

The Efficiency of Some Distribution-Free Tests

by Gottfried E. Noether *)

Samenvatting

Laten T_1 en T_2 twee toetsen zijn voor dezelfde hypothese $\theta = \theta_0$ betreffende de waarde van een parameter θ . Zij verder de onbetrouwbaarheidsdrempel van beide toetsen gelijk aan α en het onderscheidingsvermogen tegen de alternatieve hypothese $\theta = \theta_1$ gelijk aan $1 - \beta$. Indien toets T_1 nu n_1 waarnemingen vergt en toets T_2 n_2 waarnemingen, dan wordt de relatieve doeltreffendheid (Eng.: efficiency) van toets T_1 ten opzichte van toets T_2 (als toetsen voor $\theta = \theta_0$ tegen $\theta = \theta_1$) gegeven door: $e = n_2/n_1$. Indien men de waarde van θ_1 op een bepaalde wijze naar θ_0 laat convergeren bij toenemende n_1 , is het in vele gevallen, door gebruik te maken van een stelling van Pitman, mogelijk om een limiet-waarde voor e te vinden, die niet afhangt van α en β . Deze limiet-waarde wordt de asymptotische relatieve doeltreffendheid (volgens Pitman) genoemd. In dit artikel wordt een overzicht gegeven van hetgeen bekend is over de asymptotische relatieve doeltreffendheid van een aantal verdelingsvrije toetsen ten opzichte van de corresponderende standaardtoetsen.

De conclusie van de schrijver is, dat men bij het gebruik van verdelingsvrije methoden met een hoge doeltreffendheid (bijv. de symmetrietoets en de twee-steekproeven-toets van Wilcoxon, de toets van Kruskal voor k steekproeven en de methode van m rangschikkingen) slechts zeer weinig informatie kan verliezen en dat zelfs het gebruik van minder doeltreffende verdelingsvrije methoden gerechtvaardigd kan zijn.

1. Introduction

The use of distribution-free or, as they are often called, non-parametric tests has steadily increased over the last years. However, there are still many practicing statisticians who automatically rule out any use of distribution-free methods on the grounds that their use involves too great a loss of "information". It is then of some interest to have an actual comparison between distribution-free and the more familiar parametric tests, as well as comparisons of distribution-free methods among themselves.

One possibility of comparing two tests of the same hypothesis is in terms of relative efficiency. It is not the purpose of this paper to discuss theoretical aspects of the concept of efficiency of one test relative to another test. The

*) Boston University. This paper was written at the Mathematical Centre in Amsterdam (report SP 61 of the Statistical Department) while the author was conducting research with the partial support of the Office of Naval Research on the efficiency of nonparametric methods.

reader interested in such questions is referred to [8] and [13]. Rather, we are interested in bringing together in one place the considerable amount of information available in the statistical literature concerning the efficiency of the more well-known nonparametric methods. Before doing so, however, a few remarks concerning the concept of efficiency as applied to tests of hypotheses are in order.

When we say that the efficiency of test T_1 relative to test T_2 is e , we imply that the power of the first test using n observations is equal to the power of the second test using en observations. Or, putting it differently, the first test requires $1/e$ times as many observations in order to achieve as large a power with respect to a given alternative as the second test. The latter statement implies that the actual value of e may depend on the particular alternative we have in mind. Fortunately, this is often not the case, at least as long as we are dealing with large samples and alternatives which are reasonably "close" to the null hypothesis. In what follows, we shall always have this *asymptotic local efficiency*, as it is sometimes called, in mind.

In the next section, we shall take up certain classes of problems for which distribution-free solutions exist and discuss the kind of information available concerning the efficiency of these methods. When stating the efficiency of a distribution-free test relative to the corresponding parametric test, it is customary to assume that all the necessary assumptions for the validity of the latter are satisfied. Most of the time this implies that we are sampling from normal populations having equal variances. While this kind of comparison gives valuable information, it should be remembered that in practice we rarely have the assurance that our assumptions are actually satisfied. It then is important to know how the relative efficiency of the distribution-free test relative to the parametric test changes as a consequence of changed conditions. As far as it is available, we shall try to give this information.

A description of most tests to be discussed can be found, e.g., in [16]. Therefore, references will be limited to tests and results not mentioned in [16].

