

DISTRIBUTION-FREE LIMITS INVOLVING RANDOM  
VARIABLES

1050 710

by

TZE-U CHO

B. A., National Taiwan University, 1969

---

A MASTER'S REPORT

submitted in partial fulfillment of the

requirements for the degree

MASTER OF SCIENCE

Department of Statistics

KANSAS STATE UNIVERSITY

Manhattan, Kansas

1974

Approved by:



Major Professor

LD  
2668  
R4  
1974  
C56  
C-2  
Document

## TABLE OF CONTENTS

	page
INTRODUCTION . . . . .	1
CONFIDENCE INTERVAL FOR QUANTILES . . . . .	3
CONFIDENCE INTERVAL FOR MEDIAN . . . . .	9
CONFIDENCE BAND FOR POPULATION DISTRIBUTION FUNCTION .	18
TOLERANCE LIMITS . . . . .	21
REFERENCES . . . . .	27
APPENDIX : TABLES . . . . .	29
ACKNOWLEDGEMENT	

## LIST OF TABLES

Table	page
IA. Binomial Distribution, $p=.75$ . . . . .	30
IB. Binomial Distribution, $p=.50$ . . . . .	31
IC. Binomial Distribution, $p=.25$ . . . . .	32
II. Normal Distribution. . . . .	33
III. Confidence Coefficients for Confidence Limits on the Median Using Extremes of Samples of Size $n$ . . . . .	34
IV. Critical Values of the Wilcoxon Signed Ranks Test Statistic . . . . .	35
V. Some One-sided and Symmetrical Significance Tests for $n \leq 15$ . . . . .	35
VI. Critical Values of the Kolmogorov Test Statistic. . . . .	37
VIIA. Sample Sizes for One-sided Tolerance Limits. .	38
VIIIB. Sample Sizes for Two-sided Tolerance Limits. .	39
VIII A. Graphs of Population Coverage for Tolerance Level .99. . . . .	40
VIII B. Graphs of Population Coverage for Tolerance Level .95. . . . .	41
VIII C. Graphs of Population Coverage for Tolerance Level .90. . . . .	42

## INTRODUCTION

Many common statistical methods are concerned with estimating the parameters of a distribution function of known or assumed form. Thus the population may be assumed to be normal with parameters mean and variance, and these may be estimated by the sample mean and sample variance; or we may use a t-test or a F-test to make some test about them. Such statistical methods, since they deal with population parameters, are called parametric.

There are, however, other methods which do not require any assumptions about the nature of the population from which the sample is drawn. These are called distribution-free, since they are free of specific assumptions about the distribution in the parent population. These methods are valid for any parent population, hence they also could be validly applied to samples from normal distributions. It is true that these methods may be less efficient, in some sense, than the parametric methods when the population is normal. But for some kinds of data, particularly in behavioral science and engineering experiments, it would be inappropriate to assume normal distributions. In such cases, distribution-free methods should be used. Median, quartiles, etc., may be far more appropriate to describe the location and dispersion of the population than the mean and/or variance. Besides, distribution-free methods are generally simple to apply, not involving much computation.

As is well-known, it is convenient to use the t and  $\chi^2$  distributions to find confidence limits on the parameters of a

normal distribution. In this paper we will discuss some methods of finding the confidence limits when no assumption is suitable to be made about the parent population.

### CONFIDENCE INTERVAL FOR QUARTILES

There are three quartiles, usually designated  $Q_1$ ,  $Q_2$  or median,  $Q_3$ , which divide a given series of measurements into four equal parts. The first quartile,  $Q_1$ , is the value at or below which one-fourth of all items in the series fall; the second quartile,  $Q_2$ , is the value such that half of the items are smaller and half of the items exceed it aside from ties; and the third quartile,  $Q_3$ , is the value at or below which three-fourths of the items lie. The simplest way to find confidence intervals for quartiles is to use the concept of the binomial distribution. Let  $X_1, X_2, \dots, X_n$  present a random sample from a population, and arrange the  $n$  observations in increasing order:  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(r)} \leq \dots \leq X_{(s)} \leq \dots \leq X_{(n)}$ , where  $1 \leq r \leq s \leq n$ . Suppose the distribution is continuous and  $1 - \alpha$  is the desired confidence coefficient. To find the confidence interval for an unknown quartile is to find statistics  $r$  and  $s$  such that

$$(1) \quad P(X_{(r)} \leq Q_i \leq X_{(s)}) = 1 - \alpha$$

Because of continuity, (1) also can be written as

$$(2) \quad \begin{aligned} P(X_{(r)} \leq Q_i \leq X_{(s)}) &= P(Q_i \leq X_{(s)}) - P(Q_i < X_{(r)}) \\ &= P(Q_i < X_{(s)}) - P(Q_i < X_{(r)}) \end{aligned}$$

Consider the statement  $(Q_i < X_{(1)})$  where  $X_{(1)}$  is the smallest value in the sample. This statement is true only if all  $n$  observations in the sample are greater than  $Q_i$ .

