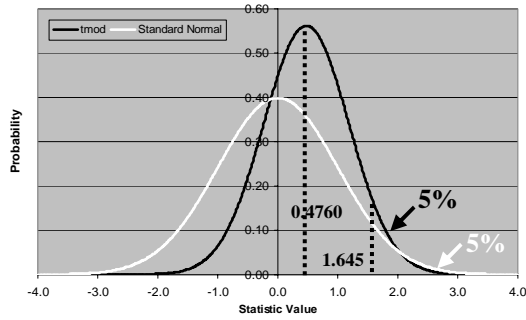


Figure 1a.  $t_{mod}$  v. Standard Normal,  $\mu_2 > \mu_1$ ,  $\sigma_2/\sigma_1 = 0.5$ ,  $n_1/n_2 = 100$

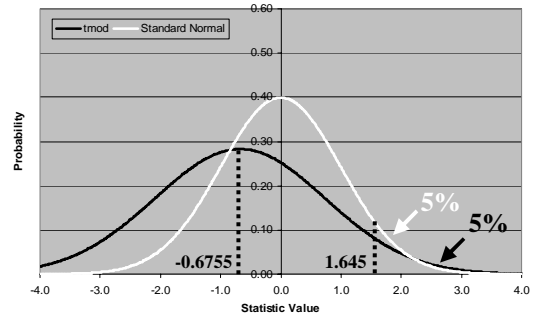


Figure 1b.  $t_{mod}$  v. Standard Normal,  $\mu_2 < \mu_1$ ,  $\sigma_2/\sigma_1 = 2$ ,  $n_1/n_2 = 100$

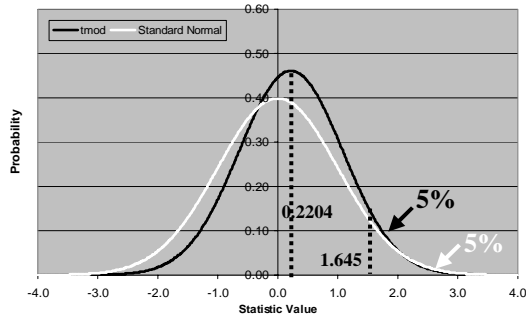


Figure 2a.  $t_{mod}$  v. Standard Normal,  $\mu_2 > \mu_1$ ,  $\sigma_2/\sigma_1 = 0.5$ ,  $n_1/n_2 = 1$

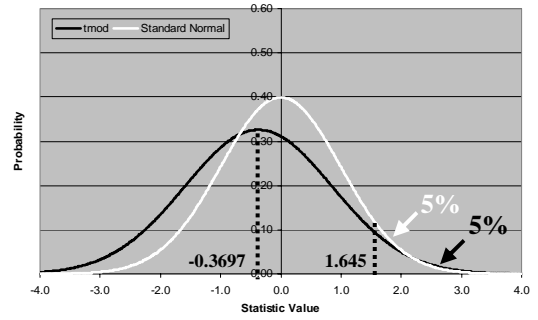


Figure 2b.  $t_{mod}$  v. Standard Normal,  $\mu_2 < \mu_1$ ,  $\sigma_2/\sigma_1 = 2$ ,  $n_1/n_2 = 1$

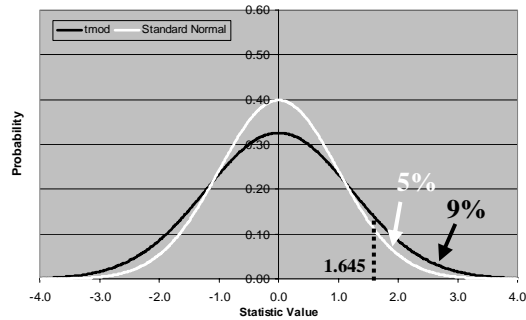


Figure 3a.  $t_{mod}$  v. Standard Normal,  $\mu_2 = \mu_1$ ,  $\sigma_2/\sigma_1 = 2$ ,  $n_1/n_2 = 1$

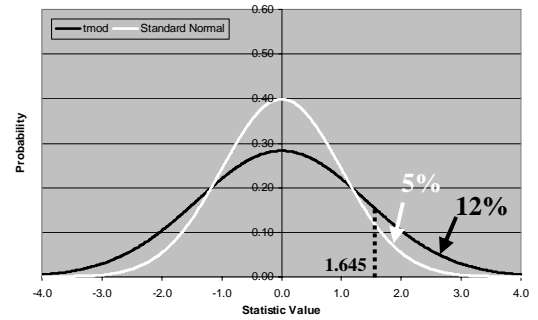


Figure 3b.  $t_{mod}$  v. Standard Normal,  $\mu_2 = \mu_1$ ,  $\sigma_2/\sigma_1 = 2$ ,  $n_1/n_2 = 100$

transition to a fully competitive economic market. The ‘modified’  $t$  statistic will be recognized as the widely used separate-variance  $t$  statistic (see Appendix) modified slightly in the denominator: the study group variance (#2) simply is swapped out and replaced with the control group variance (#1).

However, Opdyke (2004) demonstrated, via both analytic derivation and extensive simulation, that several crucial assumptions made about this statistic are false, making it inappropriate as a test of (1) in any setting. Its asymptotic distribution was shown to *not* be standard normal as previously surmised, but rather, normal with a variance that is greater than, less than, or equal to unity depending on the relative sizes of the two population variances, as shown below.

$$t_{mod} \sim N\left(0, \left(\frac{\sigma_1^2 + \sigma_2^2}{n_1 + n_2}\right) / \left(\frac{\sigma_1^2 + \sigma_1^2}{n_1 + n_2}\right)\right) \quad (4)$$

A consequence of this when using standard normal (or student’s  $t$ ) critical values, as advised in extensive expert testimony and Brownie et al. (1990), is that it allows a “trade-off” in average service for variability in service, which violates the null hypothesis of (1) with literally *zero* power to detect

these violations. This is shown very clearly in Figures 1a and 1b for (very) unbalanced sample sizes, and Figures 2a and 2b for equal sample sizes. Normalizing the ‘modified’  $t$  in an attempt to take care of this problem yields, not surprisingly, the separate-variance  $t$  statistic.

