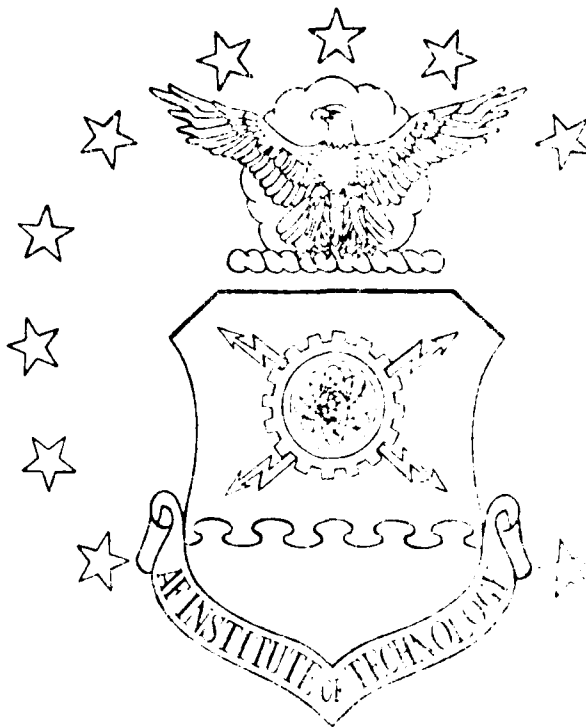


MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A

AD-A163 837



DTIC  
ELECTE  
FEB 10 1986  
S D

MODIFIED KOLMOGOROV-SMIRNOV,  
ANDERSON-DARLING, AND CRAMER-VON MISES  
TESTS FOR THE PARETO DISTRIBUTION WITH  
UNKNOWN LOCATION AND SCALE PARAMETERS

THESIS

James E. Porter III  
Captain, USAF

AFIT/GSO/MA/85D-6

DISTRIBUTION STATEMENT A

Approved for public release  
Distribution Unlimited

DEPARTMENT OF THE AIR FORCE  
AIR UNIVERSITY

**AIR FORCE INSTITUTE OF TECHNOLOGY**

Wright-Patterson Air Force Base, Ohio 45433

AFIT/GSO/MA/85D-6

①

DTIC  
ELECTE  
FEB 10 1986  
S D D

MODIFIED KOLMOGOROV-SMIRNOV,  
ANDERSON-DARLING, AND CRAMER-VON MISES  
TESTS FOR THE PARETO DISTRIBUTION WITH  
UNKNOWN LOCATION AND SCALE PARAMETERS

THESIS

James E. Porter III  
Captain, USAF

AFIT/GSO/MA/85D-6

Approved for public release; distribution unlimited

AFIT/GSO/MA/85D-6

MODIFIED KOLMOGOROV-SMIRNOV,  
ANDERSON-DARLING, AND CRAMER-VON MISES TESTS  
FOR THE PARETO DISTRIBUTION  
WITH UNKNOWN LOCATION AND SCALE PARAMETERS

THESIS

Presented to the Faculty of the School of Engineering  
of the Air Force Institute of Technology  
Air University  
in Partial Fulfillment of the  
Requirements for the Degree of  
Master of Science in Space Operations

James E. Porter III, B.S.  
Captain, USAF

December 1985



Accession For	
NTIS CRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

Approved for public release; distribution unlimited

## PREFACE

This thesis develops goodness-of-fit tests for the Pareto distribution by generating critical value tables for the modified Kolmogorov-Smirnov, Anderson-Darling, and Cramer-von Mises statistics. These tables can be used to test whether a set of observed values follows a Pareto distribution when the location and scale parameters are unspecified and must be estimated from the observed sample data. Additionally, the power of each of the three goodness-of-fit tests is studied and compared. Finally, the functional relationship between the critical values and the Pareto shape parameter is determined. Hopefully the material is presented in sufficient detail to be easily understood by those with only a passing knowledge of statistical analysis.

I wish to thank my reader and class advisor, Lieutenant Colonel Joseph Coleman, who guided me throughout my AFIT tour; and especially my thesis advisor, Dr. Albert H. Moore, who maintained my interest in statistical analysis, offered constant encouragement, and suggested the subject of this thesis. I also thank my classmates Majors Dennis Charek and Denny Danielson for their help in debugging the computer programs used in this thesis.

Above all I thank my family, especially my wife Judy, for their love and understanding during my tour at AFIT.

James E. Porter III

## TABLE OF CONTENTS

	Page
Preface . . . . .	ii
List of Figures . . . . .	vi
List of Tables . . . . .	vii
Abstract . . . . .	viii
I. Introduction . . . . .	1-1
Chapter Overview . . . . .	1-1
Background . . . . .	1-1
Problem Statement . . . . .	1-3
Research Question . . . . .	1-3
Objectives . . . . .	1-4
Presentation of Research . . . . .	1-4
II. Goodness-of-Fit Tests . . . . .	2-1
Chapter Overview . . . . .	2-1
Introduction . . . . .	2-1
Background . . . . .	2-2
Hypothesis Testing and Test Statistics . . . . .	2-4
Empirical Distribution Function . . . . .	2-6
Using Unknown Parameters . . . . .	2-9
Kolmogorov-Smirnov Statistic . . . . .	2-11
Cramer-von Mises Statistic . . . . .	2-12
Anderson-Darling Statistic . . . . .	2-13
Chapter Summary . . . . .	2-14
III. The Pareto Distribution . . . . .	3-1
Chapter Overview . . . . .	3-1
History and Application . . . . .	3-1
Origin . . . . .	3-1
Early Applications . . . . .	3-2
Recent Applications . . . . .	3-3
Air Force Applications . . . . .	3-4
The Pareto Function . . . . .	3-7
Parameter Estimation . . . . .	3-13
Various Estimators . . . . .	3-13
Best Linear Unbiased Estimator . . . . .	3-15
BLUEs for Shape $c > 2$ . . . . .	3-17
BLUEs for Shape $c \leq 2$ . . . . .	3-20
Summary of BLUEs . . . . .	3-24
Example 1 . . . . .	3-24

	Page
Modified Test Statistics . . . . .	3-27
Hypothesized Pareto CDF . . . . .	3-27
Example 2 . . . . .	3-28
Modified K-S Statistic . . . . .	3-29
Example 3 . . . . .	3-30
Modified A-D Statistic . . . . .	3-31
Example 4 . . . . .	3-31
Modified CV-M Statistic . . . . .	3-32
Example 5 . . . . .	3-32
Chapter Summary . . . . .	3-33
 IV. Methodology . . . . .	 4-1
Chapter Overview . . . . .	4-1
Basic Principles . . . . .	4-1
The Monte Carlo Method . . . . .	4-1
The Inverse Transform Technique . . . . .	4-4
Identifying Critical Values . . . . .	4-9
The Plotting Positions Technique . . . . .	4-11
Specific Procedures . . . . .	4-18
Stage 1: Generating Critical Value Tables . . . . .	4-18
Stage 2: Comparing Power . . . . .	4-21
Stage 3: Determining Relationship . . . . .	4-26
Chapter Summary . . . . .	4-27
 V. Results and Application . . . . .	 5-1
Chapter Overview . . . . .	5-1
Critical Value Tables . . . . .	5-1
Power Comparison Tables . . . . .	5-5
Regression Tables . . . . .	5-8
Use of Tables . . . . .	5-11
Using Critical Value Tables . . . . .	5-11
Using Power Comparison Tables . . . . .	5-13
Using Linear Regression Tables . . . . .	5-13
Chapter Summary . . . . .	5-15
 VI. Analysis and Discussion . . . . .	 6-1
Chapter Overview . . . . .	6-1
Critical Values . . . . .	6-1
Power Comparison . . . . .	6-2
Regression Analysis . . . . .	6-4
Verification and Validation . . . . .	6-5
Chapter Summary . . . . .	6-7
 VII. Conclusions and Recommendations . . . . .	 7-1
Conclusions . . . . .	7-1
Recommendations . . . . .	7-2

	Page
Appendix A: Computer Program for Critical Values . . .	A-1
Flow Chart . . . . .	A-2
Main Program CRITVAL . . . . .	A-4
Subroutine PARDEV . . . . .	A-12
Subroutine BXVALS . . . . .	A-14
Subroutine BLCLE2 . . . . .	A-17
Subroutine BLCGT2 . . . . .	A-19
Subroutine HYPCDF . . . . .	A-21
Subroutine TESTAT . . . . .	A-23
Subroutine CRTVAL . . . . .	A-26
 Appendix B: Computer Program for Power Comparison . .	 B-1
Flow Chart . . . . .	B-2
Main Program POWER . . . . .	B-4
Subroutine PARETO . . . . .	B-11
Subroutine BXVALS . . . . .	B-13
Subroutine BLCLE2 . . . . .	B-16
Subroutine BLCGT2 . . . . .	B-18
Subroutine HYPCDF . . . . .	B-20
Subroutine TESTAT . . . . .	B-22
Subroutine COMPAR . . . . .	B-27
 Bibliography . . . . .	 C-1
 Vita . . . . .	 D-1

## LIST OF FIGURES

Figure		Page
1	Three-Parameter Pareto Curves for Shape $c=2$ . . .	3-10
2	Two-Parameter Pareto Curves for Shape $c=2$ . . .	3-10
3	One-Parameter Pareto Curves for Several $c$ . . .	3-11
4	Probability Density of One-Parameter Pareto . . .	3-11
5	Finding Critical Values from Plotting Positions	4-17
6	Procedure for Generating Critical Values . . . .	A-2
7	Procedure for Determining Power Values . . . . .	B-2

LIST OF TABLES

Table	Page
I Calculation of BLUEs . . . . .	3-25
II Calculation of Hypothesized Pareto CDF . . . . .	3-29
III Calculation of Modified K-S Statistic . . . . .	3-30
IV Calculation of Modified A-D Statistic . . . . .	3-31
V Calculation of Modified C-VM Statistic . . . . .	3-32
VI Critical Values for the Modified K-S Test . . . . .	5-2
VII Critical Values for the Modified A-D Test . . . . .	5-3
VIII Critical Values for the Modified C-VM Test . . . . .	5-4
IX Power Test for $H_0$ : Pareto CDF ( $c = 1.0$ ) . . . . .	5-6
X Power Test for $H_0$ : Pareto CDF ( $c = 3.5$ ) . . . . .	5-7
XI K-S Critical Values vs. Pareto Shape Parameter. . . . .	5-9
XII C-VM Critical Values vs. Pareto Shape Parameter . . . . .	5-10

ABSTRACT

Modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) critical values are generated for the three-parameter Pareto distribution. The values may be used to test whether a set of observations follows a Pareto distribution when the location and scale parameters are unspecified and thus must be estimated from the sample. A Monte Carlo simulation of 5000 repetitions is used to generate critical values for sample sizes 5(5)30 (i.e., 5 to 30 in increments of 5) and Pareto shape parameters .5(.5)4.0.

A 5000-repetition Monte Carlo investigation is carried out by using 5, 15, and 25 observations from eight alternate distributions to compare the powers of the K-S, A-D, C-VM, and Chi-square tests. The power values of the tests are relatively low for a sample size of five. However, the powers of the modified K-S, A-D, and C-VM tests are considerably better than the Chi-square test at larger sample sizes. Next to the Chi-square test, the A-D test has the lowest power in most cases.

A functional relationship is identified between the modified K-S and C-VM test statistics and the Pareto shape parameter. The critical values are found to be a linear function of the shape parameters between 1.5 and 4.0.

MODIFIED KOLMOGOROV-SMIRNOV,  
ANDERSON-DARLING, AND CRAMER-VON MISES TESTS  
FOR THE PARETO DISTRIBUTION  
WITH UNKNOWN LOCATION AND SCALE PARAMETERS

I. INTRODUCTION

Chapter Overview

This chapter introduces the topic of goodness-of-fit testing and its applications. It states the problem, the research question, and the objectives of the research.

Background

Because the Air Force depends on highly complex weapons systems to perform its missions, factors such as the reliability and maintainability of equipment continue to receive a great deal of emphasis. Of particular importance to the Air Force is the ability to forecast time-to-failure of equipment components and expected maintenance service times.

In studying such phenomena, analysts often face the problem of testing agreement between probability theory and actual observations. When trying to develop a valid statistical model of observed data, the analyst performs four basic steps (5:332):

1. Collect and plot the raw data to develop a histogram (frequency distribution graph).

2. Hypothesize the underlying statistical distribution of the data by comparing the histogram to probability density functions of known distributions.

3. Use the observed data to estimate parameters that characterize the distribution.

4. Test the distributional assumption and parameter estimates for goodness-of-fit. If the hypothesis (that the data follow the assumed distribution) fails, return to step 2 (assume a different distribution) and repeat the process.

Goodness-of-fit tests measure the degree of agreement between the distribution of an observed data sample and a theoretical distribution. Three tests widely used for this purpose are the Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and the Cramer-von Mises (C-VM). Such tests have been developed for several well known distributions, including the normal, exponential, Weibull, gamma, uniform, Laplace, and others (9;19;34;35). However, there are many other distributions which have not been successfully examined for goodness-of-fit when the parameters of the distribution are unknown. One such distribution, which has significant potential for Air Force applications, is known as the Pareto distribution.

The Pareto distribution is an important function in statistical analysis, and several applications have been identified in the fields of economics and operations research. For example, the Pareto distribution has played a major role in investigations concerning the distributions of

city population sizes, natural resources, stock price fluctuations, and oil field locations (28:242). Other studies indicate that the Pareto can be used to model phenomena which may be applicable to Air Force interests, such as time-to-failure of equipment components (16), maintenance service times (22), nuclear fallout dispersion (18), and error clusters in communications circuits (7). Use of the Pareto for such practical applications would be enhanced by an accurate method to test goodness-of-fit of the Pareto distribution.

#### Problem Statement

A test to determine goodness-of-fit has not been developed for the Pareto distribution when the location and scale parameters are unknown. Such a test would be useful in determining whether a random sample of data taken from an observed phenomenon behaves as the Pareto distribution.

#### Research Question

How can the existing K-S, A-D, and C-VM tests be modified to produce new goodness-of-fit tests which can be applied to the Pareto distribution when the location and scale parameters are unknown?

## Objectives

The objectives of this thesis are to:

1. Generate and document the modified K-S, A-D, and C-VM critical value tables for the Pareto distribution. These tables can be used to test goodness-of-fit when parameters of the distribution are unknown.

2. Compare the powers of the modified K-S, A-D, and C-VM tests to determine which test can best detect a false Pareto distribution hypothesis. The power of a statistical test is the probability of correctly rejecting a false hypothesis.

3. Determine what (if any) functional relationship exists between the shape parameter and the critical values generated for the Pareto function. This relationship can then be used to interpolate critical values corresponding to parameters not found in the generated tables.

## Presentation of Research

The report on this thesis effort is presented in seven chapters. In this, the first chapter, the general topic of goodness-of-fit has been introduced and the problem, research question, and objectives have been stated.

Chapter II describes various types of goodness-of-fit tests; explains hypothesis testing and test statistics; and discusses the empirical distribution function.

Chapter III describes applications of the Pareto

distribution; presents its various forms; explores parameter estimation for the Pareto function; and develops the modified K-S, A-D, and C-VM test statistics for the Pareto.

Chapter IV describes the basic principles and specific procedures used to satisfy the research objectives.

Chapter V presents the results of the research effort, including tables of critical values, power comparisons, and regression coefficients.

Chapter VI further discusses the results of the research. Observations are made concerning the tables of critical values, power comparisons, and regression coefficients.

Chapter VII contains conclusions and recommendations based on the conduct and results of the research effort.

Finally, the flow charts and computer programs used to carry out the research are contained in the appendices.

## II. GOODNESS-OF-FIT TESTS

### Chapter Overview

This chapter briefly reviews the literature to provide a background for goodness-of-fit tests. It also describes hypothesis testing and test statistics as they relate to goodness-of-fit. Finally, it discusses the empirical distribution function and related statistics, including the exact and computational forms of the Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) test statistics.

### Introduction

Goodness-of-fit tests measure the degree of agreement between the distribution of an observed data sample and a theoretical statistical distribution (13:189). For example, a test for goodness-of-fit may involve examining a random sample from some unknown distribution to test the hypothesis that the underlying distribution is actually a known, specified function (13:345). If such tests indicate a close fit, the hypothesized distribution can then be applied in simulation modeling to predict failure and operational availability rates of Air Force systems and their components.

## Background

For years statisticians have attempted to find test statistics whose sampling distributions do not depend on certain parameter values or on the explicit form of the distribution of the population. Such tests are called non-parametric or distribution-free tests (39:68).

Two of the oldest and best known distribution-free tests for goodness-of-fit are the Chi-square and the Kolmogorov-Smirnov (K-S) tests (13:189;47:2). The Chi-square test compares frequencies of the observed data with expected frequencies of the hypothesized distribution. The test is flexible enough to allow some parameters to be estimated from the observed data, but it has some limitations. For example, it is restricted to large sample sizes (1:73). Also, it requires that the data be arbitrarily grouped, which may affect the results (13:357). The K-S test compares the cumulative distribution function (CDF) of the hypothesized distribution against the empirical distribution function (EDF) of the observed data sample. The K-S test can be used for large or small samples; however, it is restricted to distributions which are fully specified (i.e., there can be no unknown parameters that must be estimated from the sample) (13:357). The same limitation applies to two other related methods, the Anderson-Darling (A-D) and the Cramer-von Mises (C-VM) tests (19:204; 47:3-4).

In a significant development, David and Johnson (14) found that if a distribution has only a location and scale parameter, then the K-S and related goodness-of-fit tests are independent of the true parameter values when the parameters are replaced by invariant estimators. The estimators must be invariant in the sense that if each  $x$  is transformed by  $x \rightarrow ax+b$  then the estimate  $T=T(x)$  is similarly transformed by  $T \rightarrow aT+b$  (4:4). Therefore, critical values dependent only on sample size and significance level can be generated (54:5). This property also applies to a three-parameter CDF provided the shape parameter is treated as a constant. A more detailed explanation of this principle is included below in the section on "Using Unknown Parameters".