2. Efficiencies of Distribution-Free Tests

2.1. *One-Sample Tests.* The parametric form of the hypothesis tested is usually of the type $\mu = \mu_0$, where μ is the mean of a supposedly normally distributed population. The appropriate parametric test is based on Student's t . The distribution-free version of the hypothesis usually concerns the median.

The simplest distribution-free test is the sign test. In order to find its efficiency relative to the t -test, assume that the underlying distribution is

symmetric about θ with density $g(x-\theta)$, the null hypothesis being $\theta = 0$. Pitman¹⁾ has shown that the efficiency is given by

$$(1) \quad e = 4\sigma^2 g^2(0),^2$$

where σ^2 is the variance associated with $g(x)$. For the normal distribution, this reduces to the well-known value $2/\pi = .64$. Hodges and Lehmann [7] have shown that if $g(x)$ is unimodal, $e \geq \frac{1}{3}$, the minimum value being attained for the uniform distribution. On the other hand, there is no upper bound to e , since the sign test is applicable also for distributions having no finite variance.

A more powerful test of the same hypothesis is Wilcoxon's signed rank test, whose efficiency relative to the t -test according to Pitman is

$$(2) \quad e = 12\sigma^2 \left[\int g^2(x) dx \right]^2.$$

For a discussion of this quantity, see Sect. 2.3.

If observations become available sequentially, the sign test can be carried out as a sequential test of the hypothesis $p = p_0 = \frac{1}{2}$ against the alternative $p = p_1 > \frac{1}{2}$, say, where p is the parameter of a binomial distribution. Romani [15] has shown that the efficiency of this sequential sign test relative to the (nonsequential) t -test in the case of normal distributions is in the neighborhood of 1.3, the exact value depending on the probabilities of type I and type II errors associated with the sequential test. Not only does this distribution-free test require on the average a considerably smaller number of observations than the most powerful parametric test, but it is also much simpler from a computational point of view.

2.2. *Two-Sample Tests (Matched Observations)*. In practice, the tests of the previous section occur most frequently when testing the hypothesis that two populations from which we have paired observations are equal against the alternative that one population is shifted with respect to the other. If we denote the observations from the two populations by y_i and z_i , $i = 1, 2, \dots, n$, respectively, the tests are carried out by using the differences $x_i = z_i - y_i$ as the observed values. If the densities of \underline{y}_i ³⁾ and \underline{z}_i are denoted by $f(u)$ and $f(u-\theta)$, respectively, the density of \underline{x}_i may be written as

¹⁾ Professor Pitman has never published this, as well as the additional results attributed to him. They were presented in lectures given at Columbia University during 1948. Printed references can be found in [7] and [12].

²⁾ For the validity of (1), it is actually sufficient to assume that θ is the median of $g(x-\theta)$ without $g(x-\theta)$ being symmetric with respect to θ . However, the assumption of symmetry is required for the other statements made in this section.

³⁾ Chance variables are distinguished from observed values by underlining.

$$g(x - \theta) = \int f(u + x - \theta) f(u) du,$$

which is seen to be symmetric in x with respect to the point $x = \theta$. In particular, if the null hypothesis $\theta = 0$ is true,

$$(3) \quad g(x) = \int f(u + x) f(u) du.$$

The efficiencies given in the previous section now apply with $g(x)$ given by (3).

For the fixed sample sign test, the efficiency relative to the t -test can be expressed as a simple function of $f(u)$. Indeed, from (1) and (3), we easily find

$$(4) \quad e = 8\sigma^2 [\int f^2(u) du]^2,$$

where this time σ^2 refers to the variance associated with $f(u)$. For $f(u)$ normal, (4) again becomes $2/\pi$. A direct proof of (4) without the use of (1) can be found in the Appendix.

No simple expression in terms of $f(u)$ seems to be available for the efficiency of the Wilcoxon signed rank test, at least in the general case. When substituting (3) into (2), it should be remembered that the σ^2 of (2) is double the variance associated with $f(u)$. If $f(u)$ is normal, the resulting value of e is $3/\pi$.

For the sequential sign test, the efficiency statements remain the same as those of Section 2.1.