In addition to this fatal flaw, another problem with using this statistic as a test of (1) is that it has virtually no power to detect differences in variances, all else equal. For example, under equal means and a study group variance *twice as large* as that of the control group, the asymptotic power of the ‘modified’  $t$ , for  $\alpha = 0.05$ , is only 0.09 for equal sample sizes, and only 0.12 for very unbalanced ( $n_1 / n_2 = 100$ ) sample sizes, as shown in Figures 3a and 3b, respectively.

Although Brownie et al. (1990) originally proposed the ‘modified’  $t$  for use with a different pair of joint hypotheses

$$H_0: \mu_2 \leq \mu_1 \text{ and } \sigma_2 \leq \sigma_1 \text{ v. } H_a: \mu_2 > \mu_1 \text{ AND } \sigma_2 > \sigma_1 \quad (5)$$

for which the statistic usually (but not always) has more power than the pooled- and separate-variance  $t$  tests, it nonetheless remains essentially useless as a test of (1) based on the above findings, extensive expert testimony notwithstanding (see Opdyke (2004) for extensive citations).

As an alternative to the ‘modified’  $t$  statistic, Opdyke (2004) proposed the collective use of several easily implemented conditional statistical procedures. Four tests are proposed by combining O’Brien’s (1988) generalized  $t$  test (OBt) or his generalized rank sum test (OBr) with either of two straightforward tests of variances – Shoemaker’s (2003)  $F_1$  test and the ‘modified’ Levene test (Brown & Forsythe (1974)) (see Appendix for corresponding formulae). These easily calculated statistics are combined based on the relative size of the two sample means, as shown in Table 1.

Table 1. Conditional Statistical Procedures of Opdyke (2004)

Conditional statistical procedure	if $\bar{X}_2 > \bar{X}_1$ , use...	If $\bar{X}_2 \leq \bar{X}_1$ or OB fails to reject $H_0$ :, use...
OBtShoe	OBt	Shoemaker’s $F_1$
OBtLev	OBt	‘modified’ Levene
OBrShoe	OBr	Shoemaker’s $F_1$
OBrLev	OBr	‘modified’ Levene

For symmetric data, the choice of which of these four tests to use is based on two criteria – whether the data is at least as short-tailed as the normal distribution (platy- to mesokurtotic) vs. long-tailed (leptokurtotic), and whether sample sizes are balanced (or close) vs. at least moderately unbalanced, as shown in Table 2.

Table 2. Implementation of Table 1 Procedures Under Symmetry

Sample Sizes	Kurtosis of Distribution	
	platy- to mesokurtotic (OBt)	leptokurtotic (OBr)
Balanced (Shoemaker’s $F_1$ )	OBtShoe	OBrShoe
Unbalanced (‘modified’ Levene)	OBtLev	OBrLev

However, implementing Table 2 by deciding, for example, how unbalanced long-tailed samples must be before using OBrLev rather than OBrShoe requires additional, albeit straightforward, simulations not performed in Opdyke (2004). Subsequently, Opdyke (2005) bypassed this requirement by combining the four Table 1 statistics using a maximum-test approach.

“Maximum tests” – statistics whose scores ( $p$ -values) are the maximum (minimum) of two or more other statistics – have been devised and studied in a number of settings in the statistics literature with very favorable results. Neuhauser et al. (2004) compare a maximum test for the non-parametric two-sample location problem to multiple adaptive tests, finding the former to be most powerful under the widest range of data conditions. And Blair (2002) constructs a maximum test of location that is shown to be only slightly less powerful than each of its constituent tests under their respective “ideal” data conditions, but notably more powerful than each under

their respective “non-ideal” data conditions. These findings demonstrate the general purpose of maximum tests – to trade-off minor power losses under ideal data conditions for a more robust statistic with larger power gains across a wider range of possible (and usually unknown) data distributions.

To construct a maximum test for the joint mean-variance hypotheses of (1), it must be recognized that maximum tests are conditional statistical procedures, and the additional variance introduced by such conditioning will inflate the test’s size over that of its constituent statistics (and if left unadjusted, probably over the nominal level of the test as shown in Blair (2002)). But the constituent statistics in Table 1 are already conditional statistical procedures, so the  $p$ -value adjustment used to maintain validity must be large enough to take this “double conditioning” into account (this actually is “triple conditioning” since O’Brien’s tests themselves are conditional statistical procedures). The adjustment used in Opdyke (2005) is simply a multiplication of the  $p$ -values by constant factors ( $\beta$ ’s), the values of which were determined based on simulations. The  $p$ -value of the maximum test – OBMax – is defined below:

$$P_{OBMax} = \min \left( \begin{array}{l} P_{OBtShoe} \cdot \beta_{OBtShoe} \text{ ,} \\ P_{OBtLev} \cdot \beta_{OBtLev} \text{ ,} \\ P_{OBrShoe} \cdot \beta_{OBrShoe} \text{ ,} \\ P_{OBrLev} \cdot \beta_{OBrLev} \text{ ,} \\ P_{t_{sv}} \cdot \beta_{t_{sv}} \text{ ,} \\ 1.0 \end{array} \right) \quad (6)$$

where  $\beta_{OBtShoe} = \beta_{OBtLev} = \beta_{OBrShoe} = \beta_{OBrLev} = 2.8$ , and  $\beta_{t_{sv}} = 1.8$ , and  $p_{t_{sv}}$  is the  $p$ -value corresponding to the separate-variance  $t$  test with Satterthwaite’s (1946) degrees of freedom (see Appendix for corresponding formulae).

While analytic derivation of the asymptotic distribution of OBMax would be preferable to reliance on the simulation-based  $\beta$ ’s, Yang et al. (2005) show that such derivations for maximum tests are non-trivial, even under much stronger distributional assumptions than can be made with the conditional statistical procedures of Table 1. Babu and Padmanabhan (1996) describe the exact null distribution of their omnibus maximum test as “intractable” and rely on thorough simulation to demonstrate the validity and power of their statistic. Opdyke (2005) takes a similar approach to demonstrate the dramatically greater power of OBMax over the ‘modified’  $t$  under most alternate hypothesis configurations of (1). However, OBMax has two limitations: it can violate the nominal level when  $n_2 > n_1$ , as well as when, under asymmetry,  $\sigma_2 < \sigma_1$  and at least moderately large  $n_2 \approx n_1$  (the former condition was not a problem in the setting for which OBMax originally was developed – the regulatory telecommunications arena – since  $n_{ILEC} \geq n_{CLEC}$  virtually always). This paper eliminates these drawbacks with the development of a more robust statistic – OBMax2 – which maintains validity under asymmetry and any combination of sample sizes; it also retains most of the power of OBMax.