Based on this discovery by David and Johnson, critical value tables for the K-S and related tests have been modified to allow their use in several cases where parameters are estimated from observed data. In a modified test, the form of the test statistic itself remains essentially the same, except that estimates are used in place of exact parameters. However, the critical values for a modified test are considerably different. The critical value tables are no longer the same for all distributions. Instead, they are different for each different hypothesized distribution function. A modified test is still non-parametric or distribution-free because the level of significance is still independent of any untested assumptions regarding the

distribution of the underlying population. In fact, the form of the hypothesized distribution is the hypothesis being tested (13:357).

There are numerous cases for which modified tests have already been developed. For example, Lilliefors developed a modified K-S test for the normal (34) and exponential (35) distributions; Ream (43) developed another set of modified tests for the normal distribution; Woodruff, Moore, and Cortes (53) developed a modified K-S test for the three-parameter Weibull distribution; Bush (9) modified the A-D and C-VM tests to expand the goodness-of-fit tests for the Weibull distribution; Viviano (49) modified the K-S, A-D, and C-VM tests for the gamma distribution; and Yoder (54) developed a modified K-S, A-D, and C-VM test for the logistic distribution. The modified K-S, A-D, and C-VM tests have also been developed for the uniform, normal, Laplace, exponential, and Cauchy distributions (19). Using a different technique, Woodbury (52) too developed a set of modified tests for the uniform distribution.

#### Hypothesis Testing and Test Statistics

A fundamental concept in statistical testing is the hypothesis test. When studying a given phenomenon, it is often desirable to determine the distribution of the population being studied. In many cases, however, it is not practical to observe the entire population. Instead, a

relatively small sample of the population is usually selected, and observations are made from the small sample.

Hypothesis testing is the process of inferring from a sample whether to "accept" a certain statement (the null hypothesis) about the population from which the sample is drawn. Actually, "acceptance" of the null hypothesis does not imply that the null hypothesis is true, but that there is insufficient evidence from the data sample to reject the hypothesis. The null hypothesis, denoted  $H_0$ , is the hypothesis to be tested. The alternative hypothesis, denoted  $H_1$ , is equivalent to stating that  $H_0$  is not true (13:75-76).

Another key concept in statistical testing is the test statistic, a function of random variables which is used to help make the decision in a hypothesis test. In order to be useful for data analysis, the test statistic chosen should possess certain desirable properties. Most importantly, the statistic should assign real numbers to points in the sample so that the points are arranged in an order which reflects their ability to distinguish between a true  $H_0$  and a false  $H_0$  (13:77). For example, the test statistic normally assigns larger values to situations that indicate most strongly that  $H_0$  ought to be rejected, while smaller values of the test statistic usually indicate insufficient evidence to reject  $H_0$ . In this type of "one-tailed" test, if the value of the test statistic for a given set of data is greater than a certain "critical value", the analyst would reject  $H_0$ .

(13:77). The critical value is chosen so that when the null hypothesis  $H_0$  is true, the chance of erroneously rejecting  $H_0$  is some specified probability (e.g., .01 or .05) (2:193).

There are two types of errors that can be made in applying the decision criterion. The Type I error results in rejection of  $H_0$  when  $H_0$  is true. The Type II error results in acceptance of  $H_0$  when  $H_0$  is false. The probability of committing a Type I error, denoted by  $\alpha$ , is called the level of significance of the test. The probability of a Type II error is denoted  $\beta$ . The power of a statistical test, denoted  $1 - \beta$ , is the probability of correctly rejecting a false  $H_0$  (13:79).

#### Statistics Based on the Empirical Distribution Function

One class of test statistic used in goodness-of-fit testing compares an observed sample distribution function and an hypothesized theoretical distribution function. These statistics are based on the empirical distribution function (EDF), and in many cases are easily calculated and competitive in terms of power. The K-S, A-D, and C-VM test statistics are of the EDF type (45:730).

When analyzing phenomenon such as time-to-failure of equipment components,  $H(x)$ , the actual distribution function of the phenomenon, is rarely known. Often an educated guess of the form of the distribution is made, and the guess is used to approximate the true distribution function. One way

to make a "good guess" is to observe several values from random samples of the phenomenon and construct a graph that can be used to estimate the entire unknown distribution function  $H(x)$ . One widely used method of constructing such a graph is the empirical distribution function  $S(x)$ , which equals the fraction of observed values that are less than or equal to  $x$  (47:1), i.e.,

$$S(x) = \frac{\text{number of values } \leq x}{\text{total number of values}} \quad (1)$$

For a sample consisting of  $n$  observations, the EDF, which may be denoted  $S_n(x)$  to indicate the particular sample size, is a step-shaped function where each step is of height  $1/n$  and occurs only at the sample values. As  $n$  becomes larger,  $S_n(x)$  should better approximate  $H(x)$ , provided that  $H_0$  is true. When the  $n$  observations are arranged in ascending order, i.e., letting  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$  be the "order statistics" (15:4; 20:70), then  $S_n(x)$  is defined (47:1) as:

$$S_n(x) = \begin{cases} 0 & \text{for all } x < x_{(1)} \\ i/n & \text{for } x_{(i)} \leq x < x_{(i+1)}, i=1,2,\dots,n-1 \\ 1 & \text{for all } x > x_{(n)} \end{cases} \quad (2)$$

Like a CDF,  $S_n(x)$  is a nondecreasing function that ranges from zero to one in height; however,  $S_n(x)$  is determined empirically (from an observed sample), thus its name (13:70).

In a typical test for goodness-of-fit, a random sample from an unknown distribution is examined to test the null hypothesis that the unknown CDF  $H(x)$  is in fact a known, specified function  $F(x)$ , i.e.,  $H_0: H(x) = F(x)$ . The random sample is compared with the hypothesized distribution  $F(x)$  in some way to determine whether it is reasonable to conclude that  $F(x)$  is the true CDF of the random sample. Using the EDF  $S_n(x)$  is one way to compare the random sample with  $F(x)$ . The fact that  $S_n(x)$  is, by definition, the proportion of a random sample less than  $x$  implies that it should serve as a good estimate of  $F(x)$ , which is defined as the probability that the random variable  $X$  is less than the value  $x$  (47:1). Since the EDF  $S_n(x)$  may be useful as an estimator of the hypothesized CDF  $F(x)$ , then  $S_n(x)$  can be compared with  $F(x)$  to see if there is close agreement. If the level of agreement is poor, then the null hypothesis is rejected, i.e., the true but unknown CDF  $H(x)$  is not the same as the hypothesized function  $F(x)$  (13:345).

Based on this approach, the K-S, A-D, and C-VM tests use criteria that measure the discrepancy or "distance" between the hypothesized CDF  $F(x)$ , which approximates  $H(x)$  under  $H_0$ , and the EDF  $S_n(x)$ . The definitions of the three criteria relate to the full range of  $x$ , leading to integral forms of the A-D and C-VM test statistics. Conveniently, all three test statistics can be expressed in computational forms in terms of  $F$  and  $S_n$  at the observed  $x$  values (19:204).

Using Unknown Parameters. In their unmodified forms, most popular goodness-of-fit tests based on EDF statistics, including the K-S, C-VM, and A-D tests, are meant to be used only when the null-hypothesized distribution  $F(x)$  is fully specified (i.e., when all parameters are known). However, cases are rare in statistical practice when  $H_0$  is completely specified; thus, it is more realistic to have unknown parameters for the null distribution. When unknown parameters are involved, the K-S, C-VM, and A-D tests are no longer distribution-free, so that different critical values will relate to different  $F(x)$  in the null hypothesis (19:204). The reason for this is that the distributions of these and other EDF statistics depend on the sample size  $n$  and also on the values of the unknown parameters (47:4).

The K-S, C-VM, and A-D tests depend on the probability integral transformation described by David and Johnson (14). This transformation, when applied to a random sample from a distribution of specified parameters, produces ordered values from a uniform distribution over the interval from 0 to 1. These values are then used to calculate the EDF test statistic. As a result, the EDF statistic becomes a function of ordered uniform random variables. However, when parameters are unknown and must be estimated from the sample, the transformation fails to produce ordered uniform random variables (47:4). Unless appropriately modified, therefore, any EDF tests based on this transformation will generally be

restricted to cases where all parameters are specified.

An important exception occurs if the unknown parameters are location and scale only. David and Johnson (14) showed that if a distribution can be completely specified by a single parameter for location and a single parameter for scale, then goodness-of-fit tests based on the probability integral transformation are independent of the true parameter values when invariant estimators are used (38:384).

Fortunately, the Pareto distribution can be completely specified by a single location and a single scale parameter (28:239). The three-parameter form of the Pareto, presented in the next chapter, can be expressed in terms of a single location and scale parameter by treating the shape parameter as a known constant. Thus, the value of each EDF test statistic for the Pareto will depend only on the sample size and significance level, but not on the exact values of the unknown parameters (35:387). As a result, rather than having to produce a separate set of critical value tables for each set of location and scale parameters, only one set of tables is needed for each shape parameter and each sample size  $n$ . It is this principle, coupled with the fact that the Pareto possesses the necessary location and scale property, that allows the generation of valid critical value tables for the Pareto distribution (47:5).

To accomplish this goal, the existing (unmodified) K-S, A-D, and C-VM test statistics can be modified using an

invariant estimator; but first, the unmodified statistics are discussed in the following sections.

The Kolmogorov-Smirnov Statistic. The K-S statistic in its unmodified form is especially useful when sample sizes are small and when no parameters are estimated from the data. Often it is a more powerful test than the Chi-square for any sample size (34:399; 39:76). However, when parameter estimates must be made from the sample, the Chi-square test is easily modified by reducing the number of degrees of freedom, whereas the existing K-S critical values are overly conservative and must be modified using Monte Carlo techniques (5:357). In this context, the term "conservative" means that the critical values are too large so that the actual level of significance is smaller than the stated level of significance (13:90).

The K-S test statistic (36:259-260; 5:270; 19:204) is the largest (denoted "sup" for supremum) vertical distance between the completely specified hypothesized CDF  $F(x)$  and the observed EDF  $S_n(x)$ . Therefore, the test statistic is expressed as:

$$D = \sup_x |F(x) - S_n(x)| \quad (3)$$

which is equivalent to the computational form given by

$$D = \max (D^+, D^-) \quad (4)$$

$H_0$  is rejected if  $D$  exceeds a corresponding critical value (13:358).

If there are  $n$  observations,  $x_{(i)}$  is the  $i$ -th smallest observation, and  $z_i = F(x_{(i)})$  then (39:69):

$$D^+ = \sup_{1 \leq i \leq n} [(i/n) - z_i] \quad \text{and} \quad D^- = \sup_{1 \leq i \leq n} [z_i - (i-1)/n] \quad (5)$$

Thus the K-S statistic is the larger of these two values.

The Cramer-von Mises Statistic. Another way to measure the discrepancy between the hypothesized CDF  $F(x)$  and the observed EDF  $S_n(x)$  is to use statistics of the Cramer-von Mises family, based on the squared integral of the difference between the EDF and the distribution tested (47:2). One such statistic is the C-VM statistic itself (46:357):

$$W^2 = n \int_{-\infty}^{\infty} [S_n(x) - F(x)]^2 dF(x) \quad (6)$$

which in computational form is (3:766; 45:731):

$$W^2 = [1/(12n)] + \sum_{j=1}^n [z_j - (2j-1)/2n]^2 \quad (7)$$

where  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  are  $n$  ordered observations from the sample and  $z_j = F(x_{(j)})$  for  $j=1, 2, \dots, n$ .

The Anderson-Darling Statistic. Another member of the Cramer-von Mises family is the A-D statistic. To allow more flexibility in goodness-of-fit tests, Anderson and Darling (2:194) introduced the technique of incorporating a weight function into the K-S and C-VM test statistics. The result is still another method of testing the hypothesis that  $n$  observations have been drawn from a population with specified distribution function  $F(x)$ .

Anderson and Darling (3:767) suggested using a nonnegative weight function, here denoted  $\theta(u)$ , chosen by the analyst to accentuate the values of  $S_n(x) - F(x)$  in those areas where the test is desired to have greater sensitivity. This weight function serves to counteract the fact that the discrepancy between  $S_n(x)$  and  $F(x)$  becomes smaller in the tails, since each approaches 0 and 1 at the extremes (47:2). They found that choosing the weight function  $\theta$  in the form of  $\theta(u) = 1/[u(1-u)]$  has the effect of heavily weighting the discrepancy in the tails of the two distributions. The resulting A-D test statistic (2:193; 46:357) is:

$$A^2 = n \int_{-\infty}^{\infty} [S_n(x) - F(x)]^2 \theta[F(x)] dF(x) \quad (8)$$

$$\text{where } \theta[F(x)] = [F(x) \cdot (1-F(x))]^{-1}$$

Thus the C-VM statistic may be considered a special case of the A-D statistic where  $\theta[F(x)] = 1$ .

In computational form the A-D statistic is (3:765):

$$A^2 = -n - (1/n) \sum_{j=1}^n (2j-1) [\ln z_j + \ln(1-z_{n+1-j})] \quad (9)$$

where  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  are  $n$  ordered observations from the sample and  $z_j = F(x_{(j)})$  for  $j=1,2,\dots,n$ .

The A-D statistic is designed to be used when the analyst wants the test to have good power against alternatives in which  $F(x)$  and  $H(x)$ , the true distribution, disagree near the tails of  $F(x)$ , and is willing to sacrifice power against alternatives in which they disagree near the median of  $F(x)$  (3:767). Thus, the A-D statistic is used when the analyst wants to reject  $H_0$  if  $H(x)$  differs greatly from  $F(x)$ , and especially if the difference is in the tails.

### Chapter Summary

The K-S, A-D, and C-VM tests are non-parametric tests of goodness-of-fit which offer advantages over the older Chi-square test. In their usual forms, the K-S, A-D, and C-VM tests are restricted to distributions which are fully specified. However, when location and scale parameters are replaced by invariant estimators, the three tests can be modified to produce valid critical values for a given distribution. Hypothesis testing and test statistics are two statistical concepts which can be used to modify the existing tests for the Pareto distribution, which is discussed in detail in the next chapter.

### III. THE PARETO DISTRIBUTION

#### Chapter Overview

This chapter reviews the history and application of the Pareto Law; presents the Pareto distribution and its three parameters; explores parameter estimation for the Pareto function; and develops the modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) test statistics for the Pareto distribution.

#### History and Application

Origin. The Pareto distribution is an important function in statistical analysis. It is named after Vilfredo Pareto (1848-1923), a Swiss professor of economics who conducted the first extensive statistical study of the distribution of incomes. His analysis of nineteenth century income in various countries led to the development of his first law:

. . . if  $x$  signify [sic] a given income and  $N$  the number of persons with incomes exceeding  $x$ , and if a curve be drawn, of which the ordinates are logarithms of  $x$  and the abscissae logarithms of  $N$ , this curve, for all the countries examined, is approximately a straight line . . . This means that, if the number of incomes greater than  $x$  is equal to  $N$ , the number greater than  $mx$  is equal to  $N/m^{1.5}$ , whatever the value of  $m$  may be. Thus the scheme of income distribution is everywhere the same [42:647].

Therefore, "the logarithm of the percentage of units with an income greater than some value is a linear function of that value with negative slope, provided that this value is greater than an appropriate positive number" (32:6). This is known as the "strong" form of the Pareto Law, with functional form given by equation (11) below. The "weak" form of the law pertains to the asymptotic nature of a distribution's tail and implies that if  $\log [1-F_X(x)]$  is plotted against  $\log x$ , then the resulting curve should be asymptotic to a line with slope  $-c$  as  $x$  gets larger (32:6; 28:245).

Early Applications. Since the early days of its formulation, the Pareto Law and its related distribution functions have been examined primarily for potential applications in economics and operations research.

Based on his statistical observations, Pareto believed that any influence that causes an increase in the national income overall must also increase the income of the poor: "We cannot be confronted with any proposal the adoption of which would both make the dividend larger and the share of the poor smaller, or vice versa" (42:648). Pareto also believed his law to be universally inevitable, regardless of economic, social, and political conditions. Economists have since identified flaws (11:609; 17:171) in the Pareto Law to the extent that for several years the Pareto distribution became disreputable (28:233; 7:235) as an economic predictor:

The general defence of "Pareto's Law" as a law of even limited necessity rapidly crumbles. His statistics warrant no inference as to the effect on distribution of the introduction of any cause that is not already present . . . This consideration is really fatal; and Pareto is driven, in effect, to abandon the whole claim [42:654].

Nevertheless, more recent studies have shown the Pareto distribution can be very useful.

Recent Applications. Several more recent studies have revived interest in the Pareto distribution by demonstrating that it can be used to model or predict numerous empirical phenomena. For example, the Pareto distribution has played a major role in investigations concerning city population size, resources, stock price fluctuations, and oil fields (28:242). The Pareto has also been used to describe property values, inheritance, business mortality, worker migration, consumer prices, and effects of underreported income (32:7; 51).

Fisk (17:171, 174-175) showed that in some cases the Pareto distribution offers an improvement over the lognormal distribution, especially at the extremities (tails) of the distribution. Steindl (44:187-246) cited several examples of empirical economic data which follow the Pareto distribution, including the distribution of wealth, jobs by basic salary, the growth rate of firms and corporations, and several others. He also reaffirmed the Pareto Law's usefulness in economic theory:

Empirical laws are rare in economics, and the most obvious instance of such laws is the regular pattern of certain statistical distributions, such as the distribution of persons according to income or of business firms according to sales. A good many of these distributions conform to the so-called law of Pareto, i.e. the number of firms (for example) with sales in excess of X, plotted against X on logarithmic paper, is a straight line . . . The Pareto distribution is encountered in many fields and often the fit is very good [44:11].

