

A New Sequential Non-Parametric Test for the Two Sample Problem

Keerthi M. Mathad*, I. D. Shetty

Department of Statistics, Karnataka University, Dharwad-580003, INDIA.

*Corresponding Address:

keerthimathad@gmail.com

Research Article

Abstract: A simple sequential non-parametric test for the two-sample problem is proposed. A method of deriving its ASN function for uniform and exponential distribution is given and their adequacy confirmed by simulation. The test is based on the normal approximation to the distribution of U-statistic. We also consider

the small sample sensitivity. We restrict ourselves to the Lehmann alternative. The test is found to be performing well even for small sample size.

Keywords: Sequential; non-parametric; two-sample problem; average sample number; Lehmann alternative.

1. Introduction

Phatarfod and Sudbury (1988) considered a Wald-type test for the two-sample problem. Unlike in Bradley, Merchant and Wilcoxon (1966) they exploited the normal approximation to the Wilcoxon Mann-Whitney statistic. They first considered the Lehmann alternative and later on wider type of alternatives.

In this paper we consider a Wald-type test for the two-sample problem. We take $H_0 : F(x) = G(x) = F_0(x)$ $H_1 : F(.) = F_0^{k_0}(.)$ $G(.) = F_0^{\frac{1}{k_0}}(.)$, $(0 < k_0 < 1)$. Here we wish to make a judgment between two new treatments and there is no prior preference involved. There would be a region of indifference $k_0 < k < k_0^{-1}$, where, if the true situation was $G(.) = F^k(.)$, it would not matter which of the two treatments was chosen. We propose the following Sequential test of strength (α, β) . Observations are taken from the X and Y populations in pairs. At each stage, the U-statistic

$$U = \frac{1}{\binom{n}{3} \binom{n}{3}} \sum_c h(X_{i_1}, X_{i_2}, X_{i_3}; Y_{j_1}, Y_{j_2}, Y_{j_3})$$

$$h(X_1, X_2, X_3; Y_1, Y_2, Y_3) = \begin{cases} 1, & \text{if } M(X_1, X_2, X_3) < M(Y_1, Y_2, Y_3) \\ 0, & \text{otherwise} \end{cases}$$

Where

$M(X_1, X_2, X_3)$ = Median of X_1, X_2, X_3 and $M(Y_1, Y_2, Y_3)$ = Median of Y_1, Y_2, Y_3
is calculated.

Sampling is continued as long as

$$B < \frac{f_{H_1}(U)}{f_{H_0}(U)} < A$$

Where

$$A = \frac{1-\beta}{\alpha} \quad \text{and} \quad B = \frac{\beta}{1-\alpha}$$

And H_0 or H_1 is accepted according as the R.H.S or the L.H.S inequality is the first not satisfied.

2. Mean and Variance of the U-statistic when $F(x) = G(x)$

$$\text{MEAN of } U \text{ is given by } E_{H_0}(U) = \frac{1}{2}.$$

VARIANCE of U is given by

$$\text{Var}(U) = \frac{1}{\binom{n}{3}^2} \sum_{c=0}^3 \sum_{d=0}^3 \binom{3}{c} \binom{m-3}{3-c} \binom{3}{d} \binom{n-3}{3-d} \xi_{c,d}$$

Where, $\xi_{c,d}$ = Covariance between the two kernels wherein cX observations and dY observations are in common between the two kernels.