Next consider  $(Q_i < X_{(2)})$ . The statement is true either when  $n-1$  observations are greater than  $Q_i$ , in which case  $X_{(1)} \leq Q_i < X_{(2)}$ ; or all observations are greater than  $Q_i$ , in which case  $Q_i < X_{(1)} < X_{(2)}$ . A picture can make the situation clear.

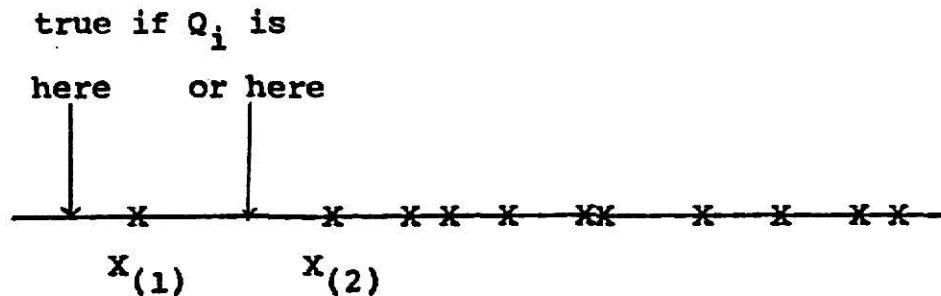


Figure 1

Each random  $X_i$  has probability  $p_1=.25$  of being less than or equal to  $Q_1$ ; has probability  $p_2=.50$  of being less or equal to  $Q_2$ ; and has probability  $p_3=.75$  of being less or equal to  $Q_3$ . Therefore, with the aid of the binomial distribution function we add the appropriate probabilities and obtain

$$\begin{aligned}
 (3) \quad P(Q_i < X_{(1)}) &= P(\text{all } X_j \text{ exceed } Q_i) \\
 &= P(\text{none of } X_j \text{ below } Q_i) \\
 &= (1 - p_i)^n \quad i=1, 2, 3.
 \end{aligned}$$

$$\begin{aligned}
 (4) \quad P(Q_i < X_{(2)}) &= P(Q_i < X_{(1)}) + P(X_{(1)} \leq Q_i < X_{(2)}) \\
 &= P(\text{at least } n-1 \text{ } X_j \text{ exceed } Q_i) \\
 &= P(1 \text{ or none of the } X_j \text{ below } Q_i) \\
 &= \sum_{a=0}^1 \binom{n}{a} p_i^a (1-p_i)^{n-a} \quad i=1, 2, 3.
 \end{aligned}$$

This result can be extended as follows:

$$\begin{aligned}
 (5) \quad P(Q_i < X_{(r)}) &= P(\text{at least } n-r+1 \text{ of } X_j \text{ exceed } Q_i) \\
 &= P(r-1 \text{ or less of } X_j \text{ below } Q_i) \\
 &= \sum_{a=0}^{r-1} \binom{n}{a} p_i^a (1-p_i)^{n-a} \quad i=1, 2, 3.
 \end{aligned}$$

$$\begin{aligned}
 (6) \quad P(Q_i < X_{(s)}) &= P(\text{at least } n-s+1 \text{ of } X_j \text{ exceed } Q_i) \\
 &= P(s-1 \text{ or less of } X_j \text{ below } Q_i) \\
 &= \sum_{a=0}^{s-1} \binom{n}{a} p_i^a (1-p_i)^{n-a} \quad i=1, 2, 3.
 \end{aligned}$$

Tables IA, IB, IC are the binomial distribution tables with  $p=.25, .50, .75$ . We can use these tables to find the confidence intervals for the three quartiles.

For the sample size  $\leq 20$ , Table I is used to find the value of  $r$  and  $s$ . From previous discuss we know that

$$\begin{aligned}
 (7) \quad 1 - \alpha &= P(X_{(r)} \leq Q_i \leq X_{(s)}) \\
 &= P(Q_i < X_{(s)}) - P(Q_i < X_{(r)})
 \end{aligned}$$

In order to obtain a two-sided confidence interval, set

$$(8) \quad P(Q_i < X_{(s)}) = 1 - \alpha/2$$

$$(9) \quad P(Q_i < X_{(r)}) = \alpha/2$$

Enter Table I with proper probability  $p$  and sample size  $n$ . Read across the row until reaching a value close to  $\alpha/2$ , the corresponding value of  $m$  is  $r-1$ . Add 1 to get  $r$ . Then continue reading across the same row until reaching a value close to  $1-\alpha/2$ , and add one to the corresponding  $m$  to get  $s$ .