Air Force Applications. Other studies have shown that the Pareto can be used to model phenomena which may be applicable to Air Force interests, such as time-to-failure of equipment, maintenance service times, nuclear fallout particles, and error clusters in communication circuits.

For example, Davis and Feldstein (16:299) showed the Pareto can be used to model survival data based on a population of items whose times-to-failure from a well defined origin are being observed. If each member of the population has a constant hazard rate based on a two-parameter gamma distribution, then the time-to-failure for the population is the Pareto type II of equation (13). Further, in some cases the Pareto competes with the Weibull distribution as a model for failure times of components. Like the Weibull, the generalized Pareto includes the exponential, and can therefore be used to test departures from the exponential (16:305-306).

Kaminsky and Nelson (30) showed how the Pareto distribution can be used in situations involving life testing,

reliability, and replacement policy. Specifically, they showed how to use the Pareto to predict the time of future failures from times of early failures in the same sample. They found, for example, that if items are put into service simultaneously, and it becomes necessary to begin replacing them when a certain percentage remain functional, then it is possible to predict the replacement time of future failures from the early failure times. In another example, "if  $n$  items form an  $n$ -component parallel system, then we can predict the time of system failure . . ." (30:145).

The Pareto distribution can also be of use in modeling queuing systems in which equipment maintenance service times are conditioned upon a random parameter. Harris (22:307) showed that if the conditional service distribution is exponential and the random parameter has a gamma density, then the resultant service times follow the Pareto distribution. Further, if a system consists of components which have exponentially distributed times-to-failure with a gamma parameter density, then the unconditional times to failure would follow the Pareto distribution (22:312). Harris also used the Pareto to develop a model which provides a means of obtaining measures of effectiveness of a large scale and complicated queuing process (22:308-309).

Freiling showed that the Pareto distribution, in the form of equation (10) with  $c = 3$ , can be used to model mass sizes of nuclear fallout particles (18:4). In addition, he

compared the usefulness of the Pareto and lognormal distributions in modeling the size distribution of particle mass in the fallout from land-surface bursts. For this specific application, Freiling found close similarities between the two distributions: "The agreement is such that if one curve is correct, the other will never be proved wrong . . . Thus it appears that the differences between the two approaches are trivial" (18:12). He concluded his study by noting that, in the case of nuclear airburst debris, the lognormal distribution has the advantage of having an "observationally confirmed theoretical basis." If the observational data is truncated, however, the Pareto distribution has the advantage of simplifying calculations of particle surface distribution.

In a study of error clusters in communication circuits, Berger and Mandelbrot (7:224) revealed still another application of the Pareto distribution. They proposed a new mathematical model to describe the distribution of the occurrence of errors in data transmission over telephone lines. They found that the statistics of communications errors can be described in terms of an error probability depending solely on the time elapsed since the last occurrence of an error. Further, they discovered that the distribution of inter-error intervals closely approximates the Pareto distribution of exponent less than one. As a result, the relative number of errors tend to zero as message lengths increase.

### The Pareto Function

Pareto's Law in its original form can be expressed as  $N = Ax^{-c}$  where  $A$  and  $c$  are parameters which characterize the function and  $N$  is the number of people having income of at least  $x$ . In a form more commonly used in statistical analysis, Pareto's Law becomes the Pareto distribution:

$$P(x) = \Pr[X \geq x] = (k/x)^c \quad \text{for } k, c > 0; x \geq k \quad (10)$$

where  $P(x)$  is the probability that the value of a random variable  $X$  (e.g., income) is at least  $x$ ,  $k$  is a lower bound on  $X$  (e.g., some minimum income), and  $c$  characterizes the shape of the graph of the distribution (28:233-234).

Accumulated probabilities over the range of values of  $x$  are given by the corresponding cumulative distribution function (CDF) of  $X$ , also known as the "Pareto distribution of the first kind" (28:234) or the "strong" Pareto law (32:50):

$$F_X(x) = 1 - (k/x)^c \quad \text{for } k, c > 0; x \geq k \quad (11)$$

The corresponding Pareto probability density function is:

$$p_X(x) = ck^c/x^{c+1} = (c/k)(k/x)^{c+1} \quad \text{for } c > 0; x \geq k > 0 \quad (12)$$

Pareto proposed two other forms of the distribution. The "Pareto distribution of the second kind" (also called the Pareto Type II or the Lomax distribution), is:

$$F_X(x) = 1 - K_1 / [(x+C)^C] \quad (13)$$

The third form proposed by Pareto, the "Pareto distribution of the third kind" (or Pareto Type III), has the CDF:

$$F_X(x) = 1 - k_2 e^{-bx} / [(x+C)^C] \quad (14)$$

which reduces to the Type II form when  $b = 0$ .

The basic difference between these various forms is in the number of parameters. The Pareto distribution of the first kind, equation (11), represents the "usual formulation" of the function and is the one most commonly found in the literature. However, the fact that it consists of only two parameters (i.e.,  $c$  and  $k$ ) may limit its usefulness in general applications. Hastings and Peacock (26) regard three types of parameters as basic to any distribution function. These three parameters are the location, scale, and shape parameters, which they denote as  $a$ ,  $b$ , and  $c$  respectively. The location parameter ( $a$ ) represents "the abscissa of a location point (usually the lower or midpoint) of the range of the variate." The scale parameter ( $b$ ) is "a parameter which determines the scale of measurement of the fractile  $x$ ".

Finally, the shape parameter (c) "determines the shape . . . of the distribution function within a family of shapes associated with a specified type of variate" (26:20).

Kulldorff and Vannan (33:218) introduced a more general form of the CDF than the two-parameter form shown in equation (11). By using the parameter notation of Hastings and Peacock, and the functional form of Kulldorff and Vannan, the generalized (three-parameter) form of the Pareto distribution is illustrated in Figure 1 and can be written as:

$$F(x) = 1 - [1 + (x-a)/b]^{-c} \text{ for } x \geq a; \quad b, c > 0 \quad (15)$$

where again a is location, b is scale, and c is shape.

In the special case when  $a = b$ , if we let  $k = a = b$  as in Figure 2, then from equation (15):

$$\begin{aligned} F(x) &= 1 - [1 + (x-a)/b]^{-c} = 1 - [1 + (x-k)/k]^{-c} \\ &= 1 - [1 + (x/k) - (k/k)]^{-c} = 1 - (1 + x/k - 1)^{-c} \\ &= 1 - (x/k)^{-c} = 1 - (k/x)^c \end{aligned}$$

where  $k, b, c > 0$  and  $x \geq k = a$ . The last expression is the "usual formulation" given by equation (11).

Another form commonly found in the literature (26:102; 51:1) is the one-parameter form (Figures 3 and 4) given by:

$$F(x) = 1 - x^{-c} \text{ for } x \geq 1; \quad c > 0 \quad (16)$$

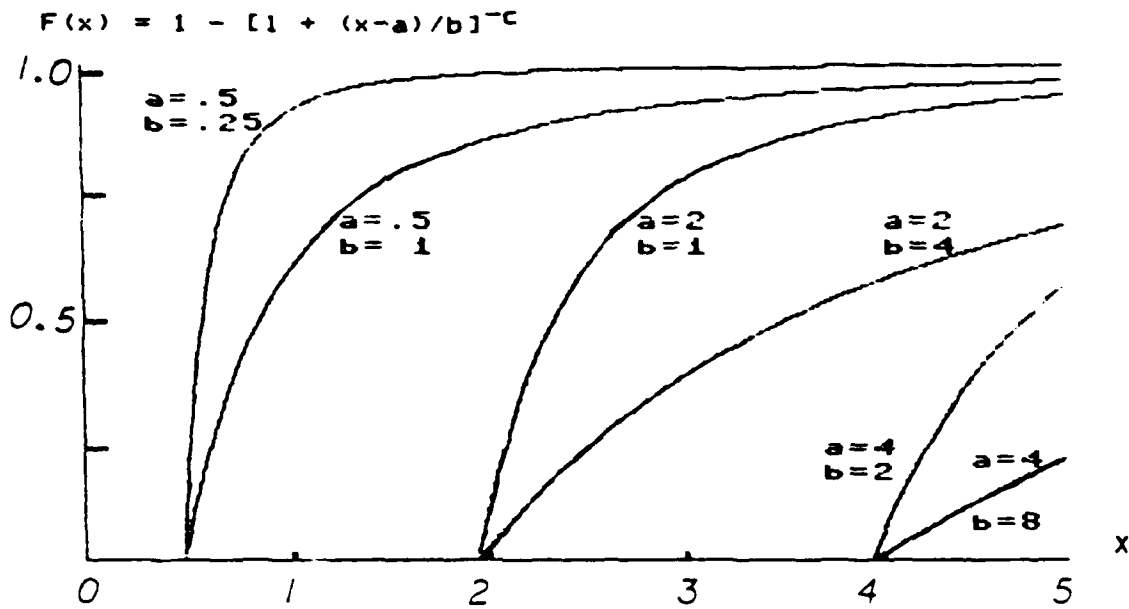


Fig 1. Three-Parameter Pareto Curves (Eqn 15) for Several Values of Location  $a$  and Scale  $b$  with Shape  $c = 2$ .

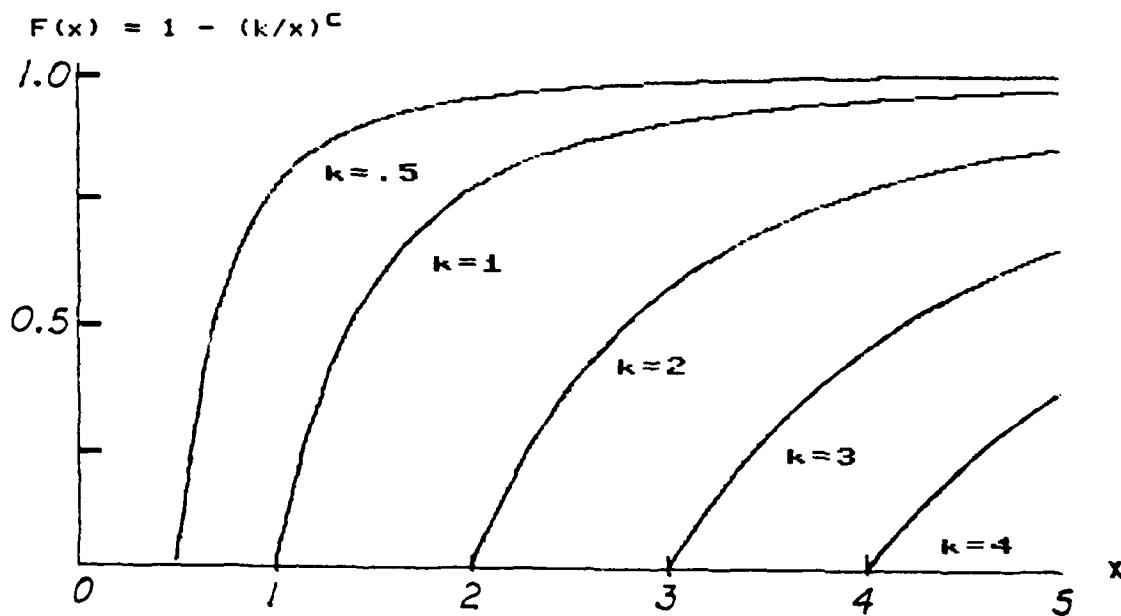


Fig 2. Two-Parameter Pareto Curves (Eqn 11) for Several Values of  $k$  with Shape  $c = 2$  and  $a = b = k$ .

$$F(x) = 1 - x^{-c}$$

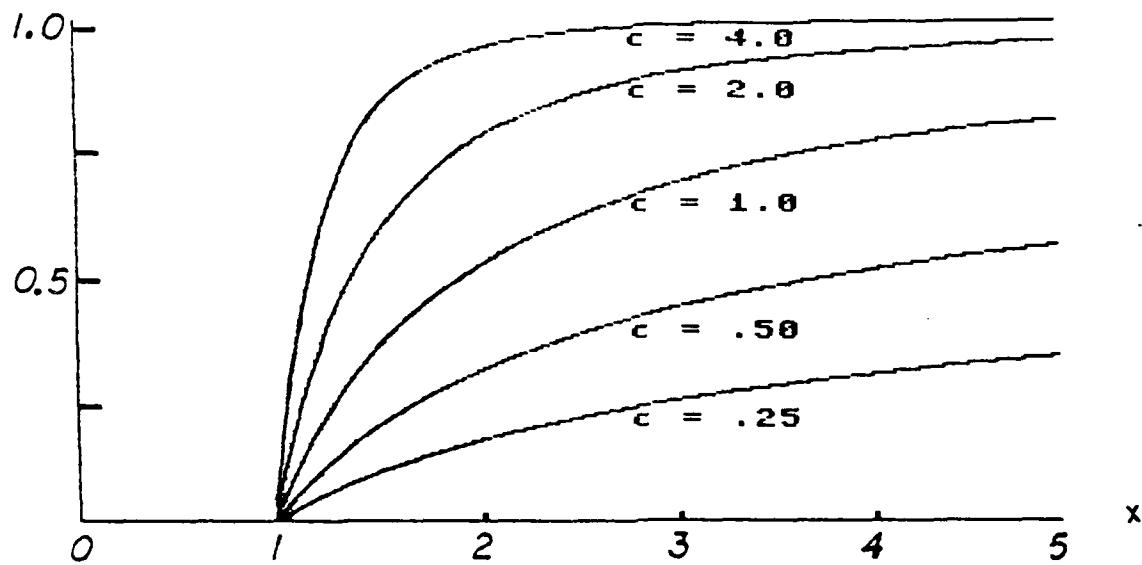


Fig 3. One-Parameter Pareto Curves (Eqn 16) for Several Values of Shape  $c$  with  $k = a = b = 1$ .

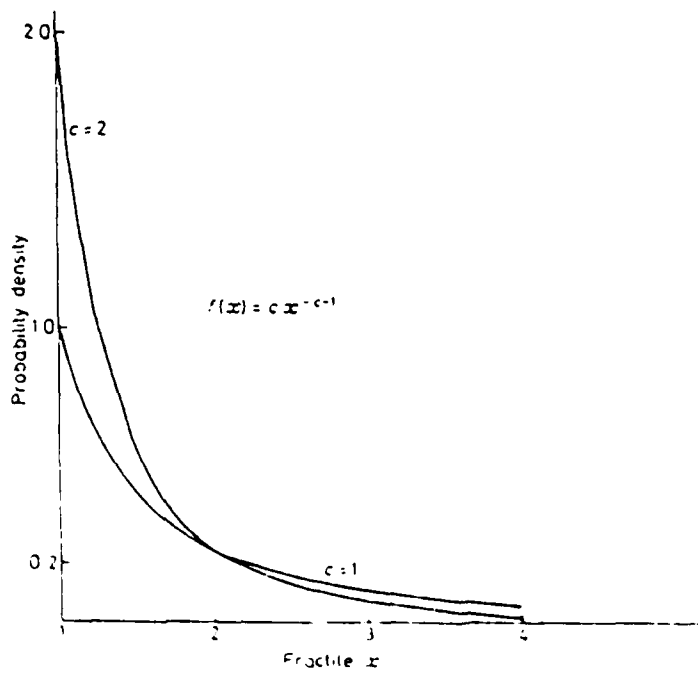


Fig 4. Probability Density (Eqn 12) of the One-Parameter Pareto with  $k = 1$  (Reprinted from 26:103).

Equation (16) is simply a special case of (15) found by setting  $a = b = 1$ . As such, it represents the least general form of the Pareto distribution.

The greater generality inherent in the three-parameter form, equation (15), allows the Pareto distribution to be more useful in practical applications. For example, in some situations the random variable represented by  $x$  may be positive by its very nature, making the assumption  $a = 0$  more realistic than  $a = b$  (33:218). In the special case where  $a = 0$ , the three-parameter Pareto distribution becomes:

$$\begin{aligned} F(x) &= 1 - [1 + (x-a)/b]^{-C} = 1 - (1 + x/b)^{-C} \\ &= 1 - (b/b + x/b)^{-C} = 1 - [(x+b)/b]^{-C} \\ &= 1 - [b/(x+b)]^C = 1 - b^C / [(x+b)^C] \end{aligned}$$

This last expression can be written as equation (13) by simply setting  $b^C = K_1$  and  $b = C$ .

Therefore, equations (11), (13), and (16) each represent special cases of the three-parameter form given by equation (15). Since (15) is a more general and hence more useful form of the Pareto distribution, this thesis uses the functional form in (15) to develop the goodness-of-fit tests for the Pareto distribution. Selecting the more general form as a basis for the test statistics will ensure the widest possible application of the goodness-of-fit tests.

## Parameter Estimation

As explained in Chapter II, the development of modified Kolmogorov-Smirnov, Anderson-Darling, and Cramer-von Mises tests depends on the use of an invariant estimator for the unspecified location and scale parameters (38:384). This section begins by briefly examining several published studies on various estimation techniques for Pareto distributions. It concludes by discussing the best linear unbiased estimator (BLUE), which is the invariant estimator used in this thesis.

Various Estimators. The two methods of invariant estimation most commonly used in modified goodness-of-fit tests are the maximum likelihood estimator (MLE) and the best linear unbiased estimator (BLUE). Various techniques for estimating the parameters of the Pareto distribution can be found in the literature. However, as Kulldorff and Vannman (33:218) point out, few studies consider the general three-parameter form of equation (15). Instead, most studies consider only "special cases", such as  $a = b$ , corresponding to equations (11) and (12).