### 3. METHODOLOGY

#### 3.1 Development of OBMax2

If not stylized for specific asymmetric distributions, most two-sample statistics lose power under asymmetric data, and the constituent tests of OBMax are no exception to this general rule. However, under certain conditions under asymmetry, OBMax fails to maintain validity: if sample sizes are large and equal (or close) and the study group variance is much *smaller* than the control group variance, OBMax (under asymmetry) will often violate the nominal level of the test. This is due to O'Brien's rank sum test (OBr) behaving badly under these conditions – surprisingly, skewed-tail outliers invalidate the Table 1 statistics that use this test under these specific conditions. Although data transformations toward symmetry can alleviate this problem to some degree, there is no guarantee this will fix the problem altogether, if much at all. Instead, Opdyke (2005) proposes the use of another maximum test – OBMax3 – if symmetry cannot be assured. OBMax3 uses only three constituent tests, eliminating the two that use O'Brien's rank sum tests, as shown below:

$$P_{OBMax3} = \min \left( \begin{array}{l} P_{OBtLev} \cdot \beta_{OBtLev} \text{ ,} \\ P_{OBtShoe} \cdot \beta_{OBtShoe} \text{ ,} \\ P_{t_{sv}} \cdot \beta_{t_{sv}} \text{ ,} \\ 1.0 \end{array} \right) \quad (7)$$

where  $\beta_{OBtLev} = \beta_{OBtShoe} = 3.0$ , and  $\beta_{t_{sv}} = 1.6$

OBMax3 maintains validity under both symmetric and asymmetric data, with maximum power losses of well under 0.10 compared to OBMax under symmetry (see Opdyke (2005)). However, it unarguably would be preferable to have, rather than two tests, a single test robust to departures from symmetry that also retains most of the power of OBMax. And that is what OBMax2 accomplishes, as defined in (8) below:

$$P_{OBMax2} = P_{OBMax3} \text{ if and only if} \quad (8)$$

- a)  $s_2^2 \leq s_1^2$  and
- b)  $\bar{X}_2 \leq (\bar{X}_1 + 0.5s_1)$  and
- c) the null hypothesis of symmetry is rejected for either sample by the test of D'Agostino et al. (1990) at  $\alpha = 0.01$  (see Appendix)

$$P_{OBMax2} = P_{OBMax} \text{ otherwise}$$

This conditioning on a), b) and c) in (8) causes minor power losses in OBMax2 ("2" for two maximum tests) compared to OBMax under symmetry, but the worst level violations, even under asymmetry, are small – far smaller than those of the 'modified'  $t$  and separate-variance  $t$  statistics. Before discussing the simulation study, however, one other adjustment to OBMax2 is presented below.

In the regulatory telecommunications arena for which OBMax originally was developed, the size of the ILEC

customer sample (the "control group" sample, #1) almost always dwarfs that of the CLEC customer sample (the "study group" sample, #2), so the behavior of OBMax under  $n_2 > n_1$  was not a concern. The present development of OBMax2, however, seeks to generalize its use under the widest range of possible conditions, making it robust and powerful not only under both symmetry and asymmetry, but also under all possible combinations of sample sizes. Since it turns out that, under  $n_2 > n_1$ , increased variation of OBMax's (and OBMax3's) constituent statistics causes its violation of the nominal level of the test, an additional adjustment is required when  $(n_2 / n_1) > 1$  for OBMax2 to maintain validity. This is accomplished simply by increasing the size of the  $\beta$  adjustments as a function of the sample size ratio, as shown below:

$$\beta_X = \beta_X + \min \left[ 2.5, \max \left( 0, \log_{2.7} \left[ n_2 / n_1 \right] \right) \right] \quad (9)$$

The maximum function in (9) ensures that the  $\beta$ 's are increased only if  $(n_2 / n_1) > 1$ , and the minimum function ensures that the largest adjustment is +2.5, which was shown in simulations of up to  $(n_2 / n_1) = (3,000 / 30) = 100$  to be adequate. The empirical level and power of OBMax2, as defined in (8) together with (9), are presented in the simulation study results below (It is important to note that when implementing OBMax2, O'Brien's tests are referenced to the  $F$  distribution, rather than Blair's (1991) size-correcting critical values, even though doing so would normally violate the nominal level of the test under some conditions, because the  $p$ -value  $\beta$  adjustments used here explicitly take this size inflation into account, as described above).

#### 3.2 Simulation Study

Under a wide range of data distributions, sample size and mean-variance configurations, this study examines the empirical level and power of seven statistics: OBMax2, as in (8) and (9); OBMax, as in (6); OBMax3, as in (7); the 'modified'  $t$  statistic ( $t_{mod}$ ), as in (3); the separate-variance  $t$  statistic ( $t_{sv}$ ) with Satterthwaite's (1946) degrees of freedom (see Appendix) to provide a well-known basis for comparison; Rosenbaum's (1954) exceedance test (Ros), which counts the number of observations in one sample beyond the maximum of the other as a test of  $H_0: F(x) = G(x)$  against the general shift alternative; and the (one-sided) Kolmogorov-Smirnov statistic (K-S) (using Goodman's (1954) Chi-square approximation – see Siegel & Castellan (1988), p.148), a widely used test of stochastic dominance whose basic structure, which relies on the difference between the samples' cumulative distribution functions, underlies many such tests. Although not designed specifically for (1), these latter two statistics are included here because researchers often turn to these and similar tests, as well as tests of distributional equality (like permutation tests), when confronted with (1), and it is important to study their behavior under carefully controlled simulations (for example, the K-S statistic has been described as being "able to detect not only differences in average but differences in dispersion between the two samples as well." (see Matlack (1980), p. 359), which could be (mis)interpreted as an endorsement of this statistic for testing (1)).