When the sample size exceed 20, use the normal distribution to approximate the binomial distribution with mean=np and standard deviation  $\sqrt{np(1-p)}$ . The formulas are

(10)  $r^* = np + Z_{\alpha/2} \sqrt{np(1-p)}$

and

(11)  $s^* = np + Z_{\alpha/2} \sqrt{np(1-p)}$

Where the Z is obtained from the cumulative standard normal distribution table (Table II). The values  $r^*$ ,  $s^*$  may not always be integers so round them upward to the next higher ingeger to get r and s.

A one-sided confidence interval with confidence coefficient 1- can be found by using the above method to find either r or s, so that

(12)  $P(Q_i < X_{(s)}) = 1 - \alpha$

or

(13)  $P(Q_i < X_{(r)}) = \alpha$

The method described above applies only for a continous variate. In practice the distribution function may be discrete. Scheffé and Tukey (1945) have investigated confidence intervals for quantiles and have noted that if in the discrete case the confidence interval bounded by the appropriate ordered statistic has confidence coefficient exactly equal 1- $\alpha$ , the corresponding closed confidence interval (with endpoints) has confidence coefficient at least 1- $\alpha$ , whereas the open interval has confidence coefficient of 1- $\alpha$  at most. Therefore we shall write the confidence interval for quartiles as

$$(14) \quad P(X_{(r)} \leq Q_i \leq X_{(s)}) \geq 1 - \alpha$$

Example 1: An experiment is made by the fibers division of a chemical company to determine the breaking strength of corn yarn, represented by pounds per square inch required to break a skein of yarn. A random sample of size 20 is taken, and the data obtained are:

98, 112, 108, 86, 124, 92, 102, 91, 95, 104, 89, 129,  
83, 98, 92, 99, 113, 116, 122 and 85.

We wish to find the 95% confidence interval on  $Q_3$ .

First arrange the observed data in an ordered array:

83, 85, 86, 89, 91, 92, 92, 93, 95, 98, 98, 99, 102,  
104, 108, 112, 113, 116, 122, 129.

Enter Table IA with  $n=20$ , reading across the row, the probability closest to .025 is .014. The  $m$  value corresponding to .014 is 10. Therefore  $s$  equals 11. In the same row the probability is chosen as being close to .975 is .976. The corresponding  $m$  value is 18, therefore  $s=19$ . The confidence interval

$$P(X_{(11)} \leq Q_i \leq X_{(19)}) = .976 - .014 = .962 \\ \geq 95\%$$

It is reasonable to consider this case as a discrete distribution, because  $X_{(11)}$  equals 98,  $X_{(19)}$  equals 122, we may say "the interval from 98 to 122 pounds is at least a 95% confidence interval for the third quartile."

Example 2: Find a 95% confidence interval for the  $Q_2$  of the population of blood lactates of patients with an anxiety

neurosis, given the following random sample of size 27 from the population:

32, 36, 42, 33, 49, 98, 51, 46, 24, 51, 56, 45, 47, 51,  
56, 24, 95, 22, 31, 34, 38, 44, 49, 52, 54, 42, 57 milligrams  
per 100 millilitres.

Arrange them in order, thus:

22, 24, 24, 31, 32, 33, 34, 36, 38, 42, 42, 44, 45, 46,  
47, 49, 49, 51, 51, 51, 52, 54, 56, 56, 57, 95, 98.

Because the sample is greater than 20, use Equation (10)  
and (11).

$$\begin{aligned} r^* &= (27)(.50) + (-1.96)(27)(.50)(.50) \\ &= 13.5 - 5.09 \\ &= 8.41 \\ s^* &= 13.5 + 5.09 \\ &= 18.59 \end{aligned}$$

Therefore  $r$  equals 9,  $s$  equals 19. Hence the 95% confidence interval on the median is:  $38 \leq Q_2 \leq 51$  mg/100ml.

## CONFIDENCE INTERVAL FOR MEDIAN

Median is identical with  $Q_2$ . In a list of values arranged in increasing magnitude, half of them at or above median, half at or below it. We will discuss three ways of finding the confidence interval for the median in this section.

The first method is based on binomial distribution. This method resembles the method shown in the previous section of finding a confidence interval for  $Q_2$ . The only difference is that we now take the same number of items on both tails. The merit of this method is ease of calculation. We can find out the confidence interval easily even without the aid of tables.

