

# Nonrobustness in classical tests on means and variances: A large-scale sampling study

JAMES V. BRADLEY

*New Mexico State University, Las Cruces, New Mexico 88003*

The robustness of the classical tests on means ( $Z$ ,  $t$ , and  $F$ ) and variances (chi square and  $F$ ) was investigated by obtaining 30,000 (or, sometimes, 10,000 or 150,000) values of the test statistic under assumption-violating conditions and comparing the actual proportion of Type I errors with the proportion expected when all assumptions are met. The sampling and testing conditions investigated were: population shape (L-shape or bell-shape), relative population variance (1 or 4), sample size (8, 16, or 24), nominal significance level (.05, .01, or .001), and location of rejection region (left-tailed, right-tailed, or two-tailed). All tests on means were nonrobust under most of the investigated conditions, and one was nonrobust under all of them. Both tests on variances were extremely nonrobust under virtually all conditions. The worst nonrobustness of the one- and two-independent-sample  $Z$  and  $t$  tests rivaled or exceeded that of the notoriously nonrobust chi-square and  $F$  tests on variances.

Unqualified, or insufficiently qualified, claims of robustness for the classical  $Z$ ,  $t$ , and  $F$  tests on means have become commonplace (Bradley, 1978), whereas the nonrobustness of the classical chi-square and  $F$  tests on variances is generally acknowledged (Scheffé, 1959). This article explores the robustness of all these tests under realistic conditions in terms of experimentation in the behavioral sciences. When the normality assumption is violated, this is accomplished by sampling from a population having an L-shape that is fairly common in certain areas of psychological and sociological research (Bradley, 1977); when the assumption of equal variances is violated, one population's standard deviation is only twice that of another. Samples are of reasonable absolute and relative size ( $8 \leq N \leq 24$ ), and standard  $\alpha$  levels and rejection regions are used. Yet under these perfectly realistic conditions, it will be seen that the robustness of tests on means depends heavily upon circumstances and that the  $Z$  and  $t$  tests are sometimes nonrobust to a degree that rivals or exceeds that of the admittedly nonrobust tests on variances.

## METHOD

### The Sampled Populations

Four different populations were constructed and labeled  $X$ ,  $Y$ ,  $x$ , and  $y$ . Population  $X$  was L-shaped and was constructed to resemble very closely a population of actual data (shown in Bradley, 1977, at the bottom of Figure B) obtained in a routine experiment (Bradley, 1980). Population  $Y$  was bell-shaped, exactly symmetric, and as nearly normal as its discreteness and limited range would permit. Populations  $X$  and  $Y$  had exactly the same mean and standard deviation and were similar in other ways. They are shown, and described in greater detail, in Bradley (1980).

By taking every unit in the  $X$  and  $Y$  populations and moving it halfway to the population mean, Populations  $x$  and  $y$  were

obtained. They are therefore essentially the same as the  $X$  and  $Y$  populations, respectively, except that they have a common standard deviation that is exactly half as large as that of the  $X$  and  $Y$  populations.

### The Tests Investigated

The robustness of eight different tests was investigated, always for the case in which the null hypothesis is true. The tests (and, in parentheses, the symbol used in this article to denote them) were: the one-sample  $Z$  test ( $Z_1$ ) and the one-sample  $t$  test ( $t_1$ ) of the hypothesis that the mean of the sampled population has a specified numerical value; the two-independent-sample  $Z$  test ( $Z_2$ ), the two-independent-sample  $t$  test ( $t_2$ ), the two-correlated-sample  $t$  test ( $t_c$ ), and the multi-independent-sample  $F$  test ( $F$ ) of the hypothesis that all sampled populations have the same unspecified mean; the one-sample chi-square test ( $\chi^2_v$ ) of the hypothesis that the variance of the sampled population has a specified numerical value; and the two-independent-sample  $F$  test ( $F_v$ ) of the hypothesis that the two sampled populations have the same unspecified variance. All of these tests assume (in addition to the usual assumptions about random sampling) that every sampled population is normally distributed, and the  $t_2$  and  $F$  tests also assume that all sampled populations have the same variance.