The simulation study data was generated from the normal, uniform, double exponential, exponential, and lognormal distributions for five different pairs of sample sizes ( $n_2 = n_1 = 30, 90, \& 300; n_2 = 30 \& n_1 = 300; \text{ and } n_2 = 300 \& n_1 = 30$ ), seven different variance ratios ( $[\sigma_2 / \sigma_1] = 0.50, 0.75, 1.00, 1.25, 1.50, 1.75, 2.00$ ) and seven different location shifts ( $\mu_2 = \mu_1 - 2\sigma_1, \mu_1 - \sigma_1, \mu_1 - 0.5\sigma_1, \mu_1, \mu_1 + 0.5\sigma_1, \mu_1 + \sigma_1, \mu_1 + 2\sigma_1$ ), making 1,225 scenarios.  $N = 10,000$  simulations were run for each scenario.

The normal distribution was chosen as a universal basis for comparison; the uniform and double exponential distributions were chosen as examples of short- and long-tailed distributions, respectively, to examine the possible effects of kurtosis on the tests; and the exponential and lognormal distributions were chosen to examine the possible effects of skewness on the tests. The simulations do not include cross-distributional comparisons: the two data samples always are generated from the same distribution. While this assumption of same (or very close) distributions commonly is made by researchers when using two-sample tests, it arguably has stronger justification in the quality control setting since the comparison is of the *quality* of two otherwise similar or identical processes or populations, not of two processes or populations that are potentially completely dissimilar by any of numerous criteria. In the latter case, other statistical tests, or in-depth distributional examinations, are more appropriate.

Sample sizes were chosen to cover a range considered to be fairly small to moderately large, as well as balanced to fairly unbalanced, and two common nominal levels were used:  $\alpha = 0.05, 0.10$ , bringing the total number of scenarios to 2,450.

#### 4. RESULTS

The rejection rates under the null and alternate hypotheses of (1) – empirical level and power, respectively – for OBMax2 are shown in Figure 4 for  $\alpha = 0.05$  (complete study results are available from the author). General patterns in power can be observed. Not surprisingly, OBMax2 has more power for detecting mean increases, all else equal, than for detecting variance increases, all else equal. Also as expected, power increases as sample sizes increase, with slightly greater power under  $n_2 = 300 \& n_1 = 30$  compared to  $n_2 = 30 \& n_1 = 300$  for  $\sigma_2 > \sigma_1$ , but often vice versa for  $\mu_2 > \mu_1$ . Power is greatest under short-tailed (uniform) data, decreasing steadily for longer-tailed data as kurtosis increases (to normal and then double exponential data). Power is lowest under skewed data, and the greater the asymmetry, the lower the power (the lognormal samples are more skewed with these mean-variance configurations than the exponential samples).

Figure 5 graphs a histogram of the differences in power between OBMax2 and the ‘modified’  $t$ . OBMax2 completely dominates the ‘modified’  $t$  whenever the study group variance is larger and the study group mean is equal or smaller (most of the 79% of the alternate hypothesis cases where there is a difference in power). This is just a demonstration of the modified  $t$ ’s inability to detect differences in variances, all else equal, as seen above in Figures 3a and 3b. However, when the study group mean is slightly larger – i.e. under slight location shifts – the ‘modified’  $t$  does have a slight power advantage over OBMax2 (most of the 21% of the alternate hypothesis cases where there is a difference in power). This is the price

paid by OBMax2 for its ability to serve simultaneously as a test of differences in both means and/OR variances. However, this modest power advantage disappears as samples increase to moderate sizes, as seen in Figures 6-8.

When there is a difference in power between OBMax2 and OBMax3, the former almost always wins, with the largest power gains under asymmetry and  $\sigma_2 > \sigma_1$  (0.23, 0.18 for  $\alpha = 0.05, 0.10$ ). OBMax3 only has more power when  $n_2 > n_1$ , but then it often substantially violates the nominal level. Compared to OBMax, the largest power loss of OBMax2 approaches 0.13 under symmetry, which is the price paid for OBMax2’s very good level control, even under asymmetry, as seen in Figure 9 compared to the ‘modified’  $t$ . Both the ‘modified’  $t$  and separate-variance  $t$  statistics substantially violate the nominal level of the test under skewed data (a well known result for the latter statistic). The worst level violations of OBMax2, on the other hand, are what most researchers would consider reasonable, if not quite good: between 0.06 and 0.07 when  $\alpha = 0.05$ , and just over 0.11 when  $\alpha = 0.10$ .

Regarding the other tests, although the Rosenbaum exceedance statistic always maintains validity, it often has dramatically less power than OBMax2, especially if the study group mean is smaller than the control group mean, when it often has absolutely no power to detect a larger study group variance (which actually is consistent with its design). This latter finding also is true of the K-S statistic which, although sometimes more powerful than OBMax2 under small location shifts, often severely violates the nominal level when means are identical but the study group variance is *smaller*. This is because the smaller variance causes the distance between the empirical cumulative distribution functions of the two samples, which determines the significance of K-S, to be sizeable on either side of the common mean (see Figures 10 and 11), thus causing high rejection rates for one-sided tests simultaneously in *both* directions. In this study, for example, when  $\mu_1 = \mu_2, \sigma_2 / \sigma_1 = 0.50$ , and  $n_1 = n_2 = 300$ , the K-S rejection rate, when  $\alpha = 0.05$ , is over 0.95 for the uniform, exponential, and lognormal distributions, and very large for the double exponential (0.35) and normal (0.57) distributions as well. This makes it and similarly structured tests of stochastic dominance that rely on the difference between cumulative distribution functions unusable for the joint hypotheses of (1).