$$\begin{aligned} \text{Var}(U) &= \frac{1}{\binom{n}{3} \binom{m}{3}} \sum_{c=0}^3 \sum_{d=0}^3 \binom{3}{c} \binom{m-3}{3-c} \binom{3}{d} \binom{n-3}{3-d} \xi_{c,d} = \frac{1}{\binom{n}{3}^2} \sum_{c=0}^3 \sum_{d=0}^3 \binom{3}{c} \binom{n-3}{3-c} \binom{3}{d} \binom{n-3}{3-d} \xi_{c,d} \\ &= \frac{1}{\binom{n}{3}^2} \left\{ \begin{array}{l} \left\{ \frac{(n-3)^2(n-4)^2(n-5)}{2} \right\} \left(\frac{131}{5775} + \frac{131}{5775} \right) + \left\{ (n-3)^2(n-4)(n-5) \right\} \left(\frac{97}{2100} + \frac{97}{2100} \right) \\ + \left\{ \frac{(n-3)(n-4)(n-5)}{3} \right\} \left(\frac{1}{12} + \frac{1}{12} \right) + \left\{ \frac{9(n-3)^2(n-4)^2}{4} \right\} \left(\frac{2323}{56700} \right) \\ + \left\{ 9(n-3)^2(n-4) \right\} \left(\frac{19}{252} + \frac{19}{252} \right) + \left\{ 9(n-3)^2 \right\} \left(\frac{3}{28} \right) + \left\{ 3(n-3)(n-4) \right\} \left(\frac{17}{140} + \frac{17}{140} \right) \\ + \left\{ 3(n-3) \right\} \left(\frac{83}{140} + \frac{83}{140} \right) + \frac{1}{4} \end{array} \right\} \end{aligned} \quad (2.1)$$

3. Mean and variance of the U-statistic when $F(x) = F^{k_0}(x)$ and $G(x) = F^{\frac{1}{k_0}}(x)$

$$\mu = E(U) = \frac{(19k_0^2 + 6)}{(k_0^2 + 1)(3k_0^2 + 2)(2k_0^2 + 3)}$$

MEAN is given by

VARIANCE of U is given by

$$\text{Var}(U) = \frac{1}{\binom{n}{3}^2} \sum_{c=0}^3 \sum_{d=0}^3 \binom{3}{c} \binom{m-3}{3-c} \binom{3}{d} \binom{n-3}{3-d} \xi_{c,d}^*$$

Where, $\xi_{c,d}^*$ = Covariance between the two kernels wherein cX observations and dY observations are in common between the two kernels.

$$\begin{aligned} \text{Var}(U) &= \frac{1}{\binom{n}{3}^2} \sum_{c=0}^3 \sum_{d=0}^3 \binom{3}{c} \binom{n-3}{3-c} \binom{3}{d} \binom{n-3}{3-d} \xi_{c,d}^* \\ &= \frac{1}{\binom{n}{3}^2} \left\{ \begin{array}{l} \left\{ \frac{(n-3)^2(n-4)^2(n-5)}{2} \right\} (\xi_{1,0}^* + \xi_{0,1}^*) + \left\{ (n-3)^2(n-4)(n-5) \right\} (\xi_{2,0}^* + \xi_{0,2}^*) \\ + \left\{ \frac{(n-3)(n-4)(n-5)}{3} \right\} (\xi_{3,0}^* + \xi_{0,3}^*) + \left\{ \frac{9(n-3)^2(n-4)^2}{4} \right\} (\xi_{1,1}^*) \\ + \left\{ 9(n-3)^2(n-4) \right\} (\xi_{1,2}^* + \xi_{2,1}^*) + \left\{ 9(n-3)^2 \right\} (\xi_{2,2}^*) + \left\{ 3(n-3)(n-4) \right\} (\xi_{1,3}^* + \xi_{3,1}^*) \\ + \left\{ 3(n-3) \right\} (\xi_{2,3}^* + \xi_{3,2}^*) + \xi_{3,3}^* \end{array} \right\} \end{aligned} \quad (3.1)$$

Where, $\xi_{1,0}^*, \xi_{0,1}^*, \xi_{2,0}^*, \xi_{0,2}^*, \xi_{3,0}^*, \xi_{0,3}^*, \xi_{1,1}^*, \xi_{1,2}^*, \xi_{2,1}^*, \xi_{2,2}^*, \xi_{1,3}^*, \xi_{3,1}^*, \xi_{2,3}^*, \xi_{3,2}^*, \xi_{3,3}^*$ are as defined earlier.
Details are omitted for the sake of brevity.