### Sampling

Sample size was always 8 for the one-sample tests and for each of the three or four samples contributing to  $F$ . For the two-sample tests, sample sizes were 8 and 16, 8 and 24, or 16 and 16. The exact combinations of sampling and testing conditions under which the tests were investigated are given in the tables presented in the Results section.

Sampling and computations were performed by an IBM 7090 computer. For each test statistic under each different sampling situation (i.e., for each row of the tables), the computer obtained a sampling distribution of that test statistic under those sampling conditions. The sampling distributions for  $Z_1$  and  $t_1$  each contained 150,000 values of the test statistic, and that for  $\chi^2_v$  contained 10,000 values. For all other test statistics, each sampling distribution was based upon 30,000 values of the test statistic.

Roughly speaking, independence of sampling existed within but not between sampling distributions. That is, nonoverlapping

series of nonrecycling pseudorandom numbers were used to draw the different observations contributing to the sampling distribution of a given test statistic in a given sampling situation. However, the same or overlapping sets of pseudorandom numbers were used to obtain the sampling distributions of the same test statistic under different sampling conditions or of different test statistics under either the same or different sampling conditions.

Correlated samples were obtained for the  $t_c$  test by using the same set of pseudorandom numbers in the same way (Bradley, 1980) to draw the sample observations from each of the two sampled populations.

As a check on the "randomness" of sampling, the accuracy of programming, and the general validity of the study, "control" sampling distributions were obtained for each test statistic (except  $t_c$ ) under each set of sample sizes for the case in which all of the samples come from the bell-shaped quasi-normal Population Y. Chi-square tests of fit to the appropriate normal-theory distribution of the test statistic showed that all was well.

RESULTS

For each test statistic under each different sampling situation, the computer obtained the proportion  $\rho$  of values of the test statistic that fell into a normal-theory rejection region corresponding to a nominal significance level of  $\alpha$ . Thus  $\rho$  is the empirical probability (and

therefore an estimate of the true probability) of a Type I error corresponding to an alleged probability  $\alpha$  (that would have been the true probability if all assumptions had been met). This was done for left-tailed, right-tailed, and two-tailed rejection regions corresponding to  $\alpha$  values of .05, .01, and .001.

Tables 1, 2, and 3 give the  $\rho$  values for the various tests under the various sets of circumstances that were investigated. The relatively robust  $\rho$  values lying in the interval from  $2\alpha/3$  to  $3\alpha/2$  are italicized; those lying outside of that interval but within the interval from  $\alpha/3$  to  $3\alpha$  are printed in light ordinary type; those lying outside of that interval but within the interval from  $\alpha/5$  to  $5\alpha$  are printed in medium type; and those lying outside of the latter interval, so that the larger of the two ratios  $\rho/\alpha$  and  $\alpha/\rho$  exceeds 5, are printed in heavy (boldface) type. This coding scheme refers to the exact values of  $\rho$ , so occasionally the tabled values, rounded to four decimal places, will appear to be misclassified. When  $\rho$  is exactly 0, it is written 0; when it becomes 0 only when rounded to the number of decimal places used, it is written .0000.

The most impressive case of robustness was that of

Table 1  
Robustness of Z (and Other) Tests: Empirical Probability  $\rho$  of a Type I Error Corresponding to Nominal Probability  $\alpha$  for Left-Tail (L), Two-Tail (T), and Right-Tail (R) Rejection Regions