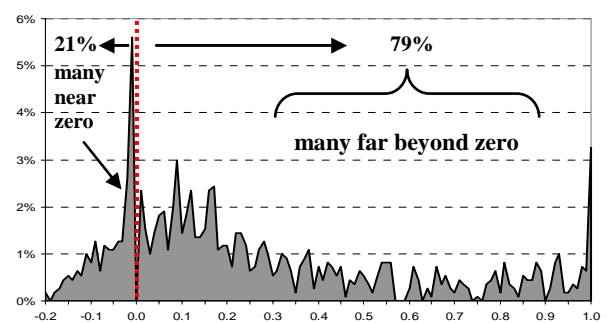


Figure 5. OBMax2 power minus  $t_{mod}$  power (1,106/1,850  $\neq 0$ )

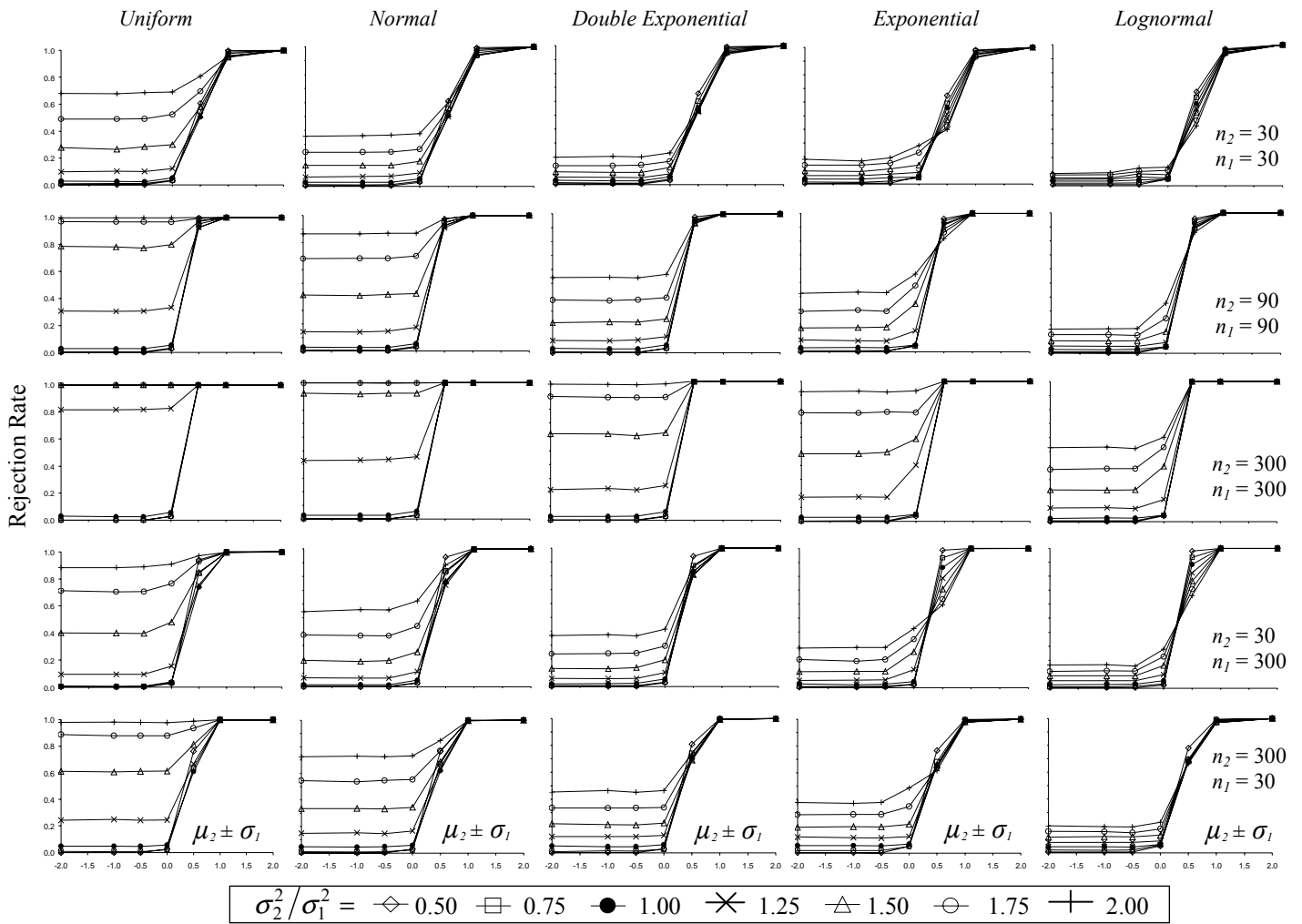


Figure 4. OBMax2 Rejection Rate (Empirical Level & Power) by Distribution by Sample Size by Mean Shift by Variance Ratio ( $\alpha = 0.05$ )