4. The Sequential Non-Parametric Test

We have developed the SPRT based on the statistic U using the normal approximation to our statistic U when n is large

for testing the simple hypothesis $H_0 : F(x) \equiv G(x)$ $H_1 : F(x) = F^{k_0}(x)$ vs $H_1 : G(x) = F^{\frac{1}{k_0}}(x)$.

I.e. the continuation region at stage m is given by

$$\log B - \log \frac{\sqrt{Var_{H_0}(U)}}{\sqrt{Var_{H_1}(U)}} < \frac{1}{2} \frac{(U - E_{H_0}(U))^2}{Var_{H_0}(U)} - \frac{1}{2} \frac{(U - E_{H_1}(U))^2}{Var_{H_1}(U)} < \log A - \log \frac{\sqrt{Var_{H_0}(U)}}{\sqrt{Var_{H_1}(U)}}$$

$$A = \frac{1-\beta}{\alpha} \quad \text{and} \quad B = \frac{\beta}{1-\alpha}$$

Where

5. The Average Sample Number

The sample size needed to reach a decision in a sequential or a multiple sampling plan is a random variable N, because at any stage of the experiment the decision to terminate the process depends on the results of the observations made earlier. The distribution of this random variable depends on the true distribution of the observations during the sampling process.

The ASN is given by

$$E(\mu(n)) = L(\mu) \log B + [1-L(\mu)] \log A / E(\mu(z)) \quad \text{where } \mu \text{ is the parameter}$$

$$\text{If } E(\mu(z)) = 0 \text{ then, } E(n) = -\log A \log B / E(z^2)$$

$$\text{Where } E(z^2) = (\mu_0 - \mu_1)^2 / \sigma^2 \quad \text{i.e. } E(n) = (-\log A \log B / (\mu_0 - \mu_1)^2) * \sigma^2$$

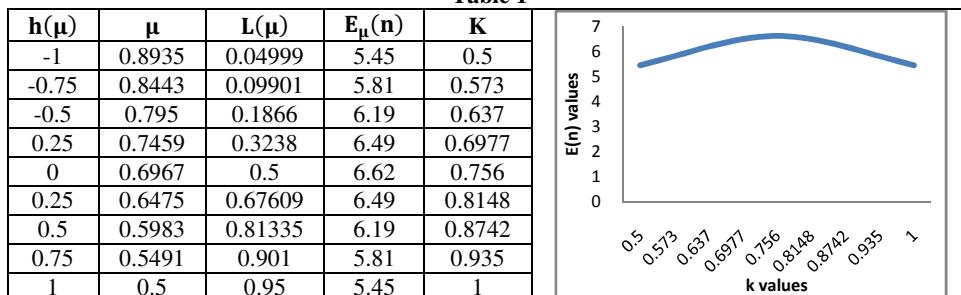
For this test the ASN function was obtained. To derive the ASN function, we note that the ASN function of our test is given by

$$E(N/K) = (L(K) \log B + [1-L(K)] \log A) / E(z) \quad \text{Where, } E(z) = \frac{[\mu_0^2 - \mu_1^2 - 2(\mu_0 - \mu_1) * E(x)]}{2 * \sigma_0^2}$$

The following tables and graphs give the ASN for different values of k_0

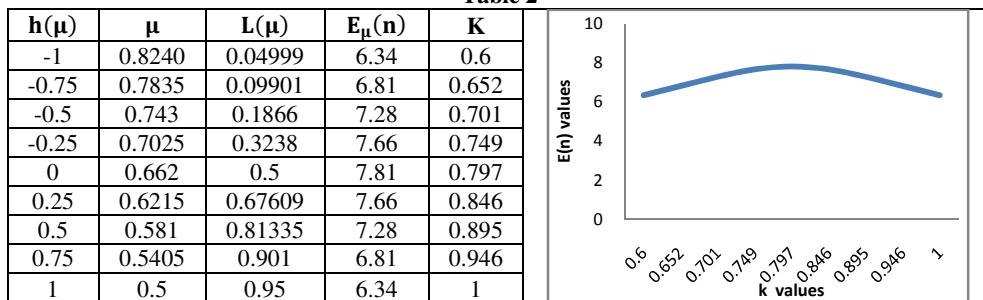
The values of ASN function for $k_0 = 0.5$ and the corresponding graph.