Statistic	Pop. 1	Pop. 2	$\sigma_1/\sigma_2$	N <sub>1</sub>	N <sub>2</sub>	$\alpha = .05$			$\alpha = .01$			$\alpha = .001$		
						L	T	R	L	T	R	L	T	R
Z <sub>1</sub>	X			8	24	<b>0</b>	<i>.0451</i>	<i>.0683</i>	<b>0</b>	.0207	.0295	<b>0</b>	<b>.0071</b>	<b>.0099</b>
						.0252	<i>.0486</i>	<i>.0646</i>	<b>.0030</b>	.0160	.0215	<b>.0001</b>	<b>.0043</b>	<b>.0059</b>
						<i>.0344</i>	<i>.0513</i>	<i>.0613</i>	.0058	<i>.0148</i>	.0189	.0005	<b>.0038</b>	<b>.0048</b>
Z <sub>2</sub>	X	X	1	16	16	<i>.0492</i>	<i>.0506</i>	<i>.0477</i>	<i>.0107</i>	<i>.0124</i>	<i>.0119</i>	<i>.0019</i>	<i>.0018</i>	<i>.0015</i>
				16	8	<i>.0613</i>	<i>.0513</i>	<i>.0344</i>	.0189	<i>.0148</i>	.0058	<b>.0048</b>	<b>.0038</b>	.0005
				24	8	<i>.0646</i>	<i>.0486</i>	.0252	.0215	.0160	<b>.0030</b>	<b>.0059</b>	<b>.0043</b>	<b>.0001</b>
	X	Y	1	8	24	.0182	<i>.0441</i>	<i>.0652</i>	<b>.0004</b>	.0162	.0232	<b>0</b>	<i>.0045</i>	<b>.0065</b>
				8	16	<i>.0253</i>	<i>.0453</i>	<i>.0619</i>	<b>.0011</b>	<i>.0141</i>	.0221	<b>0</b>	<b>.0040</b>	<b>.0060</b>
				16	16	<i>.0384</i>	<i>.0481</i>	<i>.0560</i>	.0051	<i>.0108</i>	.0162	.0004	.0018	.0024
Z <sub>2</sub>	X	x	2	16	16	<i>.0452</i>	<i>.0487</i>	<i>.0543</i>	<i>.0072</i>	<i>.0100</i>	<i>.0127</i>	.0004	<i>.0012</i>	.0020
				16	8	<i>.0468</i>	<i>.0506</i>	<i>.0528</i>	<i>.0084</i>	<i>.0101</i>	<i>.0124</i>	.0008	<i>.0010</i>	<i>.0015</i>
				24	8	<b>.0053</b>	<i>.0440</i>	<i>.0690</i>	<b>0</b>	.0178	.0256	<b>0</b>	<b>.0058</b>	<b>.0081</b>
	X	y	2	8	16	<b>.0106</b>	<i>.0455</i>	<i>.0678</i>	<b>.0002</b>	.0168	.0252	<b>0</b>	<b>.0056</b>	<b>.0076</b>
				16	16	.0312	<i>.0478</i>	<i>.0608</i>	<b>.0031</b>	<i>.0142</i>	.0196	<b>.0002</b>	.0028	<b>.0047</b>
				16	8	<i>.0434</i>	<i>.0526</i>	<i>.0540</i>	<b>.0083</b>	<i>.0142</i>	.0160	<i>.0012</i>	<i>.0027</i>	<b>.0031</b>
Z <sub>2</sub>	X	y	2	24	8	<i>.0497</i>	<i>.0538</i>	<i>.0479</i>	<i>.0126</i>	<i>.0140</i>	<i>.0115</i>	.0023	<i>.0027</i>	<i>.0017</i>
				8	24	<b>.0023</b>	<i>.0428</i>	<i>.0700</i>	<b>0</b>	.0180	.0261	<b>0</b>	<b>.0057</b>	<b>.0080</b>
				8	16	<b>.0050</b>	<i>.0432</i>	<i>.0678</i>	<b>0</b>	.0174	.0252	<b>0</b>	<b>.0057</b>	<b>.0079</b>
	X	y	2	16	16	.0264	<i>.0450</i>	<i>.0624</i>	<b>.0011</b>	<i>.0139</i>	.0216	<b>0</b>	<b>.0033</b>	<b>.0048</b>
				16	8	<i>.0334</i>	<i>.0483</i>	<i>.0587</i>	<b>.0028</b>	<i>.0130</i>	.0193	<b>.0000</b>	<i>.0026</i>	<i>.0040</i>
				24	8	<i>.0395</i>	<i>.0502</i>	<i>.0570</i>	.0054	<i>.0113</i>	.0160	<b>.0001</b>	<i>.0022</i>	<i>.0032</i>
Z <sub>2</sub>	Y	x	2	8	24	<i>.0517</i>	<i>.0495</i>	<i>.0498</i>	<i>.0107</i>	<i>.0096</i>	<i>.0088</i>	<i>.0012</i>	<i>.0011</i>	<i>.0011</i>
				8	16	<i>.0498</i>	<i>.0494</i>	<i>.0486</i>	<i>.0105</i>	<i>.0100</i>	<i>.0085</i>	<i>.0012</i>	<i>.0010</i>	<i>.0009</i>
				16	16	<i>.0513</i>	<i>.0494</i>	<i>.0473</i>	<i>.0117</i>	<i>.0101</i>	<i>.0081</i>	.0018	<i>.0014</i>	<i>.0004</i>
	Y	x	2	16	8	<i>.0529</i>	<i>.0490</i>	<i>.0407</i>	<i>.0157</i>	<i>.0121</i>	.0063	.0025	<i>.0017</i>	<b>.0003</b>
				24	8	<i>.0561</i>	<i>.0481</i>	<i>.0384</i>	.0172	<i>.0131</i>	.0043	<b>.0037</b>	.0025	<b>.0002</b>
				F <sub>v</sub>	X	X	1	16	16	<b>.2725</b>	<b>.4809</b>	<b>.2687</b>	<b>.2134</b>	<b>.3944</b>
F <sub>v</sub>	X	Y	1	16	16	<b>.2788</b>	<b>.3256</b>	.1136	<b>.2244</b>	<b>.2416</b>	<b>.0440</b>	<b>.1965</b>	<b>.1996</b>	<b>.0110</b>
$\chi^2_v$	X			8		<b>.4383</b>	<b>.5543</b>	.1608	<b>.4313</b>	<b>.4919</b>	<b>.0819</b>	<b>.4290</b>	<b>.4555</b>	<b>.0364</b>