Numerous examples of "special case" estimators can be cited. Moore and Harter (41; 23:69,86) developed a biased, single-order-statistic MLE for the Pareto shape parameter when location is specified. Harris (22:308, 310-311) considered estimation for the two-parameter form given by

equation (12): "As a first try, we can appeal to the techniques of maximum likelihood estimation. However, this particular method does not yield sufficiently simple equations (for even numerical methods)" (22:310). As a result, Harris resorted to the method of moments instead. Johnson and Kotz (28:234-240) presented MLEs for the two-parameter form in equation (11), as well as several other estimation techniques. Hastings and Peacock (26:102) gave the MLE for the one-parameter form of equation (16). In his dissertation, Koutrouvelis (32:97-115) attempted to estimate the parameters of the upper tail of Pareto distributions, but found it too difficult to calculate the Pareto MLEs, even with a computer. Instead, he developed a new method of estimating parameters based on the asymptotic theory of quantiles using only data consisting of sample values greater than some specified value. Wingo (50) wrote a FORTRAN program to calculate the MLEs from a reduced log-likelihood function for the two-parameter form in equation (12). Davis and Feldstein (16:299-300, 305) developed MLEs from progressively censored data for the Pareto Type III, equation (14). Bell, Ahmad, Park, and Lui (6:4-7) presented the MLEs, the minimum variance unbiased estimators (MVUEs), and the minimal sufficient statistic (MSS) for the two-parameter form, equation (11). Several other estimation studies are cited by Koutrouvelis (32:55) and Johnson and Kotz (28:235-240). Unfortunately, none of these studies provide

the invariant estimators of the three-parameter form in equation (15) as needed for this thesis.

Parameter estimation for the general case given by equation (15) went virtually ignored until Kulldorff and Vannman (33) derived the BLUEs of the unknown parameters on the basis of a complete Pareto sample with shape  $c > 2$ . In a follow-up paper, Vannman (48) derived the BLUEs for shape  $c \leq 2$ . Later, Kaminsky (29:7-8, 12-14) and Kaminsky and Nelson (30:148) extended the work of Kulldorff and Vannman by deriving, for equation (15), the best linear unbiased predictors of future observations from censored samples. Most recently, Charek (12) examined minimum distance estimation for the three-parameter Pareto.

Best Linear Unbiased Estimator (BLUE). The BLUE derives its name from its main properties as an estimator. It is a "linear" estimator because it can be expressed as a linear function of a random sample. It is "unbiased" because its bias term is zero; and the expected value of the estimator is equal to the true parameter value. It is considered the "best" estimator because it has the minimum variance among all other linear unbiased estimators (27:227). However, for the purposes of this thesis, the most important property of the BLUE is invariance under transformation of parameters.

The BLUE is a subset of a larger class of estimators

known as least-squares estimators. In general, least squares estimators do not possess the invariance property. However, when a least-squares estimator is also a linear function, then the invariance property holds (40:349-350). Therefore, in addition to its other properties, the BLUE is also an invariant estimator. It is this property of invariance under parameter transformations that allowed, for example, Green and Hegazy (19:205) and Woodbury (52) to use the BLUE in producing modified goodness-of-fit tests based on the findings of David and Johnson (14).

Intuitively, the property of invariance implies, for example, that if a parameter  $\theta$  is estimated, and  $\theta^2$  is also estimated from the same data, then the estimate of  $\theta^2$  should be the square of the estimate of  $\theta$  (37:434). Generally, the invariance property requires that if  $f(\theta)$  is a single valued function of a parameter  $\theta$ , and  $\hat{\theta}$  is the BLUE of  $\theta$ , then  $f(\hat{\theta})$  is the BLUE of  $f(\theta)$ , i.e.,  $f(\hat{\theta}) = \hat{f(\theta)}$  (8:94).

The studies by Kulldorff and Vannman (33; 48) derived the BLUEs of equation (15) for  $b$  when  $a$  and  $c$  are known; for  $a$  when  $b$  and  $c$  are known; and for  $a$  and  $b$  when  $c$  is known. The last case, which corresponds to invariant estimation of location and scale when shape is known, is used in this thesis to develop the modified K-S, A-D, and C-VM tests. The next two subsections use the findings of Kulldorff and Vannman to derive computational forms of the BLUEs for the Pareto location and scale parameters, assuming shape is known.

BLUEs for Shape  $c > 2$ . For the case where  $c > 2$ , Kulldorff and Vannman (33:224-226) found that the BLUEs of location  $a$  and scale  $b$  can be written in terms of the specified shape parameter  $c$  and the order statistics (15:4)  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  where  $x_{(1)}$  is the smallest and  $x_{(n)}$  the largest value in the observed random sample of size  $n$ . Thus the BLUEs for  $a$  and  $b$  are, respectively:

$$\hat{a} = x_{(1)} - Y / [(nc-1)(nc-2) - ncD] \quad (17)$$

$$\begin{aligned} \hat{b} &= Y(nc-1) / [(nc-1)(nc-2) - ncD] \\ &= [x_{(1)} - \hat{a}] (nc-1) \end{aligned} \quad (18)$$

In the special case when it is known that  $a = b$ , as in equation (11), the BLUE reduces to:

$$\hat{k} = [1 - 1/(nc)] x_{(1)} \quad (19)$$

However, before equations (17) and (18) can be used to find the BLUEs for the general case, the following terms must first be calculated:

$$B_i = \frac{\Gamma(n-i+1) \Gamma(n+1-2/c)}{\Gamma(n-i+1-2/c) \Gamma(n+1)} \quad \text{for } i = 1, 2, \dots, n \quad (20)$$

$$D = (c+1) \sum_{i=1}^{n-1} B_i + (c-1) B_n \quad (21)$$

$$Y = (c+1) \sum_{i=1}^{n-1} B_i x(i) + (c-1)B_n x(n) - Dx(1) \quad (22)$$

After these values are calculated, they can be substituted into equations (17) and (18) to find the BLUEs  $\hat{a}$  and  $\hat{b}$ .

From equations (17) to (22), it is obvious that the use of the BLUEs  $\hat{a}$  and  $\hat{b}$  involves the computation of all the coefficients  $B_i$  for  $i = 1, 2, \dots, n$ . Therefore, in order to derive a computational form of the BLUEs, the first task is to simplify equation (20). Each  $B_i$  is the ratio of a product of gamma functions. Banks and Carson (5) note that "the gamma function can be thought of as a generalization of the factorial notion which applies to all positive numbers, not just integers" (5:144). For any real  $m > 0$ :

$$\Gamma(m) = (m-1) \Gamma(m-1) \quad (23)$$

By definition  $\Gamma(1) = 1$ , so that whenever  $m$  is an integer, equation (23) becomes:

$$\Gamma(m) = (m-1)! \quad (24)$$

Applying these gamma definitions in equation (20) reveals:

$$\begin{aligned}
B_1 &= \frac{\Gamma(n-1+1) \Gamma(n+1-2/c)}{\Gamma(n-1+1-2/c) \Gamma(n+1)} = \frac{\Gamma(n) \Gamma(n+1-2/c)}{\Gamma(n-2/c) \Gamma(n+1)} \\
&= \frac{(n-1)! (n-2/c) \Gamma(n-2/c)}{n(n-1)! \Gamma(n-2/c)} = \frac{n-2/c}{n} \\
&= 1 - 2/(cn) \tag{25}
\end{aligned}$$

Similarly,  $B_2$  is found from equation (20) as follows:

$$\begin{aligned}
B_2 &= \frac{\Gamma(n-2+1) \Gamma(n+1-2/c)}{\Gamma(n-2+1-2/c) \Gamma(n+1)} = \frac{\Gamma(n-1) \Gamma(n+1-2/c)}{\Gamma(n-1-2/c) \Gamma(n+1)} \\
&= \frac{(n-2)! (n-2/c) \Gamma(n-2/c)}{\Gamma(n-1-2/c) n!} \\
&= \frac{(n-2)! (n-2/c) (n-1-2/c) \Gamma(n-1-2/c)}{n(n-1)(n-2)! \Gamma(n-1-2/c)} \\
&= \frac{(n-2/c) (n-1-2/c)}{n(n-1)} \\
&= [1 - 2/(cn)] [1 - 2/c(n-1)] \tag{26}
\end{aligned}$$

Continuing in this manner, it turns out that:

$$B_n = [1 - 2/(cn)] [1 - 2/c(n-1)] \cdots [1 - 2/c(1)] \tag{27}$$

The calculations can be simplified as follows:

Let  $g_1 = 2/(cn)$ ,  $g_2 = 2/[c(n-1)]$ ,  $\cdots$ ,  $g_n = 2/c$ .

Also let  $b_1 = 1-g_1$ ,  $b_2 = 1-g_2$ ,  $\cdots$ ,  $b_n = 1-g_n$ .

Then  $B_1 = b_1$ ,  $B_2 = b_1 b_2$ , ...,  $B_n = b_1 b_2 \dots b_n$ .

In general, then, each  $B_i$  can be expressed in computational form as:

$$B_i = \prod_{j=1}^i b_j \quad (28)$$

where  $b_j = 1 - g_j$  and  $g_j = 2/c(n-j+1)$  for  $j = 1, 2, \dots, i$ .

From these results, if we let  $B_0 = 1$ , then another way to write  $B_i$  is (48:705):

$$B_i = [1 - 2/c(n-i+1)] B_{i-1} \quad \text{for } i = 1, 2, \dots, n \quad (29)$$

As mentioned earlier, once all of the  $B_i$  are computed from equation (28) or (29), then  $D$  and  $Y$  can be computed from equations (21) and (22). Finally, these values for  $B_i$ ,  $D$ , and  $Y$  are substituted into equations (17) and (18) to find the BLUEs  $\hat{a}$  and  $\hat{b}$ .

BLUEs for Shape  $c \leq 2$ . For the case where  $c \leq 2$ , the variance of the Pareto distribution does not exist, so a different approach must be used to derive the BLUEs. In this case, Vannman (48:706-707) found that the BLUEs of location  $a$  and scale  $b$  can still be found provided that shape  $c$  satisfies  $2/n < c \leq 2$ , where again  $n$  is the sample size. Here the BLUEs  $a_k^*$  and  $b_k^*$  are based on the first  $k$  order statistics only, where  $k$  is chosen so that  $2 \leq k < n+1-2/c$ :

$$a_k^* = x_{(1)} - b_k^*/(nc-1) \quad (30)$$

and

$$\begin{aligned} b_k^* = (1/U_k) \{ & (c+1) \sum_{i=1}^{k-1} B_i x_{(i)} \\ & + [(n-k+1)c - 1] B_k x_{(k)} \\ & - [(nc-1)/(nc)] (nc-2-U_k) x_{(1)} \} \end{aligned} \quad (31)$$

where

$$U_k = \frac{(nc-2)(nc-c-2) - nc[(n-k)c - 2] B_k}{(nc-1)(c+2)} \quad (32)$$

Whenever possible,  $k$  should be chosen to achieve highest efficiency, which occurs when  $k = n - [2/c]$ , where " $[2/c]$ " denotes the integer portion of  $2/c$ . Vannman (48:707) also points out that in the case where  $2/c$  is an integer, and  $k$  is selected for highest efficiency so that  $k = n - 2/c$ , then equation (31) can be simplified to:

$$b_k^* = \frac{(c+1)(c+2)(nc-1)}{(nc-2)(nc-c-2)} \left[ \sum_{i=1}^{n-2/c} B_i x_{(i)} - \frac{nc-2}{c+2} x_{(1)} \right] \quad (33)$$

By substituting this result for  $b_k^*$  in equation (30), the BLUE for  $a$ , based on the first  $n-2/c$  order statistics, can be written in the following computational form:

$$a_k^* = x_{(1)} - \frac{(c+1)(c+2)}{(nc-2)(nc-c-2)} \left[ \sum_{i=1}^{n-2/c} B_i x_{(i)} - \frac{nc-2}{c+2} x_{(1)} \right] \quad (34)$$

Once  $a_k^*$  has been computed, it is easy to use equation (30) to find a computational form of the BLUE for  $b$ :

$$b_k^* = b_{n-2/c}^* = (nc-1) (x_{(1)} - a_k^*) \quad (35)$$

Equations (34) and (35) give the BLUEs for location  $a$  and scale  $b$  provided all of the following conditions apply:

- 1) shape parameter  $c$  is specified
- 2)  $2/n < c \leq 2$
- 3)  $2/c$  is an integer

When sample size  $n = 5, 10, 15, 20, 25,$  or  $30$ , then all three of these conditions hold for shape parameter  $c = .5, 1,$  or  $2$ . Therefore, for these values of  $n$  and  $c$ , it appears that equations (34) and (35) apply. There is, however, one important exception. As explained earlier,  $k$  must be chosen so that  $2 \leq k < n+1-2/c$ . In the case where  $n = 5$  and  $c = .5$ , notice that  $n+1-2/c = 2$ . Thus  $k$  cannot be selected as before, since it would need to satisfy  $2 \leq k < 2$ , which is not possible. As a result, the above equations fail to provide BLUEs for the special case  $c = .5$  and  $n = 5$ ; thus, when  $c = .5$ , this thesis will use  $n = 6$  instead of  $n = 5$ .

As explained in the next chapter, this thesis uses sample sizes of  $n = 5, 10, 15, 20, 25,$  and  $30$ , with shape

parameters of  $c = .5, 1, 1.5, 2, 2.5, 3, 3.5,$  and  $4$ . The preceding subsection presented the BLUEs to be used for  $c = 2.5, 3, 3.5,$  and  $4$ . This subsection has thus far shown that equations (34) and (35) provide the BLUEs for  $c = .5, 1,$  and  $2$ , except for the special case  $c = .5$  and  $n = 5$ . The one remaining case to be addressed is when  $c = 1.5$ .

When the shape parameter  $c = 1.5$ , equations (34) and (35) do not apply since condition 3) fails to hold, i.e.,  $2/c$  is not an integer. To ensure highest efficiency,  $k$  is selected so that  $k = n - [2/c]$ , where "[2/c]" denotes the integer portion of  $2/c$ . Thus:

$$k = n - [2/c] = n - [1.333] = n - 1 \quad (36)$$

According to Vannman (48:707), substituting this value of  $k$  into equations (30) to (32) gives the desired BLUEs:

$$a_k^* = a_{n-1}^* = x_{(1)} - b_{n-1}^*/(nc-1) \quad (37)$$

$$b_k^* = (1/U_{n-1}) \left\{ (c+1) \sum_{i=1}^{n-2} B_i x_{(i)} + (2c-1) B_{n-1} x_{(n-1)} - [(nc-1)/(nc)] (nc-2-U_{n-1}) x_{(1)} \right\} \quad (38)$$

where

$$U_k = U_{n-1} = \frac{(nc-2)(nc-c-2) - nc(c-2) B_{n-1}}{(nc-1)(c+2)} \quad (39)$$

Summary of BLUEs. For shape parameter  $c = .5, 1, \text{ or } 2$ , this thesis uses equations (34) and (35) to calculate the BLUEs for location parameter  $a$  and scale parameter  $b$ ; however, the case  $c = .5$  and  $n = 5$  is omitted, since then the BLUEs cannot be found. When  $c = 1.5$ , the BLUEs are given by equations (37) to (39). For  $c = 2.5, 3, 3.5, \text{ or } 4$ , equations (17), (18), (21), (22) and (29) are used to calculate the BLUEs for  $a$  and  $b$ . Once the BLUEs have been computed, the K-S, A-D and C-VM test statistics can be modified to accommodate unspecified location and scale parameters. An example will help to illustrate the calculations involved.

Example 1. In Table I the data listed under the  $x_i$  column was generated from a Pareto distribution of shape parameter  $c = 2.5$ , using equation (47) in the next chapter. Suppose it is desired to find the BLUE estimators  $\hat{a}$  and  $\hat{b}$  based on this particular random sample of size  $n = 10$ . Since in this case it is known that  $c = 2.5$ , the BLUEs will be computed from equations (17) and (18). One procedure to accomplish this is as follows:

Step 1. Arrange the  $x_i$  sample values in order from smallest to largest. The resulting order statistics (20:70) are listed under the  $x_{(i)}$  column of Table I.

Table I  
CALCULATION OF BLUES

i	$x_i$	$x_{(i)}$	$C_i$	$B_{i-1}$	$B_i$	$B_i x_{(i)}$
1	1.7986	1.0095	.9200	1.0000	.9200	.9287
2	1.0684	1.0586	.9111	.9200	.8382	.8873
3	1.3725	1.0684	.9000	.8382	.7544	.8060
4	1.1779	1.1267	.8857	.7544	.6682	.7529
5	1.4743	1.1779	.8667	.6682	.5791	.6821
6	1.0095	1.3725	.8400	.5791	.4864	.6676
7	4.8304	1.4743	.8000	.4864	.3891	.5737
8	1.0586	1.7986	.7333	.3891	.2854	.5133
9	1.1267	3.9974	.6000	.2854	.1712	.6844
10	3.9974	4.8304	.2000	.1712	.0342	.1652
$D = (c+1) \sum_{i=1}^{n-1} B_i + (c-1)B_n = 17.8733$						
$Y = (c+1) \sum_{i=1}^{n-1} B_i x_{(i)} + (c-1)B_n x_{(n)} - Dx_{(1)} = 4.9407$						
$\hat{a} = x_{(1)} - Y/[n(c-1)(n-2) - nD] = .9625$						
$\hat{b} = (x_{(1)} - \hat{a})(n-1) = 1.128$						

Step 2. Compute each  $B_i$  for  $i=1,2,\dots,n$  using equation (29). Thus:

For  $i=1$ ,  $B_1 = [1-2/2.5(10-1+1)]B_0 = (1-2/25.0)(1.000) = .9200$

For  $i=2$ ,  $B_2 = [1-2/2.5(10-2+1)]B_1 = (1-2/22.5)(.9200) = .8382$

.

.

.

For  $i=10$ ,  $B_{10} = [1-2/2.5(10-10+1)]B_9 = (1-2/2.5)(.1712) = .0342$

Table I lists all of the values of  $C_i = 1 - 2/c(n-i+1)$  and  $B_i = C_i B_{i-1}$  as computed from equation (29).