Table 1

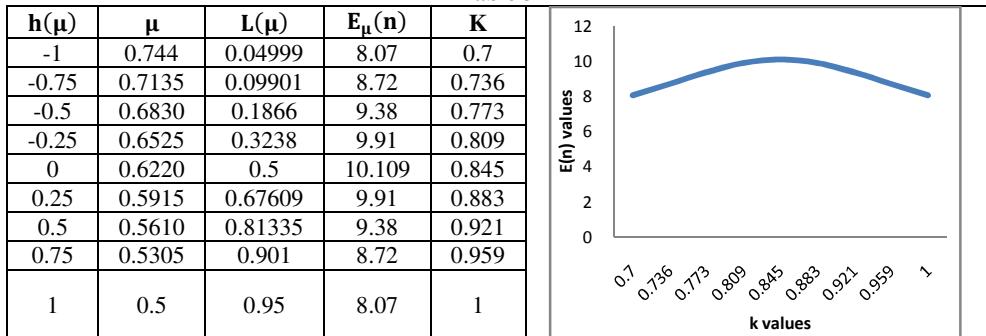


The values of ASN function for $k_0 = 0.6$ and the corresponding graph

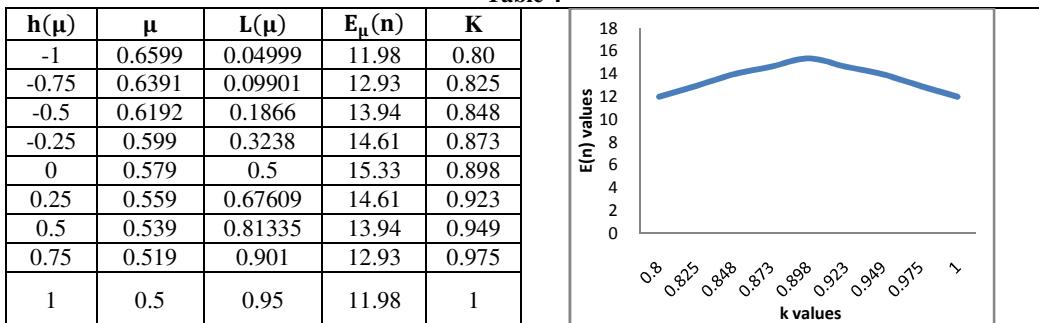
Table 2



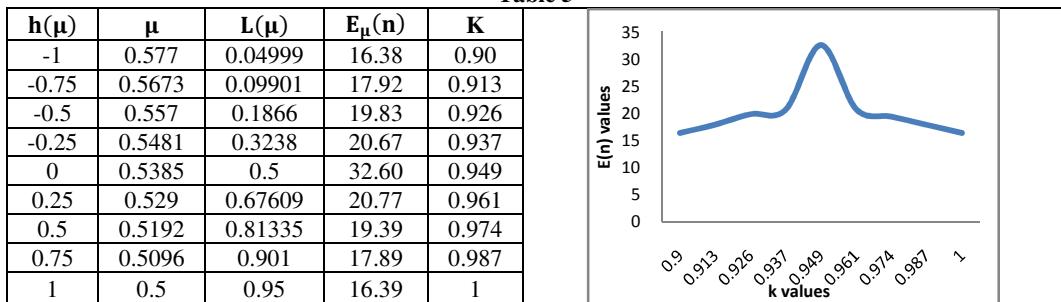
The values of ASN function for $k_0 = 0.7$ and the corresponding graph

Table 3

The values of ASN function for $k_0 = 0.8$ and the corresponding graph

Table 4

The values of ASN function for $k_0 = 0.9$ and the corresponding graph

Table 5

The test has greater ASN for the case $K_0=0.9$

6. Average Sample Number for Uniform Distribution

We have compared our test with other tests for the same problem. We have considered two distributions, namely Uniform and Exponential distributions. We have found the mean and variance for the two distributions and later found the ASN function.