Obtain a random sample  $X_1, X_2, \dots, X_n$  from a population with unknown median, arrange them in order in two ways, denoted as follow:

$$(1) \quad X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n-1)} \leq X_{(n)}$$

$$(2) \quad X^{(n)} \leq X^{(n-1)} \leq \dots \leq X^{(2)} \leq X^{(1)}$$

In Equation (1) we express the smallest value as  $X_{(1)}$ , second smallest as  $X_{(2)}$  .... etc. In Equation (2) we express the largest values as  $X^{(1)}$ , second largest as  $X^{(2)}$  .... etc.

From definition of the median we know that the probability that a random  $X_i$  will be larger than (or smaller than) the median is  $\frac{1}{2}$ . The probability that the median is between  $X_{(c)}$  and  $X^{(c)}$  for a given  $c$  can be calculated.

The statement  $(X_{(c)} \leq md \leq X^{(c)})$  (denoting the true population median by  $md$ ) is false if either less than  $c$  observations are above median, in which case  $md > X^{(c)}$ ; or less than  $c$  observations are below median, in which case  $md < X_{(c)}$ . Hence, by adding appropriate binomial probabilities we can obtain the error probability and confidence probability.

$$\begin{aligned}
 (3) \quad \text{error probability} &= P(md > X^{(c)}) + P(md < X_{(c)}) \\
 &= 2 \sum_{a=0}^{c-1} \binom{n}{a} \left(\frac{1}{2}\right)^n \\
 &= \sum_{a=0}^{c-1} \binom{n}{a} \left(\frac{1}{2}\right)^{n-1}
 \end{aligned}$$

$$\begin{aligned}
 (4) \quad \text{confidence probability} &= P(X_{(c)} \leq md \leq X^{(c)}) \\
 &= 1 - \sum_{a=0}^{c-1} \binom{n}{a} \left(\frac{1}{2}\right)^{n-1}
 \end{aligned}$$

For example, to find the confidence level of the interval  $(X_{(3)}, X^{(3)})$

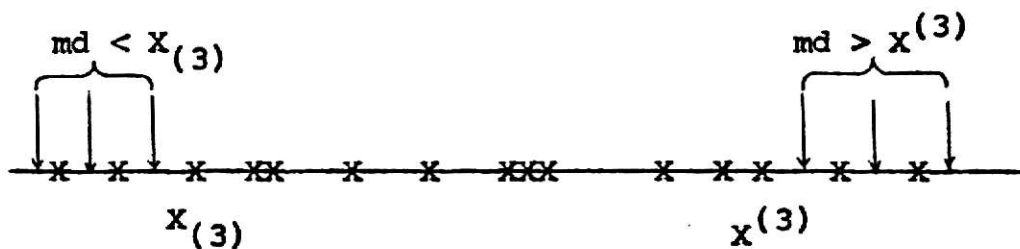


Figure 2

$$\text{error probability} = \sum_{a=0}^2 \binom{n}{a} \left(\frac{1}{2}\right)^{n-1}$$

$$\begin{aligned}
 \text{confidence probability} &= 1 - \sum_{a=0}^2 \binom{n}{a} \left(\frac{1}{2}\right)^{n-1} \\
 &= 1 - (n^2 + n + 2) \left(\frac{1}{2}\right)^n
 \end{aligned}$$

The confidence coefficient for a continuous distribution is exactly  $1-\alpha$ , but the confidence coefficient for a discrete distribution is at least  $1-\alpha$ . Therefore, the general formulas for confidence probabilities for intervals:  $X_{(c)} \leq md \leq X^{(c)}$  are:

$$\begin{aligned}
 (5) \quad c=1 & \quad \text{confidence probability} \geq 1 - \left(\frac{1}{2}\right)^{n-1} \\
 c=2 & \quad \text{confidence probability} \geq 1 - (n+1)\left(\frac{1}{2}\right)^{n-1} \\
 c=3 & \quad \text{confidence probability} \geq 1 - (n^2+n+2)\left(\frac{1}{2}\right)^n \\
 c=4 & \quad \text{confidence probability} \geq 1 - (n^2+5n+6)/3\left(\frac{1}{2}\right)^n \\
 & \quad \vdots \\
 & \quad \vdots \\
 & \quad \vdots
 \end{aligned}$$

Table III shows the confidence level for sample sizes up to 40, and all possible intervals. If  $1-\alpha$  is the chosen confidence coefficient, enter the table with sample size  $n$ , find a probability closest to but not smaller than  $1-\alpha$ , and read the corresponding  $c$  from the top of the table.