Step 3. Use the  $B_i$  to compute D from equation (21):

$$\begin{aligned} D &= (c+1)(B_1+B_2+\dots+B_9) + (c-1)B_{10} \\ &= (2.5 + 1)(.9200+ .8382+\dots+ .1712) + (2.5 - 1)(.0342) \\ &= (3.5)(5.092) + (1.5)(.0342) = 17.8733 \end{aligned}$$

Step 4. Use the  $x_{(i)}$ , D, and  $B_i$  values to compute Y from equation (22). Table I lists the values of  $B_i x_{(i)}$ :

$$\begin{aligned} Y &= (c+1)[B_1 x_{(1)} + B_2 x_{(2)} + \dots + B_9 x_{(9)}] + (c-1)B_{10} x_{(10)} - D x_{(1)} \\ &= (3.5)(.9287+ .8873+ \dots + .6844) \\ &\quad + (1.5)(.1652) - 17.8733(1.0095) \\ &= (3.5)(6.496) + .2478 - 18.0431 = 4.9407 \end{aligned}$$

Step 5. Use Y and D to compute  $\hat{a}$  from equation (17):

$$\begin{aligned} \hat{a} &= x_{(1)} - Y/[nc-1)(nc-2)-ncD] \\ &= 1.0095 - (4.9407)/[(25-1)(25-2) - 25(17.8733)] \\ &= 1.0095 - 4.9407/105.1675 = .9625 \end{aligned}$$

Step 6. Use  $\hat{a}$  to compute  $\hat{b}$  from equation (18):

$$\hat{b} = (x_{(1)} - \hat{a})(nc-1) = (1.0095 - .9625)(25 - 1) = 1.128$$

In this example, then, the BLUEs for a and b are  $\hat{a} = .9625$  and  $\hat{b} = 1.128$ . (The  $x_i$  values were actually generated from a Pareto distribution with  $a = b = 1$  and  $c = 2.5$ ). Once the BLUEs have been computed, the test statistics can be appropriately modified.

#### Modified Test Statistics

At the end of Chapter II, the standard forms of the Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) test statistics were presented. To use these "unmodified" statistics with their existing critical value tables, all parameters must be specified. When unknown location and scale parameters are involved, the test statistics must be modified to generate new critical value tables before they will produce accurate results. This section shows how to calculate the modified test statistics using an ordered sample and the BLUEs described in the preceding section. The notation and approach are adapted from Littell, McClave, and Offen (36:259-260).

Hypothesized Pareto CDF. Before computing the modified test statistics, the hypothesized Pareto CDF must be calculated for each value of the random sample. Let  $x_1, x_2, \dots, x_n$  be a random sample from the Pareto distribution with unknown location and scale parameters a and b, and known shape c; and let  $x_{(i)}$  denote the  $i$ th order statistic (20:70). The

appropriate BLUEs for location  $a$  and scale  $b$  (computed from the previous section), the specified shape  $c$ , and the  $n$  ordered Pareto deviates,  $x_{(i)}$ , are substituted into equation (15) to calculate the hypothesized Pareto CDF:

$$P_i = F(x_{(i)}; \hat{a}, \hat{b}, c) = 1 - [1 + (x_{(i)} - \hat{a}) / \hat{b}]^{-c} \quad (40)$$

for  $i = 1, 2, \dots, n$ . Note that for a given shape  $c$  (e.g.,  $c=2.5$  or  $c=4$ ) and sample size  $n$  (e.g.,  $n=10$  or  $n=30$ ), a specific, fixed pair of location and scale values (e.g.,  $a=b=1$  or  $a=0, b=1$ ) is used to produce the random Pareto deviates needed to compute the hypothesized CDF. This can be done without loss of generality because, as discussed in Chapter II, the use of invariant estimators (in this case the BLUEs) for location and scale ensures that the distribution of the test statistic depends only on the shape  $c$  and sample size  $n$ , and is independent of location and scale (36:260).

Example 2. In Example 1, the BLUEs for location  $a$  and scale  $b$  were found from a sample of size  $n=10$  generated from a Pareto distribution having shape  $c=2.5$ . In this example, the same sample of values  $x_1, x_2, \dots, x_{10}$  will be used to compute the hypothesized Pareto CDF from equation (40). Table II contains the values obtained while making the calculations. The columns for  $x_i$  and  $x_{(i)}$  are duplicated from Table I. The BLUEs  $\hat{a}$  and  $\hat{b}$  are as derived in Example 1.

Table II

CALCULATION OF HYPOTHESIZED PARETO CDF

i	$x_i$	$x_{(i)}$	$M_i$	$N_i$	$O_i$	$P_i$
1	1.7986	1.0095	.0470	.0417	.9030	.0970
2	1.0684	1.0586	.0961	.0852	.8151	.1849
3	1.3725	1.0684	.1059	.0939	.7990	.2010
4	1.1779	1.1267	.1642	.1456	.7119	.2881
5	1.4743	1.1779	.2154	.1910	.6460	.3540
6	1.0095	1.3725	.4100	.3635	.4607	.5393
7	4.8304	1.4743	.5118	.4537	.3925	.6075
8	1.0586	1.7986	.8361	.7412	.2500	.7500
9	1.1267	3.9974	3.0349	2.6905	.0382	.9618
10	3.9974	4.8304	3.8679	3.4290	.0242	.9758

$$M_i = x_{(i)} - \hat{a} = x_{(i)} - .9625$$

$$N_i = M_i / \hat{b} = M_i / 1.128$$

$$O_i = (1 + N_i)^{-c} = (1 + N_i)^{-2.5}$$

Hypothesized Pareto CDF:  $P_i = 1 - O_i$

Modified K-S Statistic. After computing all  $n$  of the values of  $P_i$  from equation (40), the modified Kolmogorov-Smirnov test statistic is found from equation (4) by substituting  $P_i$  in place of  $z_i$  in equation (5). Thus the modified test statistic in computational form is:

$$D = \max (D^+, D^-) \tag{41}$$

where

$$D^+ = \sup_{1 \leq i \leq n} [(i/n) - P_i] \quad \text{and} \quad D^- = \sup_{1 \leq i \leq n} [P_i - (i-1)/n] \tag{42}$$

Table III

## CALCULATION OF MODIFIED K-S TEST STATISTIC

i	x(i)	P <sub>i</sub>	i/n	(i-1)/n	D <sub>i</sub> <sup>+</sup>	D <sub>i</sub> <sup>-</sup>
1	1.0095	.0970	.1	.0	.0030	.0970
2	1.0586	.1849	.2	.1	.0151	.0849
3	1.0684	.2010	.3	.2	.0990	.0010
4	1.1267	.2881	.4	.3	.1119	-.0119
5	1.1779	.3540	.5	.4	(.1460)	-.0460
6	1.3725	.5393	.6	.5	.0607	.0393
7	1.4743	.6075	.7	.6	.0925	.0075
8	1.7986	.7500	.8	.7	.0500	.0500
9	3.9974	.9618	.9	.8	-.0618	(.1618)
10	4.8304	.9758	1.0	.9	.0242	.0758

$$D_i^+ = (i/n) - P_i = i/10 - P_i$$

$$D^+ = \sup [(i/n) - P_i] = .1460$$

$$D_i^- = P_i - (i-1)/n = P_i - (i-1)/10$$

$$D^- = \sup [P_i - (i-1)/n] = .1618$$

$$\text{K-S Statistic: } D = \max (D^+, D^-) = .1618$$

Example 3. Once the hypothesized Pareto CDF is computed, the values can be used to calculate the modified K-S test statistic. Table III continues the previous examples by showing the computations involved in calculating the modified K-S test statistic. As before, the calculations are based on the n=10 order statistics introduced in example 1, and the values P<sub>i</sub> of the hypothesized Pareto CDF as computed in example 2.

Table IV  
CALCULATION OF MODIFIED A-D TEST STATISTIC

j	$P_j$	$P_{n+1-j}$	$L_j$	$M_j$	$N_j$	$(2j-1)N_j$
1	.0970	.9758	-2.3330	-3.7214	-6.0544	-6.0544
2	.1849	.9618	-1.6879	-3.2649	-4.9528	-14.8584
3	.2010	.7500	-1.6045	-1.3863	-2.9908	-14.9540
4	.2881	.6075	-1.2444	-.9352	-2.1796	-15.2572
5	.3540	.5393	-1.0385	-.7750	-1.8135	-16.3215
6	.5393	.3540	-.6175	-.4370	-1.0545	-11.5995
7	.6075	.2881	-.4984	-.3398	-.8382	-10.8966
8	.7500	.2010	-.2877	-.2244	-.5121	-7.6815
9	.9618	.1849	-.0389	-.2040	-.2429	-4.1293
10	.9758	.0970	-.0245	-.1020	-.1265	-2.4035
$\sum_{j=1}^n (2j-1)N_j = -104.1559$						
$L_j = \ln P_j$		$M_j = \ln(1-P_{n+1-j})$		$N_j = L_j + M_j$		
$A^2 = -n - (1/n) \sum_{j=1}^n (2j-1) [\ln P_j + \ln(1-P_{n+1-j})]$ $= -10 - (1/10)(-104.1559) = .4156$						

Modified A-D Statistic. The modified Anderson-Darling test statistic is computed by substituting  $P_j$  from equation (40) in place of  $z_j$  in equation (9). Thus the computational form of the modified A-D test statistic is:

$$A^2 = -n - (1/n) \sum_{j=1}^n (2j-1) [\ln P_j + \ln(1-P_{n+1-j})] \quad (43)$$

Example 4. Table IV shows the calculations involved in finding the value of the modified A-D test statistic. The  $P_j$  values are as computed in example 2.

Table V

## CALCULATION OF MODIFIED C-VM TEST STATISTIC

$j$	$F_j$	$\frac{2j-1}{2n}$	$P_j - \frac{(2j-1)}{2n}$	$[P_j - \frac{(2j-1)}{2n}]^2$
1	.0970	.05	.0470	.0022
2	.1849	.15	.0349	.0012
3	.2010	.25	-.0490	.0024
4	.2881	.35	-.0619	.0038
5	.3540	.45	-.0960	.0092
6	.5393	.55	-.0107	.0001
7	.6075	.65	-.0425	.0018
8	.7500	.75	.0000	.0000
9	.9618	.85	.1118	.0125
10	.9758	.95	.0258	.0007
				$\sum_{j=1}^n = .0339$
$W^2 = [1/(12n)] + \sum_{j=1}^n [P_j - (2j-1)/2n]^2$ $= (1/120) + .0339 = .0423$				

Modified C-VM Statistic. The computational form of the modified Cramer-von Mises test statistic is found from equation (7) by substituting  $P_j$  for  $z_j$ :

$$W^2 = [1/(12n)] + \sum_{j=1}^n [P_j - (2j-1)/2n]^2 \quad (44)$$

Example 5. Table V shows the calculations involved in finding the value of the modified C-VM test statistic. The  $P_j$  values are as computed in example 2.

## Chapter Summary

Several applications for the Pareto distribution have been found in economics and operations research. It has played a major role in investigating the distributions of city population size, natural resources, stock price fluctuations, and oil field locations. Other studies show the Pareto can be used to model phenomena which may apply to Air Force interests, such as time-to-failure of equipment components, maintenance service times, nuclear fallout dispersion, and error clusters in communications circuits.

There are three basic forms of the Pareto distribution, each of which is a special case of the three-parameter form. The greater generality of the three-parameter form allows the Pareto distribution to be more useful in practical application. Various methods have been explored for estimation of Pareto parameters; but the best linear unbiased estimator (BLUE) is the only estimator known to possess the required invariance property for the three-parameter form.

For shape parameter  $c = .5, 1, \text{ or } 2$ , the BLUEs are computed from equations (34) and (35). When  $c = 1.5$ , the BLUEs are given by equations (37) to (39). For  $c = 2.5, 3, 3.5, \text{ or } 4$ , the BLUEs are computed from equations (17), (18), (21), (22), and (29). The BLUEs are used to compute the hypothesized distribution function from equation (40). The modified K-S, A-D, and C-VM test statistics can then be found using the methods presented in the next chapter.

## IV. METHODOLOGY

### Chapter Overview

This chapter describes the basic principles and specific procedures used to satisfy the research objectives of this thesis. Foremost is the Monte Carlo method used to generate the critical value tables of the modified K-S, A-D, and C-VM goodness-of-fit tests for the three-parameter Pareto distribution when only the shape parameter is specified.

### Basic Principles

This section deals with some of the basic principles used to generate critical values. It begins with an overview of the Monte Carlo method in general. Next is discussed the inverse transform technique used to generate random Pareto deviates. Then the selection of critical values is discussed. Finally, the use of plotting positions to determine percentiles is explained.

The Monte Carlo Method. Mathematics can be divided into theoretical and experimental categories. The primary distinction is that "theoreticians deduce conclusions from postulates, whereas experimentalists infer conclusions from observations" (21:1). The Monte Carlo method is a branch of experimental mathematics involving experiments using random

numbers. It has been used extensively in statistical analysis, operational research, nuclear physics, and several other fields where there are problems not easily solved by theoretical mathematics alone (21:2).

An important feature of the Monte Carlo method is its usual reliance on computers to simulate random processes (10:2). Also known as the method of statistical trials, it is basically a system of techniques which allows the modeling of random processes conveniently by digital computer. Before the advent of the computer, a study of a random process was considered to be complete when it was reduced to an analytical description. The computer has now made it convenient in many cases to solve an analytical problem by reducing it to a random process and then simulating that process (10:vii). Thus a basic principle of the method involves simulating statistical experiments through computational techniques, and then analysing numerical characteristics observed from these experiments (10:ix). For this reason, the Monte Carlo method can be defined as "the construction of an artificial random process possessing all the necessary properties, but which is in principle realizable by means of ordinary computational apparatus" (10:2).

The Monte Carlo method is typically used to solve problems of two basic types. A deterministic problem has no direct association with random processes. In this case the Monte Carlo method is often used when the problem can be

formulated in theoretical language but cannot be solved by theoretical means. Usually the approach is to recognize the underlying problem structure as resembling some apparently unrelated random process, and then solve the deterministic problem numerically by an appropriate Monte Carlo simulation.

In the case of a probabilistic problem, the Monte Carlo method is directly concerned with the behavior and outcome of random processes. The approach is to observe random variates, chosen so that they directly simulate the physical random processes of the original problem. The desired solution is then inferred from the behavior of the random numbers (21:2-4). The latter Monte Carlo approach was used in this thesis to generate the critical value tables for the goodness-of-fit tests.

The main weakness in the Monte Carlo method is that the answers it produces are to some degree uncertain since they are inferred from raw observational data consisting of random numbers. This weakness must be accounted for because:

Whenever one is inferring general laws on the basis of particular observations associated with them, the conclusions are uncertain inasmuch as the particular observations are only a more or less representative sample from the totality of all observations which might have been made. Good experimentation tries to ensure that the sample shall be more rather than less representative . . . [Monte Carlo answers] can nevertheless serve a useful purpose if we can manage to make the uncertainty fairly negligible, that is to say to make it unlikely that the answers are wrong by very much [21:4-5].

Thus there is usually no cause for concern if the uncertainty is negligible for practical purposes.

One way of reducing uncertainty is to base the Monte Carlo analysis on a larger number of observations. However, economic and time constraints must be considered. "Broadly speaking, there is a square law relationship between the error in an answer and the requisite number of observations; to reduce it tenfold calls for a hundredfold increase in the observations, and so on" (21:5). Therefore, to avoid using an inordinate amount of computer time, and to conserve financial resources, this thesis follows the common practice (9;43;49;52;54) of using 5000 repetitions rather than, say, 10000 in performing the Monte Carlo analysis.

The Inverse Transform Technique. To apply the Monte Carlo method to the problem at hand requires random samples from the Pareto distribution. The most practical way to obtain such samples is to use a computer program to produce a group of  $n$  numbers that seem to come from a Pareto population. In terminology adapted from Conover (13:323-324,360), these  $n$  numbers are called "random Pareto deviates" because they are deliberately generated to resemble observations on independent Pareto random variables. Previous AFIT theses (9;43;49; etc.) involved distributions for which computer programs to generate random samples were already available from the International Mathematical Statistics Library

(IMSL). IMSL does not contain a similar subroutine for the Pareto distribution; therefore, a computer program needed to be written to generate random Pareto deviates.

One common method of using a computer to generate random samples from a given distribution is to first generate a uniform random sample on  $(0,1)$  and then transform it into a new sample having the desired distribution. This method, called the inverse transform technique, uses the fact that the random variable  $R = F(X)$  is uniformly distributed on  $(0,1)$ , where  $X$  is a random variate (5:293-298). Thus, every variate is related to the uniform variate on  $(0,1)$  through its own inverse distribution function (26:22). Therefore, a set of uniformly distributed random numbers is required to generate a random sample from the Pareto distribution.

Conveniently, most random number generators are designed to generate random numbers which are uniformly distributed on the interval  $(0,1)$  (5:293). Hence, the inverse transform technique can be directly applied to a set of these random numbers to generate random Pareto deviates. However, the technique requires that for each random number  $r$ , the equation  $r = F(x)$  must be solved for the corresponding value of  $x = F^{-1}(r)$ . Therefore the technique is practical only when the CDF  $F(x)$  has an inverse which can be computed explicitly (5:294). Fortunately, the inverse transformation for the Pareto distribution can easily be expressed in closed form.