We have the uniform distribution

$$f(x) = 1 \quad 0 < x < 1$$

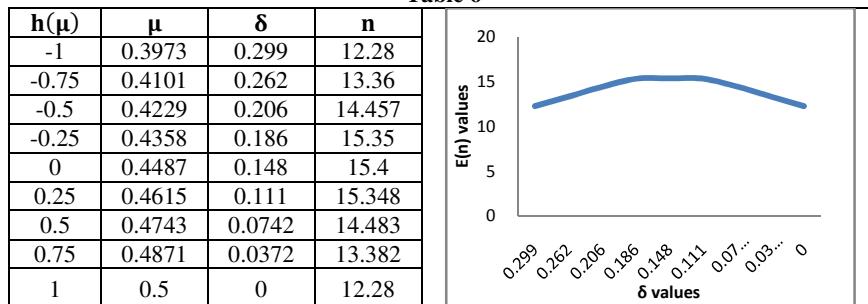
$$g(y) = 1 - \frac{\delta}{\sqrt{12}} \quad -\frac{\delta}{\sqrt{12}} < y < 1 - \frac{\delta}{\sqrt{12}}$$

We have obtained the mean and variance of the statistic U under uniform distribution. The calculation is omitted for the sake of brevity.

The ASN for different values of δ

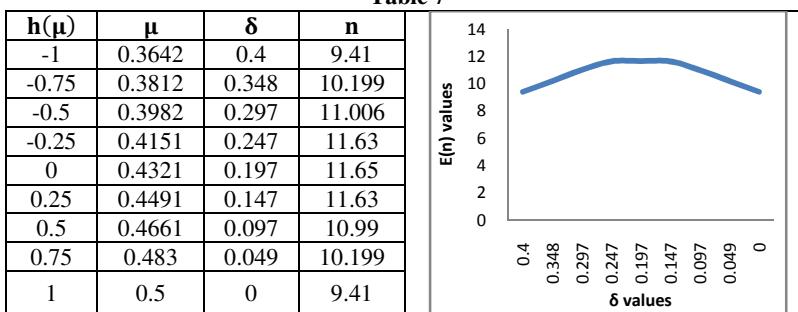
For $\delta = 0.3$

Table 6



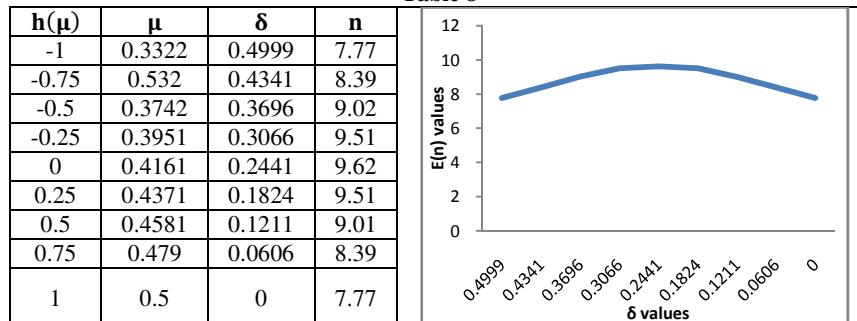
$\delta = 0.4$

Table 7



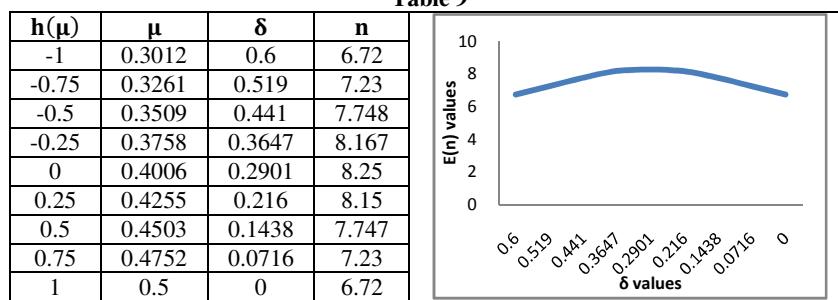
For $\delta = 0.5$

Table 8

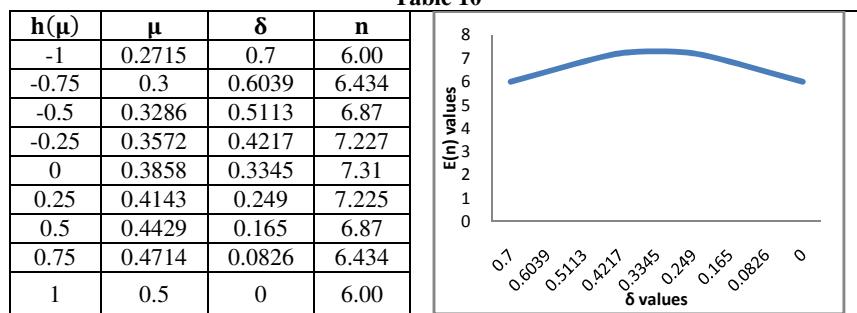