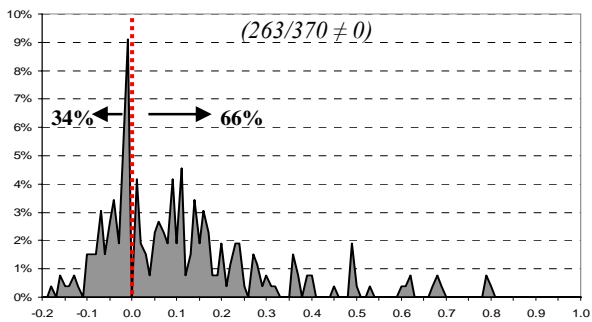


Figure 6. OBMax2 power minus  $t_{mod}$  power,  $n_1 = n_2 = 30$

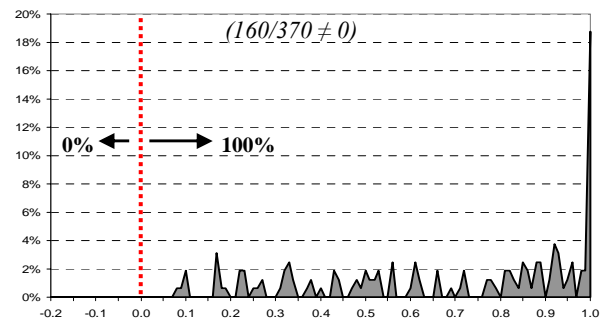


Figure 8. OBMax2 power minus  $t_{mod}$  power,  $n_1 = n_2 = 300$

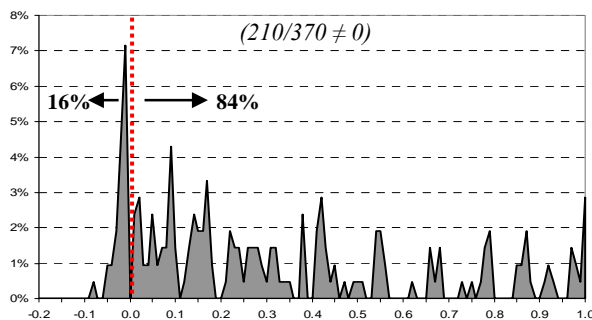


Figure 7. OBMax2 power minus  $t_{mod}$  power,  $n_1 = n_2 = 90$

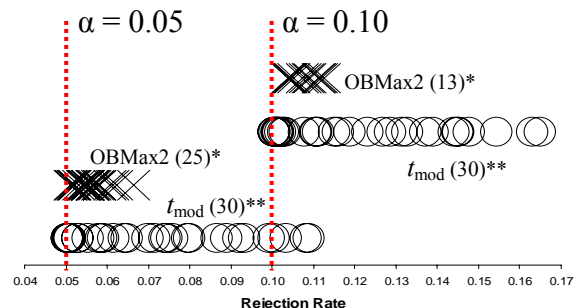


Figure 9. Level Violations:  $t_{mod}$  vs. OBMax2 (600 null hypotheses)

NOTE: \*15/38 from  $n_2 = 300$  &  $n_1 = 30$ , \*\*47/60 from skewed data

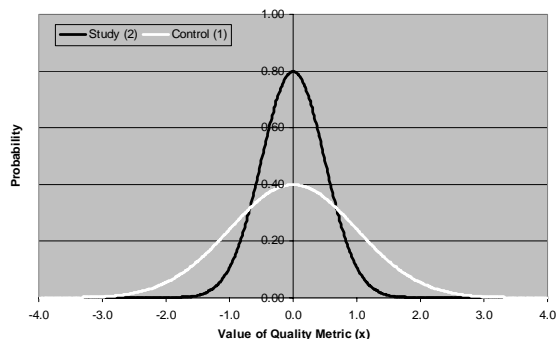


Figure 10. PDFs: K-S Violates Nominal Level when  $\mu_1=\mu_2, \sigma_2<\sigma_1$

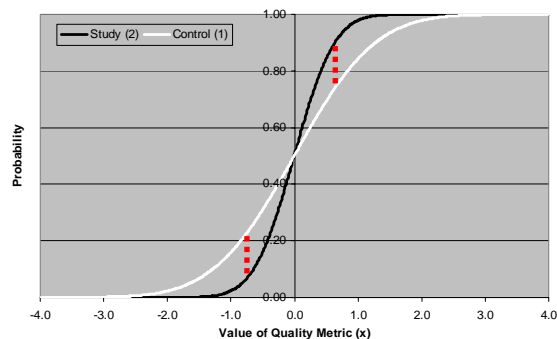


Figure 11. CDFs: K-S Violates Nominal Level when  $\mu_1=\mu_2, \sigma_2<\sigma_1$

## 5. DISCUSSION AND CONCLUSIONS

In quality control settings requiring a one-sided test of the joint mean-variance hypotheses of (1), many researchers and investigators to date have been relying on inappropriate statistics. Although supported by extensive expert testimony, the ‘modified’  $t$  was shown here and in Opdyke (2004, 2005) to be completely powerless under many, if not the majority, of alternate hypothesis scenarios under (1). Tests of stochastic dominance, too, appear to be inappropriate. Arguably the most commonly used of these – the Kolmogorov-Smirnov statistic – was shown in this simulation study not only to have no power to detect a larger study group variance when the study group mean is much smaller (which is consistent with its design), but also to severely violate the nominal level of the test when means are equal and the study group variance is *smaller*. The first of these two drawbacks also was shown to be true for a common exceedance test. The one type of test not examined here that researchers and quality investigators sometimes turn to when confronted with (1) are permutation tests. But it is well known that permutation tests of  $H_0: F(x) = G(x)$  vs.  $H_a: F(x) \neq G(x)$  have very low power to detect differences in variances, all else equal. Permutation tests of scale alone, however, have been developed (see Good (2000) and Pesarin (2000)), but only one has been developed specifically as a test of the joint mean-variance hypotheses of (1). Pesarin (2000) combines two permutation tests – one of scale and one of location – using any of several  $p$ -value combining functions (pp.147-148, 325). However, his test has one, and potentially two drawbacks relative to OBMax2: first, it is computationally intensive which, all else equal, is a limitation, and even runtime prohibitive for larger samples. Secondly, the combining function it relies upon must be chosen with care since many in common usage (e.g. Fisher and Liptak) would decrease the power of the test under many situations due to their “trade-off” nature: a small and significant  $p$ -value from one of the constituent tests can be “undone” by a large  $p$ -value from the other test when the two are combined, in which case the overall test will lose power under (1). While the Tippet function, which is itself a maximum-test approach, appears to avoid this problem, OBMax2 never suffers from a “trade-off” problem and thus, may be “safer” (i.e. more powerful, all else equal) if the Tippet function cannot be used with Pesarin’s test for some reason.

Thus, only OBMax2 remains as a generally powerful, valid, and easily implemented test of (1) – whether the mean and/OR variance of one process or population are larger than (worse than) those of another. From a quality perspective, this is the specific hypothesis that matters. The simulation study presented above demonstrates OBMax2’s notable power under symmetry, and its much more modest power under asymmetry, which unfortunately is common among two-sample statistics and is the price paid for OBMax2’s good level control under all conditions. Still, improving its power under skewed data is a worthy objective, and remains the subject of continuing research, as is the derivation of its asymptotic distribution.