The inverse transform technique can be accomplished by the following four-step procedure (5:294-295):

Step 1. Compute the cumulative distribution function (CDF) of the desired random variable  $X$ . In this case, the CDF is the three-parameter Pareto CDF, given by equation (15) and repeated here for convenience:

$$F(x) = 1 - [1 + (x-a)/b]^{-c} \quad \text{for } x \geq a; b, c > 0$$

Step 2. Set  $F(X) = R$  on the range of  $X$ , where  $X$  represents a random Pareto variable. This then becomes:

$$1 - [1 + (X-a)/b]^{-c} = R \quad \text{for } x \geq a \quad (45)$$

Since  $X$  is a random variable (with the Pareto distribution in this case), then  $R$  is also a random variable. In fact,  $R$  has a uniform distribution over the interval  $(0,1)$  (5:295).

Step 3. Solve  $F(X)$  in terms of  $R$  to find  $X = F^{-1}(R)$ . In this case the inverse is found by solving equation (45):

$$1 - [1 + (X-a)/b]^{-c} = R$$

$$[1 + (X-a)/b]^{-c} = 1 - R$$

$$[b/b + (X-a)/b]^{-c} = 1 - R$$

$$(b + X - a)/b = (1 - R)^{-1/c}$$

$$b + X - a = b(1 - R)^{-1/c}$$

$$\text{Therefore } X = (a - b) + b(1 - R)^{-1/c} = F^{-1}(R) \quad (46)$$

Equation (46) is called a "random variate generator" (5:295) for the Pareto distribution. As explained in the discussion following equation (40), a specific, fixed pair of location and scale values can be used to generate the required deviates without loss of generality. For this thesis, the Pareto deviates were generated using location and scale parameters of 1. Substituting  $a=b=1$  into equation (46) gives:

$$\begin{aligned}
 X &= a - b + b(1 - R)^{-1/c} \\
 &= 1 - 1 + 1(1 - R)^{-1/c} \\
 &= (1 - R)^{-1/c}
 \end{aligned}
 \tag{47}$$

Since  $R$  is uniformly distributed from 0 to 1, then so is  $1-R$ ; thus  $R$  can replace  $1-R$  in equation (47) to yield the particular random variate generator used to produce the random Pareto variates for this thesis:

$$X = R^{-1/c} = (1/R)^{1/c}
 \tag{48}$$

Step 4. Generate  $n$  uniform random numbers  $R_1, R_2, \dots, R_n$  and compute the  $n$  random Pareto deviates from equation (48). The random numbers used for this thesis were generated on the AFIT VAX/VMS computer system using the IMSL subroutine GGUBS. Like most random number generators (5:293), GGUBS is designed to generate random numbers which

are uniformly distributed on the interval (0,1). Therefore, the inverse transform technique was applied to these random numbers to generate random Pareto deviates.

In step 3 of the inverse transform procedure, the choice of the location and scale values is arbitrary, and 1 was used here for convenience. It should be noted, however, that the deviates can be easily transformed into deviates from a different Pareto distribution (i.e., one having the same shape  $c$  but different location  $a'$  or scale  $b'$ ). The transformation stems from the fact that all variates having the same shape can be expressed in terms of the variate having location 0 and scale 1, as follows (26:21-22):

$$X_{a,b} = b X_{0,1} + a \quad (49)$$

where  $X_{a,b}$  denotes a Pareto variate with location  $a$  and scale  $b$  and  $X_{0,1}$  is a Pareto variate with location 0 and scale 1. The transformation to the different variate is then found by expressing the given variate in terms of the 0,1 variate, since:

$$X_{a,b} = b X_{0,1} + a \quad \text{implies} \quad X_{0,1} = (X_{a,b} - a)/b$$

$$\text{Thus } X_{a',b'} = b' X_{0,1} + a' = b' [(X_{a,b} - a)/b] + a' \quad (50)$$

Therefore, given a variate having a specific pair of

values for location and scale, equation (50) can be used to transform the variate to one having a different pair of location and scale parameters. For example, the transformation from a variate having location and scale  $a=b=1$  to one having location  $a'=2$  and scale  $b'=3$  is given by:

$$X_{2,3} = 3X_{0,1} + 2 = 3[(X_{1,1} - 1)/1] + 2 = 3X_{1,1} - 1$$

The random Pareto deviates generated by the inverse transform technique were used ultimately to compute values of the modified K-S, A-D, and C-VM test statistics. However, these test statistics can only be useful if their distribution functions are at least partially known (13:31). Thus, many test statistics were computed to determine the empirical distribution. Critical values were then identified using a plotting positions technique. Before examining the plotting positions technique, it may be helpful to understand how critical values are chosen.

Identifying Critical Values. The use of random deviates to generate critical value tables is based on the concept of hypothesis testing mentioned in Chapter II. Each group of  $n$  Pareto deviates represents a simulated sample from a parameter-specified Pareto distribution. This makes the null hypothesis " $H_0: H(x) = \text{the Pareto CDF}$ " true for each sample of  $n$  random Pareto deviates. For each of the three

tests (K-S, A-D, and C-VM), equations (41) - (44) were used to compute 5000 independent values of the test statistic under the condition that  $H_0$  is true (13:361). These 5000 values were then arranged in ascending order to form sets of 5000 order statistics. To determine critical values from these 5000 statistics (15000 total for all three tests), it is necessary to identify somehow the "critical region", i.e., the set of all values of the test statistic that would result in the erroneous decision to reject the true null hypothesis (13:78). Once the critical region is identified, then the critical values can be selected according to a desired "level of significance", or  $\alpha$ , which is the maximum probability of rejecting a true null hypothesis. Since the use of random Pareto deviates to compute the test statistics ensures that  $H_0$  is true,  $\alpha$  can be found by determining the probability that the test statistic will assume a value that falls within the critical region (13:78).

Since  $H_0$  is true and  $\alpha$  is the maximum probability of rejecting  $H_0$ , then the minimum probability of correctly accepting  $H_0$  is  $1 - \alpha$ . This value of  $1 - \alpha$  represents a certain percentile of the 5000 ordered test statistic values. For example, the 99th percentile is some number that the test statistic will exceed with probability .01 or less and will be less than with probability .99 or less (13:29). It is this percentile relationship that is used to select critical values from the 5000 test statistics.

One possible method of using the percentiles to determine critical values is to simply select the test statistic value corresponding to the desired percentile level and make that the critical value. For example, under this method, out of a set of 5000 ordered test statistic values, the critical value for the 90th percentile would simply be the 4500th value (52:6). This method has some disadvantages, however, especially when the test statistics, which represent a discrete distribution, are used to determine critical values for a continuous distribution. More recently, the plotting position technique has become popular as a more accurate method of selecting critical values for continuous distributions (43:7).

The Plotting Positions Technique. The plotting positions technique is one popular method of determining percentiles of the distribution underlying a set of  $n$  ordered sample values (24:1619; 25:317). The technique involves using a large number of discrete values of the ordered test statistics and locating them on a continuous spectrum by representing the spaces between them as piecewise linear functions. This makes it possible to linearly interpolate the desired percentiles between discrete values of the test statistics, thus obtaining more accurate critical values (43:7; 52:6).

Each ordered value may be assigned a plotting position

which is its cumulative probability, thus allowing each order statistic to be mapped onto a probability scale from 0 to 1. As seen from equation (2), the distribution function of these  $n$  observations is a step function which jumps from  $(i-1)/n$  to  $i/n$  at the  $i$ th order statistic of the sample. However, if the plotting position  $i/n$  is used, the largest value cannot be plotted, while if  $(i-1)/n$  is used, the smallest value cannot be plotted (24:1615). Therefore, numerous alternative plotting conventions have been proposed, most of which have been summarized by Harter (24), who presents various arguments for and against each. Harter also conducted a Monte Carlo analysis of plotting positions for several distributions and concluded that ". . . the optimum choice of plotting positions depends not only on the purpose of the investigation, but also (definitely) on the distribution of the variable under consideration" (25:342).

While Harter made no specific recommendation for the Pareto, he did observe that, "As samples increase above a sample size of 20, the differences among the positions determined by any method of estimation decrease to the point where they are practically unimportant" (24:1621). He also noted that "in practice, plotting positions differ little compared with the randomness of the data" (24:1622). Since this thesis employed 5000 independent values of each test statistic, well in excess of the 20 cited by Harter, use of a single plotting convention seems justified.

The plotting convention selected for this thesis is the median rank, which is closely approximated by the plotting position (24:1617):

$$Y_i = (i-0.3)/(n+0.4) \quad (51)$$

where  $i = 1, \dots, n$  and for this thesis,  $n=5000$ . Thus each  $Y_i$  value lies in the interval  $(0,1)$ . The median ranks position yields median unbiased estimates of  $x_i$  for a specified  $F(x_i)$  and of  $F(x_i)$  for a specified  $x_i$  (24:1625). Also, in highly skewed distributions, the median ranks position tends to be more accurate than other conventions (31:300). Another advantage is that values of the median ranks have been tabulated for sample sizes of 1 to 50, i.e.,  $n = 1(1)50$  (31:486-489).

A detailed illustration showing how to use plotting positions to determine critical values was presented by Ream (43:11-23), and will only be summarized here. In graphical terms, the technique effectively plots the 5000 ordered test statistic values  $X_{(1)}, X_{(2)}, \dots, X_{(5000)}$  along the abscissa (horizontal) axis and the 5000 plotting position values  $Y_1, Y_2, \dots, Y_{5000}$  computed from equation (51) along the ordinate (vertical) axis. These values are assigned to positions 2 to 5001 on their respective axes. On the vertical axis, the interval  $[0,1]$  is completed by entering the endpoints  $Y_0 = 0$  at the 1st position and  $Y_{5001} = 1$  at the

5002nd position. The corresponding endpoints on the horizontal axis are found by linear extrapolation. Thus, in using the computer to program this technique, the arrays corresponding to the horizontal and vertical axes are each composed of 5002 entries, i.e., the original 5000 values and two extrapolated endpoints.

To map the collection of 5000 discrete values onto a fully continuous line between 0 and 1 requires extrapolation of the endpoints of the plotting axes. The first point on the horizontal axis,  $X_{(0)}$ , is computed by linearly extrapolating from the second and third points (i.e., the first and second order statistics), subject to a non-negativity restriction. Extrapolation is performed by using the standard linear slope-intercept formula  $Y = mX + b$  to compute the endpoints  $X_{(0)}$  and  $X_{(5001)}$ . To find the first endpoint on the horizontal axis, the slope is calculated by:

$$m = \frac{Y_2 - Y_1}{X_{(2)} - X_{(1)}} \quad (52)$$

and the intercept is:

$$b = Y_1 - m X_{(1)} \quad (53)$$

Then the lower endpoint  $X_{(0)}$  is found by:

$$X_{(0)} = (Y_0 - b)/m = (0-b)/m = -b/m$$

The nonnegativity restriction means that whenever  $-b/m < 0$ , then  $X_{(0)}$  is simply set to 0. Thus:

$$X_{(0)} = \max(0, -b/m) \quad (54)$$

The higher endpoint  $X_{(5001)}$  is found in the same way as the lower endpoint. The slope is

$$m = \frac{Y_{5000} - Y_{4999}}{X_{(5000)} - X_{(4999)}} \quad (55)$$

and the intercept is:

$$b = Y_{4999} - m X_{(4999)} \quad (56)$$

Then the second endpoint  $X_{(5001)}$  is extrapolated by:

$$X_{(5001)} = (Y_{5001} - b)/m = (1-b)/m \quad (57)$$

Once the endpoints are added to the abscissa and ordinate axes, the 5002 discrete points on the graph are "connected" by straight lines, thus producing a completely continuous, piecewise linear function. The range of this continuous function is the interval  $[0,1]$  and contains the 5000 median rank values as well as the endpoints 0 and 1.

Its domain contains the set of 5000 test statistic values and their 2 extrapolated endpoints.

As shown in Figure 5, the desired critical value for a given percentile is found by linearly interpolating between two of the 5002 points used to construct the now continuous graph. For example, to find the 95th percentile ( $\alpha = .05$ ), the largest plotting position  $Y_j$  is found such that  $Y_j \leq .95$ ; thus  $Y_{j+1}$  is the first position greater than .95. Then the critical value corresponding to the 95th percentile is found by linearly interpolating between the points  $(X_{(j)}, Y_j)$  and  $(X_{(j+1)}, Y_{j+1})$  using the formulas:

$$m = \frac{Y_{j+1} - Y_j}{X_{(j+1)} - X_{(j)}} \quad (58)$$

$$b = Y_j - m X_{(j)} \quad (59)$$

$$C_p = (p - b)/m \quad (60)$$

where  $C_p$  is the critical value for the the 100pth percentile. For this thesis, critical values were calculated for  $p = .80, .85, .90, .95,$  and  $.99,$  corresponding to the levels of significance  $\alpha = .20, .15, .10, .05,$  and  $.01.$

The specific plotting position procedure performed for this thesis is described in step 7 of the next section.

PLOTTING  
POSITIONS

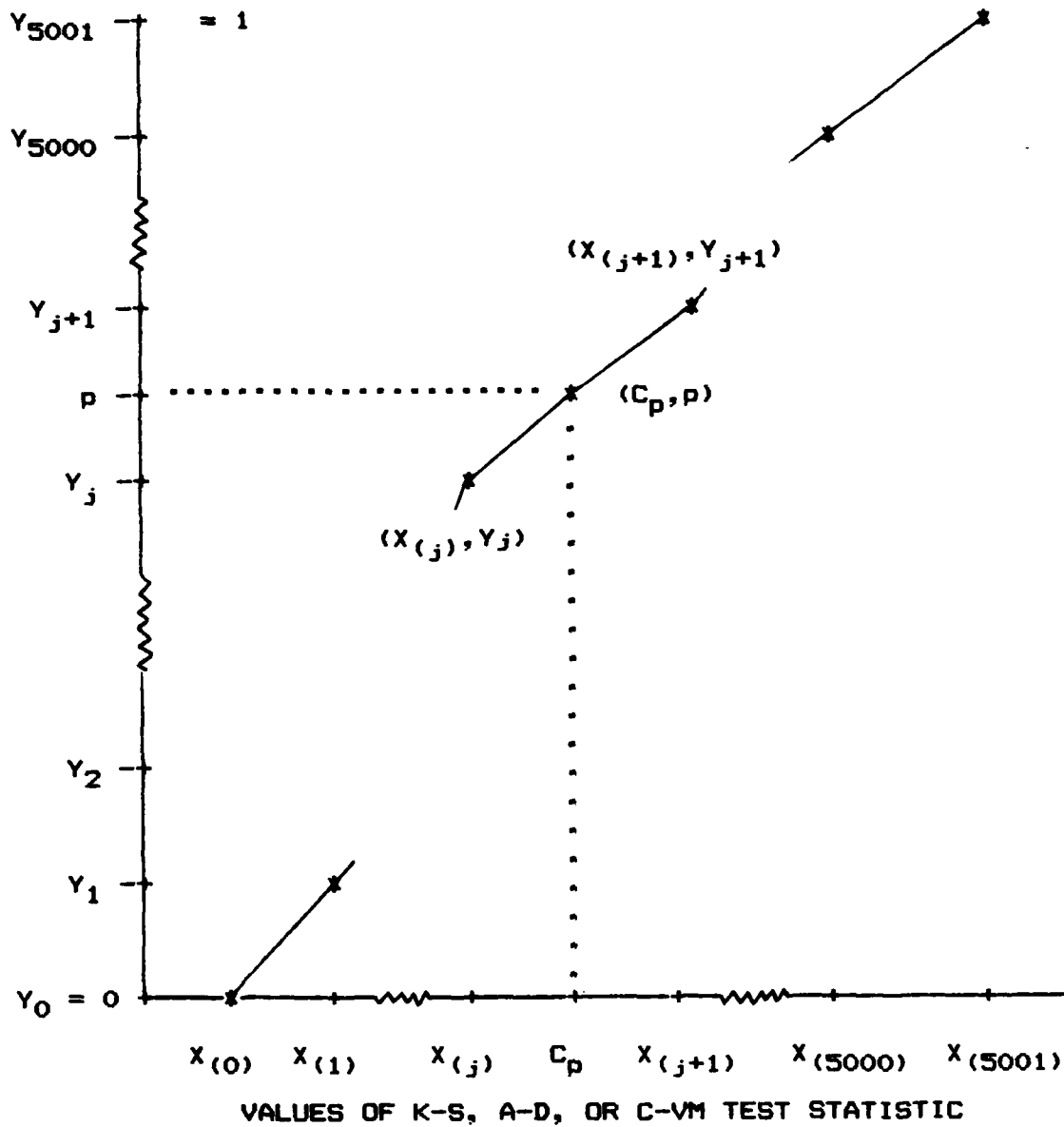


Fig 5. Using Test Statistics  $X_{(j)}$  and Plotting Positions  $Y_j$  to Find Critical Value  $C_p^j$  for the 100(p)th Percentile ( $p = .99, .95, .90, .85, .80$ ).

### Specific Procedures

By applying the basic principles and techniques described in the previous section, the K-S, A-D, and C-VM tests were modified to produce new goodness-of-fit tests for the Pareto distribution.

The research effort was performed in three stages, each corresponding to one of the three research objectives listed in Chapter I. The first stage consisted of a nine-step Monte Carlo simulation procedure to produce critical value tables for the modified K-S, A-D, and C-VM tests. The second stage of the research compared the powers of the three modified tests using eight alternative distributions. Finally, a regression analysis was performed to determine the functional relationship between the critical values and the shape parameters. Computer programs were written to accomplish the first two stages. The third stage was performed manually by using a hand calculator to compute linear relationships by the method of least squares.

Stage 1: Generating Critical Value Tables. During the first stage, critical value tables were generated using Monte Carlo simulation. A FORTRAN computer program was written for this purpose and is contained in Appendix A. The accompanying flow chart illustrates the logic flow of the program. The following nine steps outline the procedure used:

Step 1 - Generate the Data. Random deviates for a given sample size  $n$  were generated from a specified Pareto distribution by using the IMSL routine GGUBS to generate  $n$  random numbers, and then applying the inverse transform technique (equation 48).