For $\delta = 0.6$

Table 9



For $\delta = 0.7$

Table 10

7. Average Sample Number for Exponential Distribution

We now consider the exponential distribution we have found the mean and variance for the distribution and later found the ASN function.

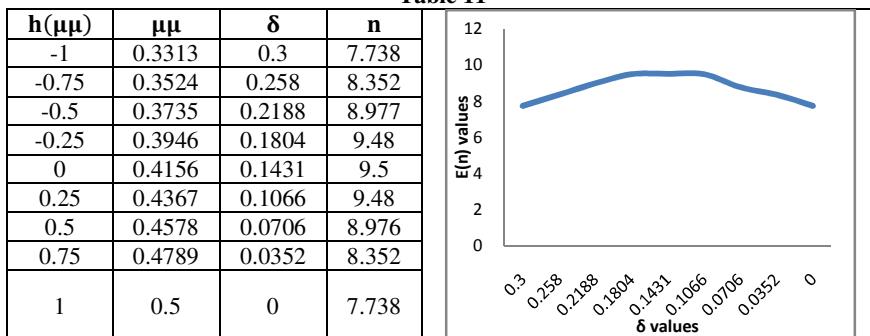
We have the exponential distribution

$$\begin{aligned} f(x) &= e^{-x} & 0 < x < \infty \\ g(y) &= e^{-(y-\delta)} & y \geq \delta \end{aligned}$$

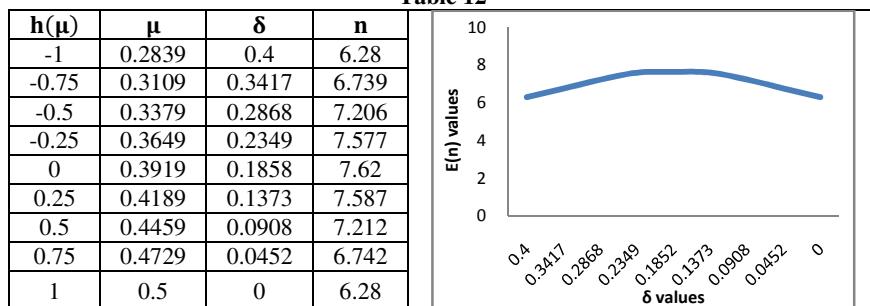
We have obtained the mean and variance of the statistic U under exponential distribution. The calculation is omitted for the sake of brevity.