2.3. *Two-Sample Tests (General Case)*. In the parametric case, we are interested in testing the hypothesis $\mu_1 = \mu_2$ on the basis of two independent samples. Again the appropriate t -test is known to be most powerful if the necessary conditions for its application are satisfied. The nonparametric hypothesis usually takes the form $F_1(u) = F_2(u)$, where $F_1(u)$ and $F_2(u)$ are the cumulative distribution functions of the two underlying populations.

Probably the most useful two-sample test — and certainly the most thoroughly investigated — is Wilcoxon's test (or Mann-Whitney's test as it is often called). The parametric case suggests the following alternative: $F_2(u) = F_1(u - \theta)$, i.e., under the alternative hypothesis the two distributions differ only with respect to a location parameter. Pitman has investigated the efficiency of the Wilcoxon test relative to the t -test for this class of alternatives and found e to be given by

$$(5) \quad e = 12\sigma^2 [\int f^2(u) du]^2,$$

where $f(u) = F_1'(u)$. In case of normality, $e = 3/\pi = .95$.

Hodges and Lehmann [7] investigated the problem of finding the minimum value of (5) considered as a function of f and showed that always $e > .86$. We may conclude that whatever the underlying distribution, when testing the hypothesis $F_1(u) = F_2(u)$ against the alternative of a shift in

location, as far as our concept of efficiency is concerned, the Wilcoxon test is never much worse than the t -test and may be infinitely better.

In addition, Hodges and Lehmann considered what they called contamination alternatives:

$$F_1(u) = F(u); F_2(u) = (1 - \theta) F(u) + \theta G(u); G(u) \leq F(u), \text{ say.}$$

Again H_0 becomes $\theta = 0$. In this case

$$e = 12\sigma^2 \left(\frac{\int (F - G) dF}{\int (F - G) du} \right)^2.$$

Since the numerator is bounded while the denominator is not, there does not exist a lower bound to e . The further G is to the right of F , the smaller is the efficiency of the Wilcoxon test relative to the t -test in discovering this kind of alternative. As Hodges and Lehmann point out, in those cases where the contamination effect is due to gross errors, this low efficiency may rather be an advantage than a disadvantage.

A two-sample test which has found considerable acceptance in statistical textbooks is the Wald-Wolfowitz test based on the total number of runs of elements in the two samples. While the test is consistent with respect to all alternatives for which $F_1'(u) \neq F_2'(u)$, Pitman has pointed out that, when the two distributions differ only with respect to the value of some parameter, the efficiency of the run test relative to a properly constructed parametric test is zero. Even if very little is known about the type of alternative to be expected, the Kolmogorov-Smirnov test based on the maximum difference of the two sample cumulative distribution functions is preferable.

2.4. *Analysis of Variance Tests.* Kruskal and Wallis have given a distribution-free equivalent for the one-way classification of the analysis of variance. For two samples, their H -test reduces to the corresponding Wilcoxon two-sample test. As far as the efficiency of the H -test is concerned, Andrews [1] and Hodges and Lehmann [7] have shown that everything said about the efficiency of the Wilcoxon test in connection with alternatives involving a shift in location applies equally to the H -test.

For the two-way classification with one observation in each cell Friedman suggested the χ^2_r -test (also known as the method of m rankings) as a distribution-free alternative to the usual analysis of variance test. Friedman's test has been extended by Benard and Van Elteren [2] to cover also the case of unequal numbers of observations per cell. A particular

important special case is that of a balanced incomplete block design, already considered by D u r b i n [5].

Let the total number of treatments (columns) to be compared be t . In a balanced incomplete block design, in each block (ranking) only $k < t$ treatments are compared on the basis of one observation for each one of the k treatments involved. Moreover, considering all m blocks, each treatment is combined with any other treatment the same number of times.

V a n E l t e r e n and N o e t h e r [6] found the efficiency (for $m \rightarrow \infty$) of the rank test relative to the corresponding analysis of variance test to be

$$(6) \quad e = \frac{12k}{k+1} \sigma^2 \left[\int f^2(u) du \right]^2,$$

independent of t . In (6), $f(u)$ refers to the distribution of the error component in the underlying model. In particular, if $f(u)$ is normal, as usually assumed in the analysis of variance, (6) becomes

$$(7) \quad e = \frac{3}{\pi} \cdot \frac{k}{k+1}.$$

The efficiency of F r i e d m a n's test is obtained by substituting t for k in (6) and (7).