To end on a cautionary note, it should be (re)emphasized that OBMax2 was developed specifically for hypothesis testing, for which it was shown to have good level control for (1) at commonly used nominal levels. In other words, its  $p$ -values appear to be uniformly distributed (at or under the appropriate “height”) in the range that matters for most hypothesis tests. It would *not* be advisable, however, to use OBMax2’s  $p$ -values for other purposes requiring a presumption of uniformity across its entire domain of zero to one. As one might expect, inflating the  $p$ -values of OBMax2’s constituent statistics to maintain its overall validity has the consequence of “bunching up” its otherwise uniform  $p$ -value distribution at 1.0, as seen in Figure 12. One positive note

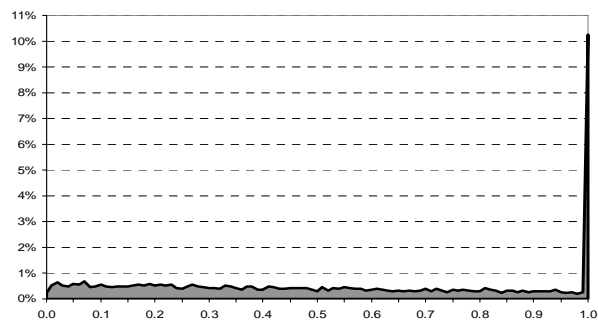


Figure 12. OBMax2  $p$ -values:  
Normal Data,  $\mu_1=\mu_2, \sigma_1=\sigma_2, n_1=n_2=30, N$  simulations=10,000

from this finding, however, is that an examination of all of the 600 null distributions of OBMax2’s  $p$ -values that were generated in this study reveals good level control for even larger nominal levels, such as  $\alpha = 0.15$  and 0.20. And power given robust level control is all that matters for a hypothesis test, and this is the criterion by which OBMax2 far surpasses any of its competitors.

## APPENDIX: Statistical Formulae

**OBt and OBr:** O'Brien's OBt test involves running the following ordinary least squares regression on pooled data including both samples:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i, \quad (10)$$

where  $y$  is a dummy variable indicating inclusion in the study group sample, and  $x$  is the performance metric variable. If the parameter on the quadratic term ( $\beta_2$ ) is (positively) statistically significant at the 0.25 level, use the critical value of the overall equation (an  $F$  test of  $\beta_1 = \beta_2 = 0$ ) to reject or fail to reject the null hypothesis; if it is not, use the critical value of the overall equation of the following ordinary least squares regression instead:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (11)$$

(Of course, this latter test of  $\beta_1 = 0$  is equivalent to a two-sample  $t$  test.) O'Brien's OBr test is identical to the OBt test except that the pooled-sample ranks of  $x$  are used in the regressions instead of the  $x$  data values themselves.

**'modified' Levene test:** The 'modified' Levene test requires a simple data transformation: take the absolute value of each data point's deviation from its respective sample median (as per Brown & Forsythe (1974)), and then calculate the usual one-way ANOVA statistic using these transformed values (as per Levene (1960)). The resulting statistic (12) is referenced to the  $F$  distribution as usual.

Let  $z_{ij} = |x_{ij} - \tilde{x}_i|$  where  $\tilde{x}_i$  is sample  $i$ 's median

$$W_o = \frac{\sum_i n_i (\bar{z}_i - \bar{z}_{..})^2 / (g-1)}{\sum_i \sum_j (z_{ij} - \bar{z}_i)^2 / \sum_i (n_i - 1)} \sim F_{(g-1), \sum_i (n_i - 1)} \quad (12)$$

where  $\bar{z}_i = \sum z_{ij} / n_i$  and  $\bar{z}_{..} = \sum \sum z_{ij} / n_i$

However, because this test is designed as a two-tailed test, and the hypotheses being tested in (1) are one-tailed, the  $p$ -value resulting from this test, when used conditionally with O'Brien's tests as in Table 1, must be subtracted from 1.0 if the study group sample variance is less than the control group sample variance.

**Shoemaker's  $F_1$  test:** Shoemaker's  $F_1$  test is simply the usual ratio of sample variances referenced to the  $F$  distribution, but using different degrees of freedom:

$$s_2^2 / s_1^2 \sim F_{df_2, df_1} \quad \text{where} \quad df_i = 2n_i / \left( \frac{\hat{\mu}_4}{\hat{\sigma}^4} - \frac{n_i - 1}{n_i - 3} \right), \quad (13)$$

$i = 1, 2$  corresponds to the two samples, and  $\mu_4$  and  $\sigma^4$  are estimated from the two samples when pooled:

$$\hat{\mu}_4 = \sum_{i=1}^2 \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^4 / (n_1 + n_2) \quad (14)$$

$$\hat{\sigma}^4 = \left[ ((n_1 - 1)s_1^2 + (n_2 - 1)s_2^2) / (n_1 + n_2) \right]^2 \quad (15)$$

Shoemaker (2003) notes that the biased estimate for  $\sigma^4$  is used for improved accuracy.

**separate-variance  $t$  test:** The separate-variance  $t$  test, also known as the Welch or Behrens-Fisher  $t$  test, is below:

$$t_{sv} = \frac{(\bar{X}_2 - \bar{X}_1) - (\mu_2 - \mu_1)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (16)$$

$$\text{where} \quad s_1^2 = \frac{\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2}{(n_1 - 1)}, \quad s_2^2 = \frac{\sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2}{(n_2 - 1)},$$

$$\bar{X}_1 = \frac{\sum_{i=1}^{n_1} X_i}{n_1}, \quad \text{and} \quad \bar{X}_2 = \frac{\sum_{i=1}^{n_2} X_i}{n_2}$$

Satterthwaite's (1946) degrees of freedom for  $t_{sv}$  is:

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{\left( \frac{s_1^2}{n_1} \right)^2}{(n_1 - 1)} + \frac{\left( \frac{s_2^2}{n_2} \right)^2}{(n_2 - 1)}} \quad (17)$$

If  $df$  is not an integer, it should be rounded down to the next smallest integer (see Zar (1999), p.129)

**test of D'Agostino et al. (1990):** The test of D'Agostino et al. (1990) is calculated as follows:

$$g_1 = \frac{k_3}{s^3} = \frac{n \sum_{i=1}^n (X_i - \bar{X})^3}{(n-1)(n-2) \sqrt{(s^2)^3}}, \quad \sqrt{b_1} = \frac{(n-2)g_1}{\sqrt{n(n-1)}} \quad (18)$$

$$A = \sqrt{b_1} \sqrt{\frac{(n+1)(n+3)}{6(n-2)}}, \quad B = \frac{3(n^2 + 27n - 70)(n+1)(n+3)}{(n-2)(n+5)(n+7)(n+9)}$$

$$C = \sqrt{2(B-1)} - 1, \quad D = \sqrt{C}, \quad E = \frac{1}{\sqrt{\ln D}}, \quad F = \frac{A}{\sqrt{C-1}}$$

$$Z_{g_1} = E \ln \left( F + \sqrt{F^2 + 1} \right) \sim \phi(0, 1)$$

For one-tailed testing of skewness to the left, check

$\Pr(Z \leq Z_{g_1})$ ; for skewness to the right, check  $\Pr(Z \geq Z_{g_1})$ .