The ASN for different values of δ

For $\delta = 0.3$

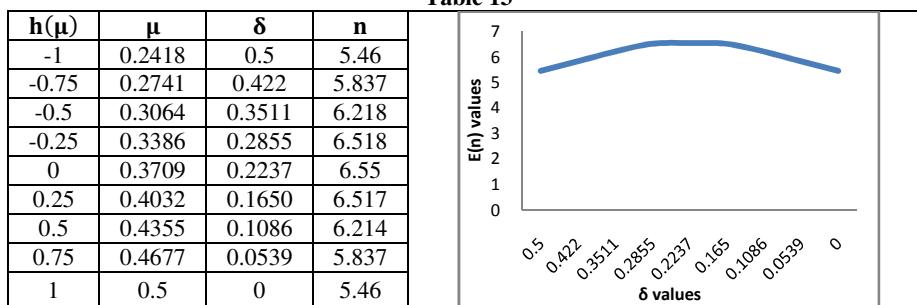
Table 11

For $\delta = 0.4$

Table 12

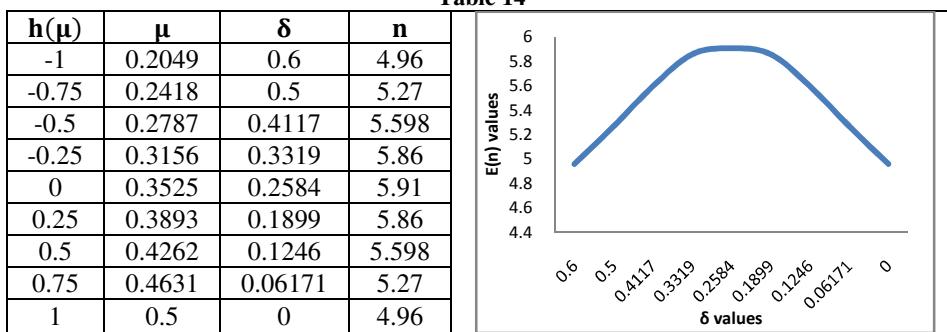
For $\delta = 0.5$

Table 13



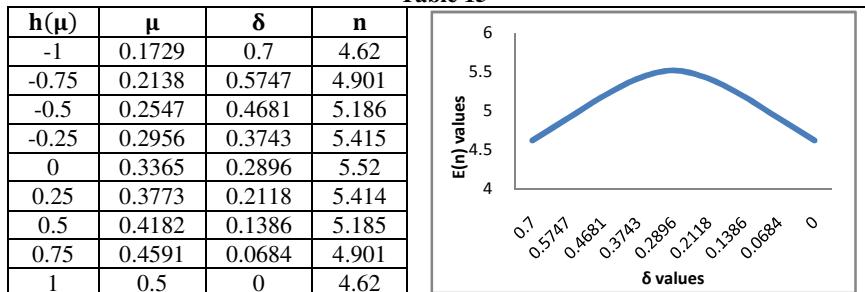
For $\delta = 0.6$

Table 14



For $\delta = 0.7$

Table 15



Conclusions

ASN for different values of δ

For $\delta = 0.3$		For $\delta = 0.4$		For $\delta = 0.5$		For $\delta = 0.6$		For $\delta = 0.7$	
Uniform	Exponential								
12.28	7.73	9.41	6.28	7.77	5.46	6.72	4.96	6.00	4.62
13.36	8.35	10.19	6.73	8.39	5.83	7.23	5.27	6.43	4.90
14.45	8.97	11.006	7.206	9.02	6.21	7.74	5.59	6.87	5.18
15.35	9.48	11.63	7.57	9.51	6.51	8.16	5.86	7.22	5.41
15.19	9.32	11.44	7.46	9.36	6.43	8.03	5.77	7.11	5.34
15.34	9.48	11.63	7.58	9.51	6.51	8.15	5.86	7.22	5.41
14.48	8.97	10.99	7.21	9.01	6.21	7.74	5.59	6.87	5.18
13.38	8.35	10.19	6.74	8.39	5.83	7.23	5.27	6.43	4.90
12.28	7.73	9.41	6.28	7.77	5.46	6.72	4.96	6.00	4.62