0 ≤ heavy < α/5 ≤ medium < α/3 ≤ light < 2α/3 ≤ italics ≤ 3α/2 < light ≤ 3α < medium ≤ 5α < heavy ≤ 1.00

**Table 2**  
**Robustness of t Tests: Empirical Probability  $\rho$  of a Type I Error Corresponding to Nominal Probability  $\alpha$  for Left-Tail (L), Two-Tail (T), and Right-Tail (R) Rejection Regions**

Statistic	Pop. 1	Pop. 2	$\sigma_1/\sigma_2$	$N_1$	$N_2$	$\alpha = .05$			$\alpha = .01$			$\alpha = .001$		
						L	T	R	L	T	R	L	T	R
$t_1$	X			8		<b>.4283</b>	<b>.4277</b>	<b>.0044</b>	<b>.4198</b>	<b>.4045</b>	<b>.0001</b>	<b>.3129</b>	<b>.2501</b>	<b>.0000</b>
						8	24	<b>.0038</b>	<i>.0414</i>	<i>.0737</i>	<b>.0002</b>	<i>.0118</i>	.0216	<b>.0000</b>
$t_2$	X	X	1	8	16	<i>.0101</i>	<i>.0350</i>	<i>.0638</i>	<b>.0009</b>	.0054	<i>.0101</i>	<b>.0001</b>	.0004	<i>.0008</i>
				16	16	<i>.0445</i>	<i>.0279</i>	<i>.0129</i>	<b>.0026</b>	<b>.0016</b>	<b>.0026</b>	<b>.0001</b>	<b>.0001</b>	<b>.0001</b>
				16	8	<i>.0638</i>	<i>.0350</i>	.0101	<i>.0101</i>	.0054	<b>.0009</b>	<i>.0008</i>	.0004	<b>.0001</b>
				24	8	<i>.0737</i>	<i>.0414</i>	<b>.0038</b>	.0216	<i>.0118</i>	<b>.0002</b>	<b>.0033</b>	.0017	<b>.0000</b>
$t_2$	X	Y	1	8	24	<i>.0441</i>	<i>.0336</i>	<i>.0377</i>	.0043	.0049	<i>.0067</i>	<b>.0002</b>	.0003	.0006
				8	16	<i>.0621</i>	<i>.0432</i>	.0293	<i>.0127</i>	<i>.0073</i>	.0035	<i>.0011</i>	.0006	<b>.0001</b>
				16	16	.0833	<i>.0607</i>	.0252	.0276	.0190	<b>.0021</b>	<b>.0066</b>	<b>.0042</b>	<b>.0000</b>
				16	8	.1055	.0853	.0293	<b>.0486</b>	<b>.0375</b>	.0041	<b>.0191</b>	<b>.0148</b>	.0005
$t_2$	X	x	2	8	24	<b>.2803</b>	<b>.1858</b>	.0756	<b>.0557</b>	<b>.0474</b>	.0188	<b>.0215</b>	<b>.0192</b>	.0020
				8	16	<b>.2849</b>	<b>.1796</b>	<i>.0410</i>	<b>.0800</b>	<b>.0615</b>	.0057	<b>.0377</b>	<b>.0322</b>	.0005
				16	16	<b>.1885</b>	<b>.1665</b>	.0138	<b>.1140</b>	<b>.0730</b>	<b>.0004</b>	<b>.0250</b>	<b>.0172</b>	<b>0</b>
				16	8	<b>.1804</b>	<b>.1513</b>	<b>.0007</b>	<b>.1029</b>	<b>.0702</b>	<b>0</b>	<b>.0222</b>	<b>.0131</b>	<b>0</b>