Example: Twenty electron tubes were tested and the following life times (hours) were reported:

7.2, 37.7, 49.6, 21.4, 67.2, 41.1, 3.8, 8.1, 23.2, 72.2,  
11.4, 17.5, 29.8, 57.8, 84.6, 12.8, 2.9, 42.7, 7.4, 33.4.

Find the 95% confidence interval for the population median.

List data in increasing order:

2.9, 3.8, 7.2, 7.4, 8.1, 11.4, 12.8, 17.5, 21.4, 23.2,  
29.8, 33.4, 37.8, 41.1, 42.7, 49.6, 57.8, 67.2, 72.1, 84.6.

Enter Table III with sample size  $n=20$ , the probability .96 expressed as a percentage is found under  $c=6$ . We find that  $X_{(6)}=11.4$ ,  $X^{(6)}=42.7$ ; therefore, the interval from 11.4 to

42.7 hours is a 96% confidence interval for the median life of electron tubes. Or, we could write:

$$CI_{96} : 11.4 \leq md \leq 42.7 \text{ hours}$$

where  $md$  is true average life of the population of electron tubes sampled.

The second method is based on the Wilcoxon signed rank test. The important assumptions of the Wilcoxon test are that the random variable  $X$  is continuous and its distribution is symmetric. Thus, the left half of the graph of the probability function is the same as the right half. So the median is also the mean, because both are located exactly in the middle of the distribution at the line of symmetry. Because of these assumptions, especially the symmetry assumption, a relatively short confidence interval can be found.

The observations  $X_1, X_2, \dots, X_n$  are a random sample from a population. To obtain the  $1-\alpha$  confidence interval, find the critical value  $W_{\alpha/2}$  from Table IV, compute all  $n(n+1)/2$  possible averages of  $(X_i+X_j)/2$  for all  $i$  and  $j$ , including  $i=j$ . A confidence interval is bounded by the  $W_{\alpha/2}$ th largest and  $W_{\alpha/2}$ th smallest of these averages. When the sample size is large, it is not necessary to count all the averages, compute only the averages near the largest and smallest.

This method is equivalent to a graphic procedure. Mark the heavy dots representing the sample data on the vertical axis as in Figure 3.  $X_{(n)}$  is the top point denoted by A,  $X_{(1)}$  is the bottom point denoted by B. A point half-way between A and B is

found and marked C. On the horizontal line segment through C mark point D at some convenient distance. The line segments AD, AB, and BD form an isosceles triangle. Through each dot on the vertical draw two lines parallel to the sides of the triangle and extending to them. Each intersection is marked with a heavy dot. The number of such dots is  $n(n+1)/2$ . Triangle ABD is an isosceles triangle, and each smaller triangle on the graph is an isosceles triangle similar to triangle ABD. Then the ordinate of each point of intersection is the average of pairs of original data points. So the n-th largest average is represented by the n-th dot from the top.