Step 2 - Order the Data. Next, the  $n$  random deviates  $x_1, x_2, \dots, x_n$  were converted to order statistics  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$  by arranging them in ascending order using the IMSL subroutine VSRTA.

Step 3 - Estimate the Parameters. The ordered Pareto deviates were then used to find the best linear unbiased estimates of the scale and location parameters as explained in the "Summary of BLUEs" section of Chapter III.

Step 4 - Compute the Hypothesized CDF. The estimated parameters found in step 3 were used with the  $n$  ordered Pareto deviates from step 2 to calculate the hypothesized cumulative distribution function (CDF)  $P_i$  for  $i=1, 2, \dots, n$  (equation 40 in chapter III).

Step 5 - Calculate the Test Statistics. Based on the hypothesized CDF and the BLUEs, the modified K-S, A-D, and C-VM statistics were next calculated using equations (42), (43), and (44).

Step 6 - Generate 5000 Statistics. Each of these five steps were repeated 5000 times to generate 5000 independent K-S, A-D, and C-VM statistical values  $x_1, x_2, \dots, x_{5000}$ .

Step 7 - Find the Critical Values. For each of the three tests, the 5000 statistics were ordered as in step 2. Using the median ranks plotting position technique (equation 51), the 80th, 85th, 90th, 95th, and 99th percentiles of the distributions of each test statistic were calculated by linear interpolation. These percentiles correspond, respectively, to the .20, .15, .10, .05, and .01 levels of significance and served as the critical values for the modified K-S, A-D, and C-VM goodness-of-fit tests. The specific step-by-step process was to:

a. Use the IMSL subroutine VSRTA to order the 5000 test statistics, thus forming the 5000 order statistics  $X_{(1)}, X_{(2)}, \dots, X_{(5000)}$ .

b. Use equation (51) to compute the 5000 plotting positions  $Y_1, Y_2, \dots, Y_{5000}$ . Also, set  $Y_0 = 0$  and  $Y_{5001} = 1$ .

c. Use equations (52), (53), and (54) to find  $X_{(0)}$ . Similarly, use equations (55), (56), and (57) to find  $X_{(5001)}$ .

d. For a given  $p$ , find the largest  $Y_j$  such that  $Y_j \leq p$ ; then use equations (58), (59), and (60) to find the critical value  $C_p$  representing the 100(p)th percentile. Repeat this step for  $p = .80, .85, .90, .95, \text{ and } .99$ .

Step 8 - Repeat for Sample Sizes. To evaluate the effect of sample size on the critical values, steps 1 through 7 were repeated for each sample size  $n$ . This thesis

followed the common practice (9:15) of using sample sizes of  $n$  equal to 5, 10, 15, 20, 25, and 30.

Step 9 - Repeat for Shape Parameters. Steps 1 through 8 were repeated for specified shape parameters 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5, and 4.0. The critical values were then arranged into tabular form and appear in Chapter V, Tables VI - VIII.

Stage 2: Comparing Power. The second stage of the research compared the powers of the modified K-S, A-D, and C-VM tests against the Chi-square to determine which test can best detect a false Pareto distribution hypothesis. As explained in Chapter II, the power of a statistical test is the probability of correctly rejecting a false null hypothesis. The null hypothesis that a set of sample deviates follows a Pareto distribution with a specified shape parameter was tested against the alternative hypothesis that the sample deviates follow some other distribution:

$H_0$ : Sample deviates follow a Pareto CDF with shape  $c$

$H_1$ : They follow some other distribution

For this thesis, the power study was conducted for both  $c = 1$  and  $c = 3.5$  in the null hypothesis.

The Chi-square portion of the study was performed as described by Banks and Carson (5:352-356) using five

equiprobable (ie,  $p = .20$ ) class intervals (or cells) with expected frequencies of 3 observations per cell for  $n = 15$  and 5 per cell for  $n = 25$ . The endpoints of each cell were computed from the Pareto CDF (equation 15) as follows:

$$F(e_i) = 1 - [1 + (e_i - a)/b]^{-c} \quad (61)$$

where  $e_1, e_2, e_3, e_4$  represent the right endpoints (maximum value) of the first four cells. Since  $F(e_i)$  is the cumulative area from 0 to  $e_i$ , then  $F(e_i) = ip = .2i$ , so equation (61) leads to:

$$\begin{aligned} .2i &= 1 - [1 + (e_i - a)/b]^{-c} \\ [1 + (e_i - a)/b]^{-c} &= 1 - .2i \\ 1 + (e_i - a)/b &= (1 - .2i)^{-1/c} \\ b + e_i - a &= b(1 - .2i)^{-1/c} \\ e_i &= a - b + b(1 - .2i)^{-1/c} \end{aligned}$$

After substituting the BLUEs for location and scale into this last expression, the right endpoints were found by:

$$e_i = \hat{a} - \hat{b} + \hat{b}(1 - .2i)^{-1/c} \quad (62)$$

Assuming a true Pareto null hypothesis, the four endpoints  $e_1, \dots, e_4$  essentially divide the real line into five equiprobable class intervals. Given a random sample, the

number of observations occurring within each cell were counted. The Chi-square test statistic was then computed by (5:350):

$$\chi^2 = \sum_{i=1}^5 [(O_i - E)^2]/E \quad (63)$$

where  $O_i$  is the number of observations occurring in cell  $i$  and  $E = n/5$  is the expected frequency in each interval. The distribution of this test statistic approximately follows a chi-square CDF with  $s-1-k$  degrees of freedom (13:194) where  $s$  is the number of cells (i.e.,  $s = 5$ ) and  $k$  is the number of parameters estimated from the sample (i.e.,  $k = 2$ ).

Using the IMSL subroutines GGWIB, GGAMR, GGBTR, GGEXN, and GGNML, random deviates from different distributions of sample size  $n$  were generated. The alternate distributions used were, respectively, the Weibull at shape parameter 3.5, the Gamma at shape parameter 2.0, the Beta at parameters  $P = 2$  and  $Q = 3$ , the exponential with mean = 2, and the normal distribution. Also tested were three sets of Pareto deviates generated by a FORTRAN subroutine. The first Pareto deviate set was generated using  $a = b = c = 1.0$ ; the second set used  $a = 2$ ,  $b = 3$ , and  $c = 3.5$ ; the third used  $a = 10$ ,  $b = 5$ , and  $c = 2.0$ . Five thousand random samples of size  $n$  were generated for each of the alternate distributions.

The K-S, A-D, C-VM, and Chi-square test statistics were then calculated under the null hypothesis that the random

deviates follow the Pareto distribution with specified shape  $c = 1.0$  or  $3.5$ . To determine whether to reject the null hypothesis, the calculated K-S, A-D, and C-VM statistics were compared to the corresponding critical value obtained in stage one. The computed Chi-square test statistic was compared against two sets of critical values. The first set was taken from a standard table of Chi-square critical values (13:432) based on 2 degrees of freedom. The second set of critical values was generated by using equations (62) and (63) and applying the 9-step, 5000-repetition Monte Carlo procedure described in the previous section.

This procedure of comparing test statistics against critical values was repeated 5000 times for each distribution and test. The number of times each statistic exceeded the respective critical value was counted for each sample size. This total, representing the number of rejections of the null hypothesis, was divided by the total number of tests performed (5000), to yield an hypothesis rejection quotient. For a random sample generated from the hypothesized Pareto distribution, the quotient represents the rate of erroneous rejection of a true null hypothesis; thus, it is expected to be approximately the level of significance  $\alpha$ , which is the probability of committing a Type I error (13:78). In those cases involving random samples generated from an alternative distribution, the quotient represents the power of the test, since it approximates the probability of correctly rejecting

a false null hypothesis (13:79).

A FORTRAN program, written to compute the hypothesis rejection rates and accomplish the power study, is contained in Appendix B. Figure 7 in Appendix B shows how the program used the following 9-step process:

Step 1. Use IMSL or inverse transform to generate  $n$  random deviates from a selected distribution.

Step 2. Assume the null hypothesis that this set of  $n$  deviates follows the Pareto of given shape  $c = 1.0$ . Then perform steps 2-5 of the previous section to compute the values of the Chi-square (eqn 63) and modified K-S, A-D, and C-VM test statistics (eqns 42-44).

Step 3. For a given level of significance  $\alpha$ , compare the test statistic value against the appropriate critical value found in the previous section. If the test statistic value equals or exceeds the critical value,  $H_0$  is rejected.

Step 4. Repeat steps 1-3 5000 times, each time using a different seed to generate the deviates.

Step 5. Count the number of times  $H_0$  was rejected and divide by 5000 to obtain the power.

Step 6. Repeat steps 1-5 for each alternative distribution considered.

Step 7. Repeat steps 1-6 for sample sizes  $n = 5$ , 15, and 25.

Step 8. Repeat steps 1-7 for  $\alpha = .05$  and  $.01$ .

Step 9. Repeat steps 1-8 using hypothesized Pareto shape  $c = 3.5$ . The power values were then arranged into tabular form and appear in Chapter V, Tables IX and X.

Stage 3: Determining Functional Relationship. The third and final stage of the research was to determine what (if any) functional relationship exists between the shape parameter and the critical values generated. This relationship can then be used to interpolate critical values corresponding to parameters not found in the generated tables.

To accomplish this stage, shape parameters and critical values were examined for linear relationships. In an attempt to "fit" the data to a line, a linear regression was performed using the method of least squares (13:263-271), which minimizes the sum of the squares of the deviations of the actual data points from the straight line of "best" fit (5:359-363). Where applicable, the correlation coefficient (13:250-251) was also found.

Linear regression is a capability available on many hand calculators currently on the market, so it was unnecessary to write a separate computer program to perform this function. For each level of significance and sample size, critical values from Tables VI - VIII were paired against a corresponding Pareto shape parameter. The

regression and correlation coefficients were then obtained manually by using the linear regression keys on a Texas Instruments TI-55-II calculator. The results are contained in Chapter V, Tables XI and XII.

#### Chapter Summary

The research for this thesis was performed by applying the Monte Carlo method using 5000 repetitions to generate critical value tables and a power study.

In stage 1, random Pareto deviates were generated by using the inverse transform technique, and 5000 test statistics were computed for each test. The median ranks plotting positions technique was then used to select critical values from the 5000 test statistics. In stage 2, the powers of the modified K-S, A-D, and C-VM tests were compared against the power of the Chi-square test. The calculations were performed by computer programs written to accomplish a 9-step Monte Carlo procedure. Stage 3 involved manual calculations based on the method of least squares to find linear relationships between shape parameters and critical values.

The results of this research are presented in the next chapter.

## V. RESULTS AND APPLICATION

### Chapter Overview

This chapter shows the results obtained from carrying out the methodology described in Chapter IV. In response to the three research objectives listed in Chapter I, tables of critical values for the modified K-S, A-D, and C-VM tests are presented. Also included are tables comparing powers of the K-S, A-D, and C-VM statistics against the Chi-square. Tables of regression coefficients are presented as well. The use of the tables is explained, and an example is described.

### Critical Value Tables

Table VI contains critical values for the modified Kolmogorov-Smirnov Test. The modified Anderson-Darling critical values appear in Table VII. In Table VIII, the modified Cramer-von Mises critical values are presented. Critical values are presented for each level of significance  $\alpha = .20, .15, .10, .05, \text{ and } .01$ ; sample sizes  $n = 5, 10, 15, 20, 25, \text{ and } 30$ ; and Pareto shape parameters  $.5, 1, 1.5, 2, 2.5, 3, 3.5, \text{ and } 4$ . It is important to note that for shape  $c = 0.5$ , the presented critical values correspond to sample size  $n = 6$  instead of  $n = 5$ . As explained in Chapter III, this exception is necessary since the BLUEs could not be computed for the case where  $c = .5, n = 5$ .

Table VI

## CRITICAL VALUES FOR THE MODIFIED KOLMOGOROV-SMIRNOV TEST

$\alpha$	n	Pareto Shape Parameter c							
		0.5*	1.0	1.5	2.0	2.5	3.0	3.5	4.0
.20	5*	.400	.318	.289	.286	.283	.286	.293	.297
	10	.255	.222	.217	.219	.222	.225	.228	.231
	15	.204	.184	.184	.185	.187	.191	.192	.197
	20	.175	.160	.160	.163	.167	.168	.170	.171
	25	.155	.144	.146	.148	.149	.153	.154	.155
	30	.142	.133	.135	.135	.138	.139	.142	.141
.15	5*	.426	.328	.296	.294	.293	.298	.306	.309
	10	.268	.230	.226	.228	.232	.236	.239	.242
	15	.214	.191	.191	.193	.196	.199	.203	.207
	20	.184	.167	.168	.172	.175	.176	.178	.179
	25	.163	.150	.152	.155	.157	.160	.161	.163
	30	.149	.138	.142	.142	.145	.146	.150	.149
.10	5*	.467	.341	.305	.306	.307	.314	.323	.327
	10	.284	.241	.236	.239	.245	.251	.253	.258
	15	.227	.200	.201	.204	.208	.212	.216	.219
	20	.196	.176	.176	.182	.185	.187	.188	.191
	25	.173	.158	.160	.164	.166	.171	.170	.173
	30	.159	.145	.149	.151	.153	.155	.161	.159
.05	5*	.525	.368	.321	.323	.328	.335	.349	.353
	10	.308	.257	.254	.258	.265	.272	.277	.282
	15	.248	.216	.217	.223	.227	.231	.238	.239
	20	.213	.188	.191	.197	.201	.206	.205	.209
	25	.189	.170	.174	.177	.180	.186	.189	.192
	30	.173	.156	.162	.165	.167	.169	.175	.174
.01	5*	.609	.407	.378	.363	.361	.369	.382	.391
	10	.348	.297	.290	.300	.308	.314	.322	.326
	15	.289	.247	.251	.258	.265	.266	.274	.282
	20	.247	.216	.221	.233	.233	.237	.238	.249
	25	.222	.201	.201	.208	.210	.220	.218	.225
	30	.204	.180	.187	.189	.196	.199	.207	.207

\*NOTE: For shape  $c = 0.5$ , critical values correspond to sample size  $n = 6$  instead of  $n = 5$ .

Table VII

## CRITICAL VALUES FOR THE MODIFIED ANDERSON-DARLING TEST

$\alpha$	n	Pareto Shape Parameter c							
		0.5*	1.0	1.5	2.0	2.5	3.0	3.5	4.0
.20	5*	1.344	.736	.568	.546	.503	.494	.499	.497
	10	.780	.587	.544	.535	.541	.545	.540	.551
	15	.706	.589	.562	.559	.562	.568	.581	.588
	20	.684	.582	.571	.586	.591	.586	.599	.604
	25	.664	.588	.591	.585	.600	.608	.624	.621
	30	.674	.598	.607	.600	.606	.621	.638	.625
.15	5*	1.668	.835	.628	.602	.545	.532	.538	.537
	10	.875	.646	.594	.589	.588	.601	.597	.610
	15	.789	.645	.621	.612	.626	.630	.650	.659
	20	.764	.639	.629	.646	.656	.659	.661	.673
	25	.750	.653	.655	.652	.660	.672	.692	.694
	30	.756	.665	.679	.665	.678	.688	.708	.690
.10	5*	2.100	.966	.709	.671	.606	.585	.590	.599
	10	1.031	.726	.675	.655	.654	.678	.677	.691
	15	.917	.727	.705	.704	.705	.707	.748	.756
	20	.862	.732	.718	.734	.740	.747	.751	.755
	25	.853	.748	.742	.766	.750	.769	.788	.801
	30	.862	.756	.777	.768	.774	.776	.822	.791
.05	5*	2.903	1.237	.849	.791	.702	.683	.684	.687
	10	1.311	.886	.808	.783	.788	.805	.818	.835
	15	1.154	.891	.849	.853	.836	.852	.899	.927
	20	1.053	.874	.866	.898	.902	.917	.917	.925
	25	1.055	.915	.910	.940	.904	.926	.952	.987
	30	1.070	.913	.952	.947	.960	.937	.999	.990
.01	5*	4.877	2.076	1.145	1.100	.932	.883	.913	.903
	10	1.872	1.303	1.102	1.113	1.100	1.147	1.169	1.200
	15	1.705	1.250	1.229	1.154	1.256	1.316	1.269	1.358
	20	1.535	1.245	1.255	1.318	1.326	1.353	1.330	1.398
	25	1.543	1.312	1.286	1.358	1.253	1.427	1.450	1.441
	30	1.631	1.337	1.361	1.368	1.401	1.413	1.500	1.475

\*NOTE: For shape c = 0.5, critical values correspond to sample size n = 6 instead of n = 5.

Table VIII

## CRITICAL VALUES FOR THE MODIFIED CRAMER-VON MISES TEST

$\alpha$	n	Pareto Shape Parameter c							
		0.5*	1.0	1.5	2.0	2.5	3.0	3.5	4.0
.20	5*	.212	.103	.078	.078	.077	.079	.082	.083
	10	.121	.083	.081	.082	.086	.088	.089	.092
	15	.108	.086	.086	.086	.090	.092	.096	.098
	20	.104	.085	.086	.091	.094	.094	.097	.099
	25	.101	.086	.090	.090	.094	.097	.101	.102
	30	.100	.086	.091	.092	.095	.097	.102	.101
.15	5*	.251	.112	.083	.084	.084	.085	.089	.092
	10	.135	.093	.090	.091	.096	.099	.099	.102
	15	.120	.094	.094	.097	.100	.103	.108	.111
	20	.115	.095	.096	.102	.105	.106	.108	.110
	25	.112	.097	.101	.102	.105	.108	.112	.115
	30	.112	.096	.102	.105	.106	.109	.116	.114
.10	5*	.304	.123	.093	.093	.093	.096	.100	.103
	10	.154	.105	.102	.104	.108	.114	.113	.119
	15	.139	.109	.109	.114	.116	.120	.126	.130
	20	.133	.109	.111	.118	.120	.123	.125	.129
	25	.130	.112	.115	.120	.121	.127	.130	.133
	30	.129	.110	.121	.121	.124	.127	.135	.131
.05	5*	.381	.139	.113	.111	.111	.112	.119	.120
	10	.184	.127	.125	.126	.131	.139	.142	.147
	15	.172	.131	.134	.140	.142	.144	.156	.161
	20	.163	.133	.136	.145	.148	.155	.153	.161
	25	.157	.139	.142	.149	.148	.157	.162	.168
	30	.159	.137	.151	.150	.156	.154	.169	.166
.01	5*	.508	.174	.157	.148	.149	.150	.157	.163
	10	.251	.191	.174	.182	.194	.199	.202	.209
	15	.255	.192	.193	.198	.207	.222	.222	.238
	20	.233	.195	.201	.212	.220	.225	.226	.239
	25	.245	.199	.206	.224	.215	.238	.240	.250
	30	.247	.202	.217	.218	.230	.243	.251	.251

\*NOTE: For shape  $c = 0.5$ , critical values correspond to sample size  $n = 6$  instead of  $n = 5$ .