$t_2$	X	y	2	8	24	<b>.2288</b>	<b>.1783</b>	<i>.0651</i>	<b>.0710</b>	<b>.0469</b>	<i>.0119</i>	<b>.0081</b>	.0050	.0016
				8	16	<b>.2192</b>	<b>.1636</b>	<i>.0376</i>	<b>.0812</b>	<b>.0507</b>	.0046	<b>.0153</b>	<b>.0085</b>	.0003
				16	16	<b>.1550</b>	<b>.1262</b>	.0136	<b>.0904</b>	<b>.0694</b>	<b>.0004</b>	<b>.0358</b>	<b>.0261</b>	<b>0</b>
				16	8	.1389	.1123	<b>.0025</b>	<b>.0861</b>	<b>.0708</b>	<b>.0001</b>	<b>.0435</b>	<b>.0348</b>	<b>.0001</b>
$t_2$	Y	x	2	8	24	.1141	<b>.1711</b>	.1341	<b>.0431</b>	<b>.0721</b>	<b>.0613</b>	<b>.0108</b>	<b>.0239</b>	<b>.0226</b>
				8	16	.0877	<b>.1266</b>	.1098	.0261	<b>.0463</b>	<b>.0448</b>	<b>.0051</b>	<b>.0141</b>	<b>.0150</b>
				16	16	<i>.0441</i>	<i>.0520</i>	<i>.0582</i>	<i>.0079</i>	<i>.0123</i>	<i>.0147</i>	.0006	.0016	.0024
				16	8	.0214	.0183	.0239	<b>.0021</b>	<b>.0018</b>	<b>.0029</b>	<b>.0000</b>	<b>.0001</b>	<b>.0002</b>
$t_2$	Y	y	2	8	24	.0156	<b>.0080</b>	.0101	<b>.0017</b>	<b>.0007</b>	<b>.0007</b>	<b>.0001</b>	<b>.0001</b>	<b>0</b>
				8	16	.1180	<b>.1566</b>	.1168	<b>.0460</b>	<b>.0612</b>	.0451	<b>.0115</b>	<b>.0154</b>	<b>.0117</b>
				8	16	.0945	<b>.1141</b>	.0928	.0294	<b>.0371</b>	<b>.0312</b>	<b>.0056</b>	<b>.0069</b>	<b>.0058</b>
				16	16	<i>.0508</i>	<i>.0516</i>	<i>.0510</i>	<i>.0105</i>	<i>.0110</i>	<i>.0102</i>	<i>.0013</i>	<i>.0015</i>	<i>.0012</i>
$t_2$	Y	y	2	16	8	.0246	.0198	.0232	<b>.0028</b>	<b>.0022</b>	<b>.0029</b>	<b>.0002</b>	<b>.0001</b>	<b>.0001</b>
				24	8	<b>.0146</b>	<b>.0091</b>	<b>.0133</b>	<b>.0012</b>	<b>.0009</b>	<b>.0012</b>	<b>.0001</b>	<b>.0001</b>	<b>.0000</b>
				8	24	.1180	<b>.1566</b>	.1168	<b>.0460</b>	<b>.0612</b>	.0451	<b>.0115</b>	<b>.0154</b>	<b>.0117</b>
				8	16	.0945	<b>.1141</b>	.0928	.0294	<b>.0371</b>	<b>.0312</b>	<b>.0056</b>	<b>.0069</b>	<b>.0058</b>
$t_c$	X	Y	1	16	16	<i>.0645</i>	<i>.0527</i>	<i>.0375</i>	.0182	<i>.0139</i>	.0057	.0037	.0023	.0005