See Zar (1999), pp.115-116, for further details.

## REFERENCES

- Babu, G. & Padmanabhan, A.R. (1996), "A Robust Test for Omnibus Alternatives," *Research Developments in Probability and Statistics*, E. Brunner and M. Denker (eds.), VSP, Utrecht, 319-327.
- Blair, R.C. (2002), "Combining Two Nonparametric Tests of Location," *Journal of Modern Applied Statistical Methods*, 1(1), 13-18.
- Blair, R.C. (1991), "New Critical Values for the Generalized  $t$  and Generalized Rank-sum Procedures," *Communications in Statistics*, 20, 981-994.
- Brown, M. & Forsythe, A. (1974), "Robust Tests for the Equality of Variances," *Journal of the American Statistical Association*, 69, 364-367.
- Brownie, C., Boos, D. D. & Hughes-Oliver, J. (1990), "Modifying the  $t$  and ANOVA  $F$  Tests When Treatment is Expected to Increase Variability Relative to Controls," *Biometrics*, 46, 259-266.
- Buning, H., & Thadewald, T. (2000), "An Adaptive Two-sample Location-scale Test of Lepage Type for Symmetric Distributions," *Journal of Statistical Computation and Simulation*, 65(4), 287-310.
- Costa, A.F.B., & Rahim, M.A. (2004), "Monitoring Process Mean and Variability with One Non-central Chi-square Chart," *Journal of Applied Statistics*, 31(10), 1171-1183.
- D'Agostino, R.B., A. Belanger, & R.B. D'Agostino, Jr. (1990), "A Suggestion for Using Powerful and Informative Tests of Normality," *The American Statistician*, 44, 316-321.
- Gan, F.F., Ting, K.W., & Chang, T.C. (2004), "Interval Charting Schemes for Joint Monitoring of Process Mean and Variance," *Quality and Reliability Engineering International*, 20(4), 291-303.
- Good, P., (2000), *Permutation Tests*, 2<sup>nd</sup> ed. New York: Springer-Verlag New York, Inc.
- Goodman, L.A. (1954), "Kolmogorov-Smirnov Tests for Psychological Research," *Psychological Bulletin*, 51, 160-168.
- Hawkins, D. M. & Zamba, K.D. (2005), "Statistical Process Control for Shifts in Mean or Variance Using a Change-point Formulation," *Technometrics*, 47(2), 164-173.
- Levene, H. (1960), "Robust Tests for Equality of Variances," in *Contribution to Probability and Statistics: Essays in Honor of Harold Hotelling*, I. Olkin et al., eds., Stanford University Press, 278-292.
- Manly, B., & Francis, C. (2002), "Testing for Mean and Variance Differences with Samples from Distributions that may be Non-normal with Unequal Variances," *Journal of Statistical Computation and Simulation*, 72(8), 633-646.
- Matlack, W.F., (1980), *Statistics for Public Policy and Management*, Belmont, CA: Duxbury Press.
- Neuhäuser, M. Büning, H., & Hothorn, L. A. (2004), "Maximum Test Versus Adaptive Tests for the Two-sample Location Problem," *Journal of Applied Statistics*, 31(2), 215-227.
- O'Brien, P. (1988), "Comparing Two Samples: Extensions of the  $t$ , Rank-sum, and Log-rank Tests," *Journal of the American Statistical Association*, 83, 52-61.
- Opdyke, J.D. (2004), "Misuse of the 'modified'  $t$  Statistic in Regulatory Telecommunications," *Telecommunications Policy*, 28, 821-866.
- Opdyke, J.D. (2005), "A Single, Powerful, Nonparametric Statistic for Continuous-data Telecommunications 'Parity Testing'," *Journal of Modern Applied Statistical Methods*, 4(4), forthcoming.
- Pesarin, F. (2001), *Multivariate Permutation Tests with Applications in Biostatistics*, John Wiley & Sons, Ltd., New York.
- Podgor, M.J., & Gastwirth, J.L. (1994), "On Non-parametric and Generalized Tests for the Two-sample Problem with Location and Scale Change Alternatives," *Statistics in Medicine*, 13(5-7), 747-758.
- Reynolds, M.R., & Stoumbos, Z.G. (2005), "Should Exponentially Weighted Moving Average and Cumulative Sum Charts be Used with Shewhart Limits?" *Technometrics*, 47(4), 409-424.
- Rosenbaum, S. (1954), "Tables for a Nonparametric Test of Location," *Annals of Mathematical Statistics*, 25, 146-150.
- Satterthwaite, F. W. (1946), "An Approximate Distribution of Estimates of Variance Components," *Biometrics Bulletin*, 2, 110-114.
- Siegel, S. & Castellan, N. John, (1988), *Nonparametric Statistics for the Behavioral Sciences*, 2nd ed., New York: McGraw-Hill.
- Shoemaker, L. H. (2003), "Fixing the  $F$  test for Equal Variances," *The American Statistician*, 57, 105-114.
- Telecommunications Act of 1996, Pub. LA. No. 104-104, 110 Stat. 56 (1996). [available at <http://www.fcc.gov/telecom.html>]
- Wu, Zhang, Tian, Y., & Zhang, Sheng (2005), "Adjusted-loss-function Charts with Variable Sample Sizes and Sampling Intervals," *Journal of Applied Statistics*, 32(3), 221-242.
- Yang, Song, Li Hsu, & Lueping Zhao (2005), "Combining Asymptotically Normal Tests: Case Studies in Comparison of Two Groups," *Journal of Statistical Planning and Inference*, 133(1), 139-158.
- Zar, J.H., (1999), *Biostatistical Analysis*, 4th ed., Upper Saddle River, NJ: Prentice-Hall.