For $k = 2$, we have the case of a paired comparison test, which is asymptotically equivalent to the B r a d l e y-T e r r y test [3], whose efficiency thus turns out to be equal to (4). This is not surprising since, for $t = 2$, the F r i e d m a n test is equivalent to the sign test applied to matched observations.

2.5. *Tests of Independence and Regression.* The best known distribution free tests of independence are those based on K e n d a l l's and S p e a r m a n's rank correlation coefficients. K o n i j n [10] has investigated their efficiency relative to the product moment correlation coefficient for the following situation. Let \underline{y} and \underline{z} be two independent chance variables and define

$$\begin{aligned} v &= \lambda_1 \underline{y} + \lambda_2 \underline{z}, \\ \underline{w} &= \lambda_3 \underline{y} + \lambda_4 \underline{z}, \end{aligned}$$

where the λ_i , $i = 1, 2, 3, 4$, are constants. The hypothesis of independence of \underline{v} and \underline{w} can now be stated as $\lambda_2 = \lambda_3 = 0$. The following examples illustrate the kind of results obtained by K o n i j n. If \underline{y} and \underline{z} are normally distributed, the efficiency of either rank correlation coefficient relative to the product moment correlation coefficient is $\left(\frac{3}{\pi}\right)^2 = .91$. For the uniform distribution,

the efficiency is 1, while for the Laplace (symmetric exponential) distribution it is $\left(\frac{9}{8}\right)^2 = 1.27$.

Either rank correlation coefficient is also highly effective in testing for randomness against the alternative of a downward (upward) trend. S t u a r t [17] has investigated the efficiency of a great many distribution-free tests of randomness against normal regression alternatives. In particular, the efficiency of the rank correlation coefficients relative to the usual parametric test based on the regression coefficient is $\left(\frac{3}{\pi}\right)^{\frac{1}{3}} = .98$. The efficiency of the sign test comparing the observations in the first third of the sample with the corresponding observations in the last third is .83. On the other hand, the efficiency of the test based on the total number of runs up-and-down is zero. The same is true of the W a l d - W o l f o w i t z test based on the serial correlation coefficient.

If it is possible to carry out the test sequentially, the author's S_j -test [14] has efficiency slightly higher than 1.

2.6. *Goodness of Fit Tests.* The term goodness of fit test is customarily used if the hypothesis to be tested concerns the functional form of the underlying distribution, and not only the specific value of a given parameter. The classical test statistic for such a problem is χ^2 . However, there exist other test procedures which can be used under certain conditions. We shall mention only the K o l m o g o r o v test based on the maximum distance of the sample cumulative distribution function from the hypothetical population cumulative distribution function, since the distribution of this test statistic is better tabulated than the distributions of other similar statistics. The K o l m o g o r o v test is applicable if the hypothesis tested is simple. Investigations by M a s s e y [11] indicated that the K o l m o g o r o v test is more powerful than the corresponding χ^2 -test. These indications were borne out by a study by K a c, K i e f e r and W o l f o w i t z [9] in the case of tests of normality. While it would require too much detail to state their results precisely, the following statement should be sufficient. The three writers have shown that if the most powerful χ^2 -test requires n observations in order to have power $\geq \frac{1}{2}$ against a certain class of alternatives close to the null hypothesis, the number of observations required by the corresponding K o l m o g o r o v test has to be only of the order $n^{4/5}$. This implies that the efficiency of the χ^2 -test with respect to the K o l m o g o r o v test for this kind of problem is zero.

It may be objected that the K o l m o g o r o v test can only be used to

test normality when the mean and variance are assumed to be known. Actually, K a c, K i e f e r and W o l f o w i t z have extended the K o l m o g o r o v test to cover the case when the mean and variance are estimated from the available data. Nothing definite is known about how the two tests compare under this more general situation, since the power of the χ^2 -test is unknown in this case. However, it seems extremely plausible to assume that the situation is not very much changed from that when testing a simple hypothesis.

3. Conclusion

As pointed out in the Introduction, the efficiencies given in the previous section are asymptotic local efficiencies. The question then arises as to how these efficiencies change when we are dealing with small samples and/or alternatives relatively far removed from the null hypothesis.