0 ≤ heavy <  $\alpha/5$  ≤ medium <  $\alpha/3$  ≤ light <  $2\alpha/3$  ≤ *italics* ≤  $3\alpha/2$  < light ≤  $3\alpha$  < medium ≤  $5\alpha$  < heavy ≤ 1.00

the two-tailed  $Z_1$  and  $Z_2$  tests at  $\alpha = .05$ , in which  $\rho$  usually fell within a 10% deviation from  $\alpha$  and always within 20%. The  $t_c$  test was also reasonably robust at  $\alpha = .05$ ; and when sample sizes were equal and only the assumption of equal variances was violated, the  $t_2$  test was robust at all  $\alpha$  values. However, it is clear from the tables that there are few easily and succinctly generalizable patterns of robustness.

Extreme departures of  $\rho$  from  $\alpha$  were common for most of the tests on means. Often  $\rho$  was many times greater than  $\alpha$  (more than 312 times as great in the worst case) or  $\alpha$  was many times greater than  $\rho$  (infinitely greater in the 20 cases in which  $\rho$  was 0). Yet even at the relatively innocuous  $\alpha$  level of .05, multiples exceeding 3 occurred for all the table entries for  $t_1$ , for a third of the entries for  $Z_1$  and  $t_2$ , and for some of the entries for  $Z_2$ . And only for the  $t_c$  test did all of the table entries at  $\alpha = .05$  fall in the interval from  $2\alpha/3$  to  $3\alpha/2$ . The latter is an ultraliberal criterion for robustness, since a lower boundary only three times as far from  $\alpha$

would include  $\rho = 0$  so that all  $\rho$  values below  $\alpha$  would be called robust.

Perhaps the most startling result of this study, however, is that there were numerous instances in which the nonrobustness of the reputedly robust  $Z$  and  $t$  tests on means rivaled or exceeded that of the notoriously nonrobust  $\chi^2_v$  or  $F_v$  tests on variances. If we measure nonrobustness by the larger of the two ratios  $\rho/\alpha$  and  $\alpha/\rho$  and make comparisons at the same  $\alpha$  value, we find that at all  $\alpha$  values the  $t_1$  test is about as nonrobust and the left-tailed  $Z_1$  test is more nonrobust than the  $\chi^2_v$  test at the same sample size; both tests are usually more nonrobust than the  $F_v$  test based on larger samples. Furthermore, there are 22 cases for  $t_2$  and 15 for  $Z_2$  in which their nonrobustness exceeds that of the median table entry for  $F_v$  tests at the same  $\alpha$  value, and 20%-40% of these cases occur at the most robustness-conducive  $\alpha$  value of .05. Finally, the nonrobustness of about 1/2 the table entries for  $t_2$ , 1/10 for  $Z_2$ , and 1/7 for  $F$  exceeds that of the smallest table entry for  $F_v$  at the

**Table 3**  
**Robustness of the Multisample F Test When N Observations Are Drawn From Each of Three or Four Sampled Populations**