The results appear in the table. It is seen that for the exponential distribution ASN behaves better than the uniform in having smaller error probabilities and lower ASN, the ASN in general being less than half the fixed sample size for the test of the difference of two means from normal populations with equal known variance.

References

1. Abraham Wald (1947). Sequential Analysis, New York: John Wiley and Sons, Inc. Chapman and Hall, Ltd, London.
2. Bradley, R.A, Merchant, S.D. and Wilcoxon, F. (1966). Sequential rank tests II modified two sample procedures. *Technometrics*, 8, 615-623.
3. Cox, D.R. (1952). Sequential tests for composite hypotheses *Proceedings of the Cambridge Philosophical society*, 48, 290-299.
4. Ghosh, B.K.(1970). Sequential tests of Statistical Hypotheses. Reading, M.A; Addison-Wesley.
5. Gibbons, J.D and Subhabrato Chakraborti. (1992). Non-parametric Statistical Inference. Marcel Dekker, Inc. New York.
6. H.B. Mann and D.R. Whitney. (1947). on a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Statist*, 18, 50-60.
7. H.F. Dodge and H.G. Romig. (1929). A Method of Sampling Inspection, *The Bell System Technical Journal*, Vol.8, pp.613-631.
8. H.Scheffe. (1943). Statistical Inference in the Nonparametric case. *Ann. Math. Statist*. 14, 305-32.
9. Hoeffding, W. (1948). A class of statistics with asymptotic normal distribution, *Ann. Math. Statist*, 19, 293-325.
10. Lai, T.L. (1975). On Chernoff-Savage Statistics and Sequential rank tests. *Ann. Statist*. 3, 825-845.
11. Lehmann, E.L. (1951). Consistency and Unbiasedness for certain non-parametric tests, *Ann. Math. Statistics*, 22, 165-179.
12. Lehmann, E.L. (1975). Non-parametric, San Francisco: Holden-Day.
13. Mathisen, H.C. (1943). A method of testing the hypothesis that two samples are from the same population. *Ann. Math. Statistics*, 14, 188-194.
14. Miller, R. G. Jr. (1969). Sequential signed-rank test. *J. Amer. Statist. Assoc.* 65, 1554-1561.
15. Miller, R. G. Jr. (1972). Sequential rank tests-one sample case. *Proc. Sixth Berkeley Symp.Math Statist.Prob.* 1, 97-108.
16. Phatarfod, R.M. and Aidan Sudbury. (1988). A simple Sequential Wilcoxon test. *Austral. J. Statist*, 30(1), 93-106.
17. Pitman, E. J.G. (1948). Notes on non-parametric statistical inference. Columbia University.
18. Randles, R. H. and Wolfe, D. A. (1979). Introduction to the Theory of Non parametric Statistics, Wiley, New York.
19. Savage, I. R, and Sethuraman, J. (1966). Stopping time of a rank order sequential probability ratio test based on Lehman alternatives. *AMS* 37, 1154-1160.
20. Shetty, I.D, and Z. Govindarajalu. (1988). A two-sample test for location, *Comm. in Statistics-Theory and Methods*, 17, 2389- 2401.
21. Shetty, I.D. and Bhat, S.V. (1993). Some Competitors of Mood's Median test for the location alternative, *Journal of Karnataka University-Science*, 37, 11, 138-146.
22. Wilcoxon, F, Rhodes, L.J. and Bradley, R.A. (1963). Two Sequential two-sample grouped rank tests with applications to screening experiments, *Biometrics*, 19, 58-84.