Not much information is available on these two points except in the case of the sign test when the underlying distributions are assumed to be normal. For large samples, H o d g e s and L e h m a n n [7] have shown that the efficiency of the sign test relative to the t -test decreases very slowly from $2/\pi$ for alternatives close to the null hypothesis to .50 for alternatives far removed from the null hypothesis. Thus the value $2/\pi$ of the efficiency given in Section 2.1 turns out to be a good approximation for most alternatives of practical interest.

According to an investigation by D i x o n [4], the efficiency of the sign test decreases as the sample size increases, so that the asymptotic values given above actually represent minimum values of the efficiency.

While it is, of course, dangerous to generalize from results like these, this author feels that similar statements are true with respect to the efficiency of many tests discussed in Section 2.

One final remark seems to be in order before concluding the paper. The concept of efficiency as discussed in the Introduction is certainly of value when designing an experiment and deciding on the method of analysing the data to be obtained. On the other hand, once a given number of observations is available, it matters little how many more observations are required to make a less efficient method of analysis as powerful as a more efficient method. What matters then is the actual difference in power of the two methods of analysis. It is often true, particularly when large samples are involved, that the actual difference in power for a given sample size is quite small, even though there is a considerable difference in the efficiency, due to the fact that it may take a large number of additional observations to make up for even a small difference in power.

In conclusion, then, we may say that when using distribution-free methods

aving high efficiency instead of the more customary parametric methods, very little information, if any, is lost, and even low-efficiency distribution-free tests may have very legitimate uses.

4. Appendix

The purpose of this Appendix is to make somewhat more precise the meaning of asymptotic local efficiency and indicate a method for computing its value under conditions which are satisfied for many of the more customary tests. We shall also give an example of the computations involved.

4.1. *Pitman's Theorem.* Let $\underline{T}_n = T(x_1, \dots, x_n)$ be a test statistic for testing the hypothesis that a certain parameter θ has the value θ_0 . Let $ET_n = \psi_n(\theta)$ and $\text{var } T_n = \sigma_n^2(\theta)$. Pitman has called the quantity

$$R_n^2(\theta_0) = \frac{(\psi'_n(\theta_0))^2}{\sigma_n^2(\theta_0)}$$

the *efficacy* of \underline{T}_n for testing the hypothesis $\theta = \theta_0$ against the alternative $\theta = \theta_n = \theta_0 + k/\sqrt{n}$ where k is an arbitrary, but fixed constant. It is seen that the alternative θ_n approaches the null hypothesis θ_0 closer and closer as the sample size n increases. This is the reason for the term *asymptotic local efficiency*.

If we have two tests of the same hypothesis with efficacies $R_{1,n}^2(\theta_0)$ and $R_{2,n}^2(\theta_0)$, we can state Pitman's Theorem¹): The asymptotic efficiency of the first test relative to the second test is given by the limit of the ratio of the two efficacies,

$$e = \lim_{n \rightarrow \infty} \frac{R_{1,n}^2(\theta_0)}{R_{2,n}^2(\theta_0)}$$

4.2. *Example.* In Section 2.2, we considered the sign test for testing the hypothesis that two populations from which we have paired observations are equal against the alternative that the second population is shifted with respect to the first. More precisely, let the first population have density function $f_1(u) = f(u)$ where $f(u)$ is an arbitrary density function while the second population has density function $f_2(u) = f(u - \theta)$. We want to compute the efficiency of the sign test relative to the t -test for testing the hypothesis $\theta = \theta_0 = 0$, where both tests are based on the differences $x_i = z_i - y_i$, $i = 1, \dots, n$, of the paired observations y_i and z_i from $f_1(u)$ and $f_2(u)$, respectively.

¹) For the necessary regularity conditions, see [13]. A more general approach can be found in [8].