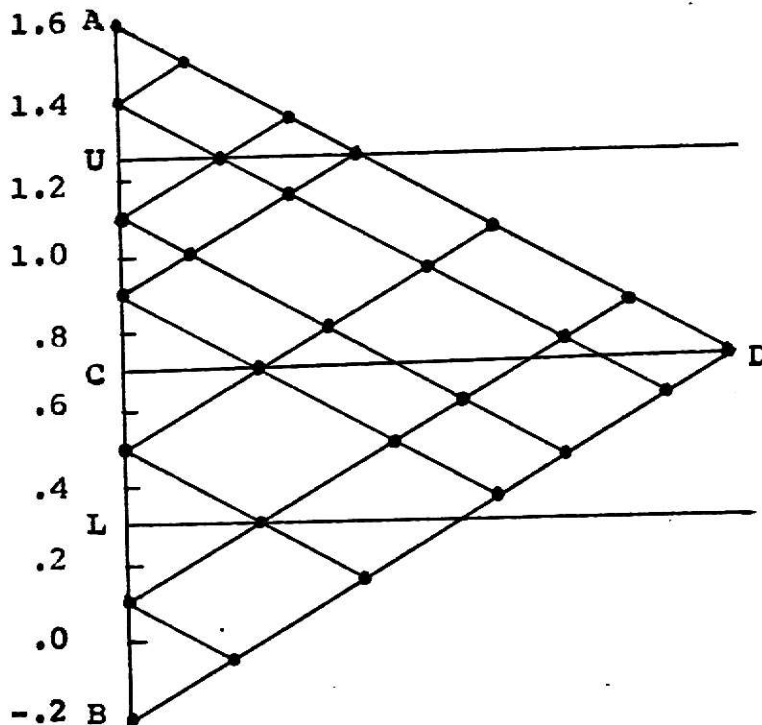


Figure 3



A confidence interval will be an interval on the vertical scale. To find it, read from Table IV the critical value  $W_{\alpha/2}$  for given  $n$ . As usual,  $1-\alpha$  is the confidence coefficient. Count down from the top dot A to the  $W_{\alpha/2}$ th dot, considering all dots in the figure. Draw a horizontal line through that dot until it intersects the vertical axis, then indicate this intersection as U (see Figure 3 where  $W_{\alpha/2}=5$ ). Now find the  $W_{\alpha/2}$ th of all dots from the bottom of the figure. Draw a horizontal line through it, and indicate the intersection of this line with the vertical axis as L. The points U and L are the endpoints of the confidence interval with coefficient  $1-\alpha$ . If two or more observations are equal, graph with a tiny distance between them. Say if  $r$  values are equal to the same value  $d$ , each dot on the line segment extend from  $d$  is counted  $r$  times, and dot  $d$  is counted  $r(r+1)/2$  times.

Example: Suppose that 10 samples of a particular type of wire are obtained and resistances in ohms are measured, and the following values obtained:

9.0, 15.0, 13.5, 7.5, 10.5, 9.5, 12.5, 10.5, 17.5, 11.5.

Find a 90% confidence interval for the median resistance.

For a 90% confidence interval, from Table IV,  $W_{\alpha/2}$  for  $n=10$  is found to equal 11. The algebraic method of finding the confidence interval is to find the eleventh largest and the eleventh smallest averages. The largest averages are:

17.50, 16.25, 15.50, 15.00, 15.00, 14.50, 14.25, 14.00,  
14.00, 13.75, 13.50, 13.50, 13.25, .....;



The eleventh dot from the top is at 13.5, and the eleventh dot from the bottom is at 9.75. Therefore, the 90% confidence interval on the median is from 9.75 to 13.50, the same answer as before.

Another method for obtaining a confidence interval on the median is to use the binomial distribution. Table III is entered with  $n=10$ . The probability 0.89 is selected as the nearest to 0.90 in Table III. The value of  $c$  associated with it is  $c=3$ .  $X_{(3)}$  is equal to 9.5,  $X^{(3)}$  is equal to 13.5. Therefore, The confidence interval found by this method is from 9.5 to 13.5, which agrees quite well with Wilcoxon method.

The third method is provided by John E. Walsh (1949). He proposed a significance test for the population median, which can also be used to find the confidence interval on the median. This method is efficient for small sample size, and can be applied with little computation. If the  $n$  sample observations are arranged in increasing order, the two-sided confidence interval has the form

$$(6) \quad P(\min(X_h, \frac{1}{2}(X_m + X_n)) \leq md \leq \max(X_r, \frac{1}{2}(X_s + X_t))) \geq 1 - \alpha$$

The one-sided confidence intervals are of the forms

$$(7) \quad P(md \geq \min(X_h, \frac{1}{2}(X_m + X_n))) \geq 1 - \alpha$$

and

$$(8) \quad P(md \leq \max(X_r, \frac{1}{2}(X_s + X_t))) \geq 1 - \alpha$$

where  $h, m, n, r, s, t$  are the values used in Table V with proper sample size and confidence level.

Example: An experiment is made to study the effect of certain hormones on the growth of cucumber seeds. Following are the lengths in centimeters of the roots developed from the seeds grown in the presence of a certain small concentration of the hormone, in increasing order of size:

7.2, 8.0, 8.1, 8.7, 11.5, 12.1, 12.8, 13.2, 15.2.

Find the 95% confidence interval for the median length of those roots.

Table V is entered with  $n=9$ , in the symmetrical column we find 4.3% is closest to 5%. Reading from that row we see that the median is expected to be between  $\max(X_7, \frac{1}{2}(X_5+X_9))$  and  $\min(X_3, \frac{1}{2}(X_1+X_5))$ . It is found that:

$$\min(X_3, \frac{1}{2}(X_1+X_5)) = \min(8.1, 9.35) = 8.1$$

$$\max(X_7, \frac{1}{2}(X_5+X_9)) = \max(12.8, 13.35) = 13.35$$

Therefore the approximate 95% confidence interval on the true median length of roots is between 8.1 and 13.35 cms.