Pop.1	Pop.2	Pop.3	Pop.4	Empirical Probability $\rho$ of a Type I Error Corresponding to Nominal Probability $\alpha$ When N = 8		
				$\alpha = .05$	$\alpha = .01$	$\alpha = .001$
X	X	X		.0248	.0036	.0003
X	X	X	X	.0314	.0052	.0006
X	Y	Y		.0604	.0156	.0026
X	Y	Y	Y	.0562	.0138	.0025
Y	X	X		.0563	.0216	.0083
Y	X	X	X	.0536	.0185	.0063
X	x	x		.1275	.0480	.0210
X	x	x	x	.0988	.0318	.0115
x	X	X		.1093	.0349	.0078
x	X	X	X	.0714	.0204	.0032
X	y	y		.1281	.0451	.0084
X	y	y	y	.1082	.0324	.0057
y	X	X		.0916	.0481	.0206
y	X	X	X	.0682	.0292	.0129
Y	x	x		.0674	.0236	.0061
Y	x	x	x	.0727	.0267	.0083
x	Y	Y		.0597	.0158	.0027
x	Y	Y	Y	.0601	.0161	.0024
Y	y	y		.0652	.0184	.0026
Y	y	y	y	.0709	.0222	.0040
y	Y	Y		.0573	.0151	.0021
y	Y	Y	Y	.0579	.0158	.0021

Larger of  $\rho/\alpha$  and  $\alpha/\rho$  is:  $\leq 1.5$ , 1.5 to 3, 3 to 5,  $> 5$

same  $\alpha$  value, and always at least 1/5 of these cases occur at  $\alpha = .05$ . Thus there were occasions when each of the tests on means except  $t_c$  was nonrobust to a degree encountered in the tests on variances, and these occasions were quite frequent for  $t_1$ ,  $t_2$ , and  $Z_1$ .

**DISCUSSION**

Objections are sometimes raised against using the same

relative criterion for robustness at all three  $\alpha$  levels rather than using increasingly lax criteria at increasingly remote rejection regions. However, if one has properly chosen his  $\alpha$  level as the maximum risk of Type I error that he is willing to countenance (after due consideration of power requirements and tradeoffs between risks of Type I and Type II errors), it is about as disastrous for that risk to be double what he thought when  $\alpha = .001$  as when  $\alpha = .05$ .

It is sometimes claimed that it is less serious for  $\rho$  to fall below  $\alpha$  than to fall above it because in the former case the actual probability of a Type I error is on the "safe" or "conservative" side of the nominal probability  $\alpha$ . However, this is no cause for complacency, since, in that case, the power of the test is reduced accordingly and the probability of a Type II error must therefore be on the "unsafe" or "radical" side of  $\beta$ . Therefore, ironically, a nonparametric test using the same nominal  $\alpha$  level and the same sample sizes may have considerably more power than the test actually selected because of the robustness and superior efficiency claimed for it.

As mentioned earlier, none of the conditions investigated in this study were unrealistic to actual experimentation in the behavioral sciences. On the contrary, most of the investigated conditions were quite mild. Consequently, the present study cannot be regarded as having explored the outer reaches of nonrobustness—neither the mathematical boundaries nor the extremes that might actually be encountered in practice—nor as having remotely approached such an exploration. Yet the majority of the  $\rho$  values found for each of the tests failed to meet a rather lax criterion of robustness. Worse, what robustness was found occurred sporadically and was highly peculiar to specific sampling conditions. Consequently, it is seldom possible to formulate an adequately, but not profusely, qualified statement of any set of conditions under which robustness is invariably attained by a given test, not only in general but even within the limited range of (usually mild) conditions investigated in the present study. This resistance of robustness to succinct generalization greatly weakens the usefulness of the term.

**REFERENCES**

BRADLEY, J. V. A common situation conducive to bizarre distribution shapes. *American Statistician*, 1977, 31, 147-150.  
 BRADLEY, J. V. Robustness? *British Journal of Mathematical and Statistical Psychology*, 1978, 31, 144-152.  
 BRADLEY, J. V. Nonrobustness in one-sample Z and t tests: A large-scale sampling study. *Bulletin of the Psychonomic Society*, 1980, 15, 29-32.  
 SCHEFFÉ, H. *The analysis of variance*. New York: Wiley, 1959. Pp. 336-337, 360.

(Received for publication February 18, 1980.)