The sign test is based on, say, the number $T_{1,n}$ of positive values among the x_i 's. We easily find $\psi_{1,n}(\theta) = np$, $\sigma_{1,n}^2(\theta) = np(1-p)$ where $p = p(\theta)$ is the probability that a y -observation is smaller than the corresponding z -observation. Thus

$$\begin{aligned} p(\theta) &= P\{\underline{y} < \underline{z}\} = \int_{-\infty}^{\infty} \int_{-\infty}^z f(y)dy f(z-\theta)dz = \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{u+\theta} f(y)dy f(u)du, \end{aligned}$$

so that

$$\left. \frac{dp(\theta)}{d\theta} \right|_{\theta=0} = \int_{-\infty}^{+\infty} f^2(u)du.$$

Since $p(0) = \frac{1}{2}$,

$$R_{1,n}^2(0) = \frac{n^2[\int f^2(u)du]^2}{n/4} = 4n[\int f^2(u)du]^2.$$

Asymptotically, the t -test is equivalent to the test based on the statistic

$$T_{2,n} = \bar{x} = \frac{1}{n} \sum_{i=1}^n (z_i - y_i). \text{ We easily find}$$

$$\psi_{2,n}(\theta) = E\bar{x} = \theta,$$

$$\sigma_{2,n}^2(\theta) = \frac{2\sigma^2}{n},$$

where σ^2 is the variance associated with $f(u)$. It follows that

$$R_{2,n}^2(0) = \frac{1}{2\sigma^2/n} = \frac{n}{2\sigma^2},$$

and finally,

$$e = \lim_{n \rightarrow \infty} \frac{R_{1,n}^2(0)}{R_{2,n}^2(0)} = 8\sigma^2[\int f^2(u)du]^2,$$

as stated in (4).

References

- [1] F. C. A n d r e w s, "Asymptotic behavior of some rank tests for analysis of variance", Ann. Math. Stat., 25 (1954), 724-735.
- [2] A. B e n a r d and Ph. v a n E l t e r e n, "A generalization of the method of m rankings", Proc. Kon. Ned. Ak. v. Wet. A, 56, Indagationes Mathematicae, 15 (1953), 358-369.
- [3] R. A. B r a d l e y and M. E. T e r r y, "Rank analysis of incomplete block designs, I", Biometrika, 39 (1952), 324-345.
- [4] W. J. D i x o n, "Power functions of the sign test and power efficiency for normal alternatives", Ann. Math. Stat., 24 (1953), 467-473.

- [5] J. D u r b i n, "Incomplete blocks in ranking experiments", *British Journal of Psychology*, 4 (1951), 85—90.
- [6] Ph. v a n E l t e r e n and G. E. N o e t h e r, "The asymptotic efficiency of the χ^2 -test for a balanced incomplete block design", to be published.
- [7] J. L. H o d g e s, J r. and E. L. L e h m a n n, "The efficiency of some nonparametric competitors of the t -test", *Ann. Math. Stat.* 27 (1956), 324—335.
- [8] W. H o e f f d i n g and J. R o s e n b l a t t, "The efficiency of tests", *Ann. Math. Stat.*, 26 (1955), 52—63.
- [9] M. K a c, J. K i e f e r and J. W o l f o w i t z, "On tests of normality and other tests of goodness of fit based on distance methods", *Ann. Math. Stat.*, 26 (1955), 189—211.
- [10] H. S. K o n i j n, "On the power of certain tests of independence in bivariate populations", *Ann. Math. Stat.*, 27 (1956), 300—323.
- [11] F. J. M a s s e y, J r., "The Kolmogorov-Smirnov-test for goodness of fit", *Journal Am. Stat. Ass.*, 46 (1951), 68—78.
- [12] A. M o o d, "On the asymptotic efficiency of certain non-parametric two-sample tests", *Ann. Math. Stat.*, 25 (1954), 514—522.
- [13] G. E. N o e t h e r, "On a theorem of Pitman", *Ann. Math. Stat.*, 26 (1955), 64—68.
- [14] G. E. N o e t h e r, "Two sequential tests against trend", *Journal Am. Stat. Ass.*, 51 (1956), 440—450.
- [15] J. R o m a n i, "Tests no parametricos en forma secuencial", *Trabajos de Estadística*, 7 (1956), 43—96.
- [16] S. S i e g e l, *Nonparametric Statistics*, McGraw Hill, New York (1956).
- [17] A. S t u a r t, "The efficiencies of tests of randomness against normal regression", *Journal Am. Stat. Ass.*, 51 (1956), 285—287.