## CONFIDENCE BAND FOR POPULATION DISTRIBUTION FUNCTION

Let  $X_1, X_2, \dots, X_n$  be mutually independent random variables with the common but unknown distribution function  $F(x)$ . Here the interest is to obtain a confidence band on  $F(x)$ . Let  $S(x)$  be the empirical cumulative probability distribution function for the  $n$  observations. Thus  $S(x)$  is a step function of  $x$  which shows the proportion of  $X_i$ 's whose values are less than or equal to  $x$  for arbitrary  $x$ . If  $k$  of the  $X_i$ 's are equal to or less than  $x$  then  $S(x)=k/n$ . Clearly, for any value of  $x$  from  $-\infty$  to  $+\infty$ ,  $S(x)$  must be an integer multiplied by  $1/n$ , and ranges from 0 to 1.

A confidence band can be based on the largest distance between  $S(x)$  and  $F(x)$  measured in a vertical direction. The statistic is in the form

$$(1) \quad d_\alpha = \max |S(x) - F(x)|$$

This is the statistic suggest by Kolmogorov (1933). The distance of  $d_\alpha$  (also the distribution of  $\max |S(x)-F(x)|$ ) is tabled as Table VI.

To form a confidence band graphically for  $F(x)$  with confidence coefficient  $1-\alpha$ , first draw a praph of the empirical distribution function  $S(x)$  based on the random sample. Find the critical value  $d$  of the Kolmogorov test statistic from Table VI for the two-sided test and for the appropriate sample size  $n$ . Then  $P(\max |S(x)-F(x)| \geq d_\alpha) = \alpha$  or contraiwise,  $P(\max |S(x)-F(x)| < d_\alpha) = 1-\alpha$ . This means that we can state with  $1-\alpha$  confidence that  $F(x)$  will

be within  $+d_\alpha$  of  $S(x)$ . Thus we may plot two step functions  $S(x)+d_\alpha$ , called  $U(x)$ , and  $S(x)-d_\alpha$ , called  $L(x)$ ; and expect that the unknown function  $F(x)$  lies entirely within the enclosing area with confidence level  $1-\alpha$ . The two step functions  $U(x)$  and  $L(x)$  together constitute a  $1-\alpha$  confidence band on  $F(x)$ . Because  $F(x)$  and  $S(x)$  both are cumulative functions,  $U(x)$  and  $L(x)$  must be bounded by 1 and 0 so the band obtained is limited to that range. The formal mathematical definitions of  $U(x)$  and  $L(x)$  are as follows:

$$\begin{array}{lll}
 (2) & U(x) = S(x) + d_\alpha & \text{if } S(x) + d_\alpha \leq 1 \\
 & U(x) = 1 & \text{if } S(x) + d_\alpha > 1 \\
 (3) & L(x) = S(x) - d_\alpha & \text{if } S(x) - d_\alpha \geq 0 \\
 & L(x) = 0 & \text{if } S(x) - d_\alpha < 0
 \end{array}$$

The Kolmogorov statistic for the discrete data was investigated by John Walsh (1963) who showed that the confidence coefficient for a band is exactly  $1-\alpha$  if  $F(x)$  is continuous, otherwise the confidence coefficient is at least as large as  $1-\alpha$ .

Let  $d_\alpha^+$  be the critical value for a one-sided test corresponding exactly to the significance level  $\alpha$ . Then by the same reasoning as above, we may state with confidence level  $1-\alpha$  that  $F(x)$  lies entirely below  $S(x)+d_\alpha^+$ .

Example: A random sample of size 10 is drawn and arranged in increasing order of magnitude as follows:

9.1, 10.3, 11.6, 12.4, 12.5, 13.0, 13.6, 14.2, 16.1  
and 18.7.

We wish to find a 90% confidence band for  $F(x)$ .

The empirical cumulative probability function  $S(x)$  is shown by the solid line in Figure 5. Entering Table VI with  $n=10$ ,  $\alpha=.10$ , we find that  $d_{.10}=.369$ . Then one dotted line is drawn parallel to the sample line in Figure 5 and 0.369 above it. A second dotted line is drawn parallel to the sample line and 0.369 below it. We may now assert with confidence measured by .90 that the cumulative distribution from which our sample was drawn lies inside the band bounded by these dotted lines.

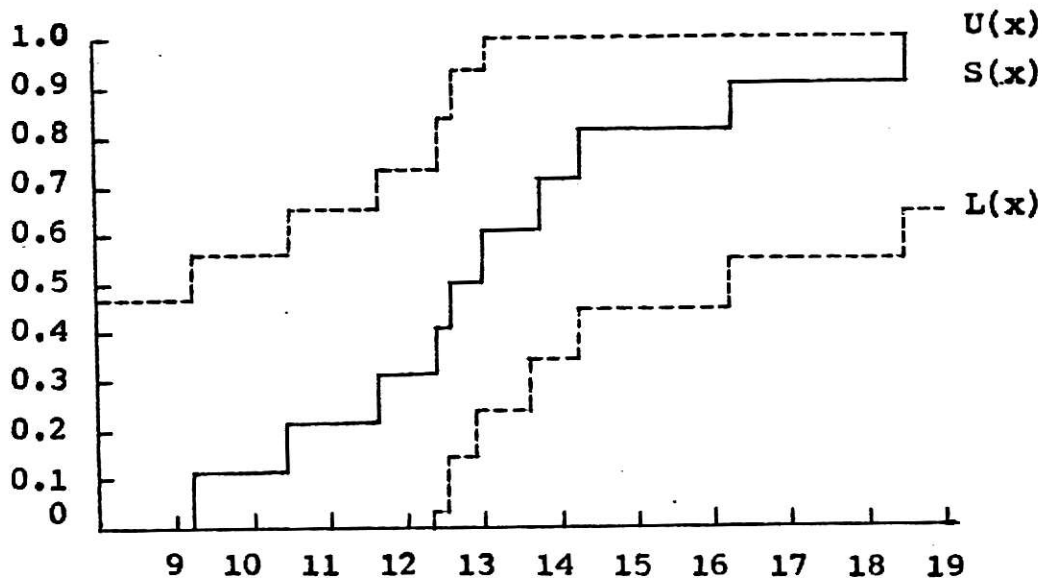


Figure 5

## TOLERANCE LIMITS

Suppose a random sample of size  $n$  is drawn from a population having a continuous cumulative distribution. Let the sample values arranged in order of increasing magnitude be  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ . Let the proportion of the population which is included between  $X_{(r)}$  (the  $r$ -th smallest value in the sample) and  $X_{(n-s+1)}$  (the  $s$ -th largest value) be greater or equal to  $p$ . Also, let  $1-\alpha$  be the probability that the interval  $(X_{(r)}, X_{(n-s+1)})$  will contain the proportion  $p$  or more of the population. Then the random interval  $(X_{(r)}, X_{(n-s+1)})$  has probability  $1-\alpha$  of containing at least 100p% of the population and is called the  $1-\alpha$  tolerance interval for 100p per cent of the population.  $X_{(r)}$  and  $X_{(n-s+1)}$  are called the  $1-\alpha$  tolerance limits for 100p per cent of the population. The quantity  $p$  is between 0 and 1 and has been called the population coverage. The restriction to the continuous distribution function has been removed by Scheffe and Tukey (1942); hence if the distribution is not continuous the statement of the tolerance limit must be "the probability is at least  $1-\alpha$  that at least the proportion  $p$  of the population is between  $X_{(r)}$  and  $X_{(n-s+1)}$ ".

The one-sided tolerance limit can be obtained by using the same theory as in finding the confidence interval of quartiles, already shown in previous sections. The one-sided tolerance limits are of the form:



$$(1) \quad P(\text{at least } p \text{ of the population} \leq X_{(n-s+1)}) \geq 1 - \alpha$$

or

$$(2) \quad P(X_{(r)} \leq \text{at least } p \text{ of the population}) \geq 1 - \alpha$$

The statement "at least  $p$  of the population  $\geq X_{(n-s+1)}$ " is saying each  $X$  has the probability  $p$  of being smaller than or equal to  $X_{(n-s+1)}$ . Therefore, we can write:

$$(3) \quad P(\text{at least } p \text{ of the population} \leq X_{(n-s+1)}) \\ = \sum_{a=0}^{n-s} \binom{n}{a} p^a (1-p)^{n-a}$$

Rewrite the right side of (3) as

$$(4) \quad \sum_{a=0}^{n-s} \binom{n}{a} p^a (1-p)^{n-a} = 1 - \sum_{a=n-s+1}^n \binom{n}{a} p^a (1-p)^{n-a}$$

because the sum of the binomial probabilities equals unity.

Again, rewrite the right side of equation (4) as

$$(5) \quad 1 - \sum_{a=n-s+1}^n \binom{n}{a} p^a (1-p)^{n-a} = 1 - \sum_{a=0}^{s-1} \binom{n}{a} (1-p)^a p^{n-a}$$

because the probability of  $n-s+1$  or more successes equals the probability of  $s-1$  or fewer failures. Combine equations (3), (4) and (5) we get

$$(6) \quad 1 - \sum_{a=0}^{s-1} \binom{n}{a} (1-p)^a p^{n-a} \geq 1 - \alpha$$

and also can be written as

$$(7) \quad \sum_{a=0}^{s-1} \binom{n}{a} (1-p)^a p^{n-a} \leq \alpha$$

The other one-sided tolerance limit is shown in Equation (2), and equals