### Power Comparison Tables

Tables IX and X display the results of the power analysis. For sample sizes  $n = 5, 15,$  and  $25,$  the tables indicate relative power of the K-S, A-D, and C-VM tests to reject a null hypothesis when the hypothesis claims that a random sample of data follows a Pareto distribution. For sample sizes  $n = 15$  and  $25,$  the power of the Chi-square test is also included. Table IX shows power values when the null hypothesized Pareto CDF has shape parameter  $c = 1.0.$  In Table X, the hypothesized shape parameter is  $c = 3.5.$  Both tables examine power performance against eight different distributions, including three variations of the Pareto distribution having different sets of parameters.

The power tables are divided into two levels of significance,  $\alpha = .05$  and  $.01.$  In Table IX, the first column corresponds to a Pareto distribution with shape  $c = 1.0.$  Thus, the values in the first column of Table IX approximate the level of significance  $\alpha,$  since they represent rejection rates of the null hypothesis when  $H_0$  is true. Similarly in Table X, the second column represents a true null hypothesis since the underlying data was generated from a Pareto distribution with shape parameter  $c = 3.5.$  Aside from these two exceptions, all other columns represent power values since they indicate rejection rates of the null hypothesis when  $H_0$  is in fact false. A note following the tables indicates parameters of the alternate distributions.

Table IX

POWER TEST FOR THE PARETO DISTRIBUTION  
 $H_0$ : Pareto Distribution at Shape  $c = 1.0$   
 $H_1$ : The data follow another distribution

Level of Significance = .05

n	Test	Alternate Distributions*							
		Par.1	Par.2	Par.3	Weibl	Gamma	Beta	Expon	Norml
5	K-S	0.046	0.061	0.050	0.288	0.123	0.227	0.074	0.311
	A-D	0.048	0.014	0.022	0.007	0.006	0.008	0.009	0.007
	CVM	0.050	0.063	0.051	0.283	0.127	0.224	0.076	0.307
15	K-S	0.048	0.145	0.107	0.979	0.657	0.933	0.290	0.979
	A-D	0.052	0.126	0.083	0.966	0.644	0.898	0.266	0.965
	CVM	0.052	0.173	0.121	0.974	0.697	0.915	0.329	0.973
	$\chi^2$	0.043	0.118	0.086	0.860	0.480	0.738	0.235	0.878
25	K-S	0.052	0.248	0.138	1.000	0.927	1.000	0.503	1.000
	A-D	0.049	0.250	0.128	1.000	0.937	0.998	0.528	1.000
	CVM	0.050	0.256	0.143	0.999	0.926	0.996	0.504	1.000
	$\chi^2$	0.045	0.178	0.105	0.999	0.823	0.996	0.377	0.999

Level of Significance = .01

5	K-S	0.010	0.021	0.021	0.171	0.067	0.115	0.034	0.172
	A-D	0.009	0.002	0.004	0.000	0.000	0.000	0.000	0.001
	CVM	0.010	0.019	0.019	0.155	0.059	0.098	0.030	0.160
15	K-S	0.015	0.059	0.035	0.941	0.448	0.852	0.150	0.937
	A-D	0.011	0.038	0.021	0.875	0.356	0.716	0.103	0.878
	CVM	0.016	0.062	0.034	0.906	0.439	0.777	0.139	0.910
	$\chi^2$	0.006	0.031	0.016	0.645	0.172	0.400	0.064	0.669
25	K-S	0.010	0.086	0.039	0.999	0.774	0.992	0.250	0.998
	A-D	0.009	0.080	0.032	0.997	0.778	0.982	0.247	0.997
	CVM	0.010	0.100	0.046	0.997	0.792	0.982	0.274	0.998
	$\chi^2$	0.011	0.061	0.033	0.964	0.594	0.884	0.172	0.971

\* Key to Alternate Distributions:

Par.1 - Pareto (a=1, b=1, c=1)  
 Par.2 - Pareto (a=2, b=3, c=3.5)  
 Par.3 - Pareto (a=10, b=5, c=2)  
 Weibl - Weibull (shape = 3.5)

Gamma - Gamma (shape = 2)  
 Beta - Beta (P=2, Q=3)  
 Expon - Exponential (mean = 2)  
 Norml - Normal distribution

Table X

POWER TEST FOR THE PARETO DISTRIBUTION  
 $H_0$ : Pareto Distribution at Shape  $c = 3.5$   
 $H_1$ : The data follow another distribution

Level of Significance = .05

		Alternate Distributions*							
n	Test	Par.1	Par.2	Par.3	Weibl	Gamma	Beta	Expon	Norml
5	K-S	0.120	0.048	0.051	0.160	0.065	0.108	0.051	0.156
	A-D	0.182	0.054	0.072	0.153	0.052	0.098	0.045	0.153
	CVM	0.122	0.051	0.050	0.212	0.074	0.148	0.055	0.208
15	K-S	0.312	0.048	0.072	0.673	0.211	0.428	0.060	0.690
	A-D	0.389	0.046	0.100	0.813	0.262	0.605	0.065	0.823
	CVM	0.332	0.043	0.080	0.814	0.278	0.602	0.076	0.826
	$\chi^2$	0.136	0.037	0.044	0.707	0.169	0.480	0.060	0.717
25	K-S	0.472	0.045	0.086	0.928	0.387	0.763	0.084	0.942
	A-D	0.559	0.051	0.122	0.983	0.531	0.924	0.092	0.985
	CVM	0.511	0.049	0.099	0.980	0.527	0.907	0.098	0.982
	$\chi^2$	0.245	0.036	0.048	0.940	0.317	0.784	0.071	0.948

Level of Significance = .01

5	K-S	0.075	0.009	0.017	0.033	0.011	0.026	0.005	0.034
	A-D	0.096	0.012	0.026	0.014	0.004	0.012	0.003	0.011
	CVM	0.064	0.011	0.014	0.053	0.015	0.041	0.010	0.056
15	K-S	0.198	0.011	0.027	0.379	0.067	0.177	0.014	0.412
	A-D	0.261	0.009	0.038	0.578	0.081	0.317	0.012	0.603
	CVM	0.224	0.011	0.031	0.609	0.101	0.347	0.017	0.630
	$\chi^2$	0.078	0.015	0.013	0.576	0.094	0.346	0.026	0.582
25	K-S	0.341	0.010	0.040	0.794	0.167	0.479	0.021	0.807
	A-D	0.402	0.008	0.052	0.912	0.209	0.685	0.016	0.915
	CVM	0.377	0.011	0.046	0.922	0.246	0.706	0.023	0.924
	$\chi^2$	0.130	0.009	0.011	0.878	0.148	0.614	0.022	0.882

\* Key to Alternate Distributions:

Par.1 - Pareto (a=1, b=1, c=1)  
 Par.2 - Pareto (a=2, b=3, c=3.5)  
 Par.3 - Pareto (a=10, b=5, c=2)  
 Weibl - Weibull (shape = 3.5)

Gamma - Gamma (shape = 2)  
 Beta - Beta (F=2, Q=3)  
 Expon - Exponential (mean = 2)  
 Norml - Normal distribution

AD-A163 837

MODIFIED KOLMOGOROV-SHIRNOV ANDERSON-DARLING AND  
CRAMER-VON MISES TESTS F.. (U) AIR FORCE INST OF TECH  
WRIGHT-PATTERSON AFB OH SCHOOL OF ENGI.. J E PORTER  
DEC 85 AFIT/GSO/MA/85D-6 F/G 12/1

2/2

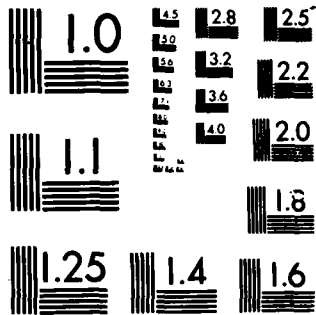
UNCLASSIFIED

NL

END

FILMED

GPO



MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A

### Linear Regression Tables

Tables XI and XII indicate the linear relationships existing between critical values and Pareto shape parameters. Table XI pertains to Kolmogorov-Smirnov critical values, while Table XII pertains to Cramer-von Mises critical values. No consistent linear relationship was identified for Anderson-Darling critical values.

The two tables contain linear coefficients and correlation values for each combination of sample sizes  $n = 10, 15, 20, 25,$  and  $30$  and levels of significance  $\alpha = .20, .15, .10, .05,$  and  $.01$ . No consistent linear relationship could be found for sample size  $n = 5$ . Further, the linear relationships apply only for values of the shape parameter  $c$  in the range  $1.5 \leq c \leq 4.0$ . Critical values for  $c < 1.5$  failed to display any consistent linear trend.

Each combination of sample size and significance level has its own linear coefficients and correlation value. In each case, the relationship between critical value  $Y$  and shape parameter  $c$  is given by the simple linear regression equation  $Y = b_0 + b_1c$  where  $b_0$  corresponds to the  $Y$ -axis intercept and  $b_1$  represents the slope of the described line. The correlation value  $R^2$  indicates the percent of total variation explained by the regression line. Thus,  $R^2$  is a measure of the strength of the linear relationship, with values near 1 indicating a strong linear tendency (13:250).

Table XI

COEFFICIENTS AND  $R^2$  VALUES OF THE RELATIONSHIP\*  
 BETWEEN KOLMOGOROV-SMIRNOV CRITICAL VALUES AND  
 PARETO SHAPE PARAMETERS FOR  $1.5 \leq c \leq 4.0$

n	Coeff	Level of Significance				
		.20	.15	.10	.05	.01
10	$b_0$	.2080	.2154	.2222	.2359	.2704
	$b_1$	.0057	.0067	.0090	.0117	.0144
	$R^2$	0.998	0.997	0.993	0.997	0.993
15	$b_0$	.1752	.1804	.1896	.2042	.2339
	$b_1$	.0051	.0065	.0074	.0091	.0117
	$R^2$	0.977	0.993	0.999	0.990	0.987
20	$b_0$	.1544	.1630	.1699	.1828	.2102
	$b_1$	.0044	.0042	.0054	.0068	.0091
	$R^2$	0.973	0.969	0.964	0.960	0.935
25	$b_0$	.1403	.1461	.1535	.1623	.1885
	$b_1$	.0038	.0043	.0050	.0075	.0091
	$R^2$	0.980	0.991	0.963	0.994	0.964
30	$b_0$	.1302	.1362	.1418	.1542	.1728
	$b_1$	.0030	.0034	.0047	.0053	.0090
	$R^2$	0.944	0.947	0.946	0.967	0.979

\* Relationship between K-S critical values Y  
 and Pareto shape parameter c is approximately

$$Y = b_0 + b_1 c \quad \text{where } 1.5 \leq c \leq 4.0$$

Table XII

COEFFICIENTS AND  $R^2$  VALUES OF THE RELATIONSHIP\*  
 BETWEEN CRAMER-VON MISES CRITICAL VALUES AND  
 PARETO SHAPE PARAMETERS FOR  $1.5 \leq c \leq 4.0$

n	Coeff	Level of Significance				
		.20	.15	.10	.05	.01
10	$b_0$	.0741	.0825	.0915	.1089	.1556
	$b_1$	.0045	.0050	.0067	.0095	.0137
	$R^2$	0.986	0.970	0.973	0.985	0.981
15	$b_0$	.0769	.0832	.0964	.1170	.1643
	$b_1$	.0053	.0069	.0083	.0106	.0178
	$R^2$	0.982	0.996	0.993	0.965	0.980
20	$b_0$	.0805	.0905	.1031	.1252	.1833
	$b_1$	.0047	.0051	.0065	.0089	.0135
	$R^2$	0.966	0.957	0.978	0.957	0.974
25	$b_0$	.0806	.0910	.1045	.1264	.1831
	$b_1$	.0055	.0059	.0072	.0102	.0166
	$R^2$	0.979	0.989	0.992	0.978	0.932
30	$b_0$	.0834	.0936	.1116	.1372	.1907
	$b_1$	.0047	.0055	.0054	.0074	.0161
	$R^2$	0.964	0.945	0.899	0.872	0.967

\* Relationship between C-VM critical values  $Y$   
 and Pareto shape parameter  $c$  is approximately

$$Y = b_0 + b_1 c \quad \text{where } 1.5 \leq c \leq 4.0$$

### Use of Tables

This section explains how to use the research results contained in Tables VI - XII.

Using Critical Value Tables. The critical values contained in Tables VI - VIII can be used to test whether a random data sample of size  $n = 5, 10, 15, 20, 25,$  or  $30$  follows a three-parameter Pareto distribution having specified shape parameter  $c = .5, 1, 1.5, 2, 2.5, 3, 3.5,$  or  $4$ . Given a random sample of observed data, the following steps outline basic elements of the procedure used in testing goodness-of-fit (13:357-367):

Step 1. Determine  $n$ , the number of observations contained in the random data sample.

Step 2. Identify the null and alternative hypotheses to be tested. In this case, the hypothesized shape parameter  $c$  must also be specified. Thus, the hypotheses are:

$H_0$ : The sample observations follow a Pareto distribution of specified shape  $c$ .

$H_1$ : At least one of the observations does not follow the Pareto of shape  $c$ .

Step 3. Determine the desired probability of committing a Type I error, i.e., the probability of erroneously rejecting the null hypothesis when  $H_0$  is true. This probability is the level of significance,  $\alpha$  (13:78).

Step 4. Order the  $n$  observations from smallest to largest.

Step 5. Assume  $H_0$  is true and estimate the unknown location and scale parameters using an invariant estimator. If the BLUE is selected as the estimator, and the sample size is small, the estimates can be computed manually from equations (34) and (35) for  $c = .5, 1, \text{ or } 2$ ; equations (37) to (39) for  $c = 1.5$ ; or equations (17), (18), (21), (22), and (29) for  $c = 2.5, 3, 3.5, \text{ or } 4$ . For larger sample sizes, or if several samples are involved, use the FORTRAN subroutines BXVALS, BLCLE2, and BLCGT2 in Appendix A.

Step 6. Use the estimates of location  $a$  and scale  $b$ , the hypothesized shape  $c$ , and the  $n$  ordered sample observations to compute the hypothesized Pareto CDF from equation (40). Subroutine HYPcdf in Appendix A can be used if manual calculations are not practical.

Step 7. Select the type of test to be performed and compute the corresponding test statistic. Use equation (42) for the modified Kolmogorov-Smirnov test, equation (43) for the modified Anderson-Darling test, or equation (44) for the modified Cramer-von Mises test. Subroutine TESTAT in Appendix A can be used to compute test statistics for all three tests.

Step 8. Identify the critical value from Table VI, VII, or VIII, based on test type, level of significance, sample size, and hypothesized shape parameter.

Step 9. Reject the null hypothesis if the value of the test statistic exceeds the critical value. If the test statistic does not exceed the critical value, conclude that there is insufficient evidence to reject the null hypothesis (13:76).

Using Power Comparison Tables. Tables IX and X can be used to draw conclusions regarding the relative ability of a test to correctly reject a false null hypothesis. This information can then be used to select the best test for a given situation. The higher the power value, the better are the chances against committing a Type II error because the probability of erroneously accepting a false null hypothesis is lessened (13:78).

Using Linear Regression Tables. Tables XI and XII can be used to estimate critical values for shape parameters which are not specifically listed in Tables VI and VIII, provided the hypothesized shape parameter  $c$  satisfies  $1.5 \leq c \leq 4.0$ . Given the sample size and specified level of significance, the linear slope and intercept values contained in Table XI can be substituted into the regression equation  $y = b_0 + b_1 c$  to find the Kolmogorov-Smirnov critical value  $y$ . If the Cramer-von Mises test is involved, the values should be taken from Table XII.

### Example

Suppose a maintenance unit wants to model the failure rate of a certain equipment component. Based on 10 independent random samples, the unit observes the following failure times of the component (expressed in months following initial use): 1.178, 1.127, 1.373, 1.068, 1.059, 1.010, 1.474, 4.830, 3.997, 1.799. The unit desires to test the hypothesis that the component failure times follow the Pareto distribution with shape  $c = 2.5$ . One specified requirement is that the test be designed so that the probability of erroneously rejecting a true null hypothesis must not exceed .05.

Since there are 10 random observations in the data sample,  $n = 10$  for this example. The required level of significance is  $\alpha = .05$ . The hypotheses are:

$H_0$ : The observed failure times follow the Pareto distribution of shape  $c = 2.5$ .

$H_1$ : At least one of the observations does not follow the Pareto of shape 2.5.

The next step is to arrange the random sample in ascending order: 1.010, 1.059, 1.068, 1.127, 1.178, 1.373, 1.474, 1.799, 3.997, 4.830. These values are input into subroutine BXVALS which yields  $B_i$  values of .920, .838, .754, .668, .579, .486, .389, .285, .171, and .034. These values are then input into subroutine BLCGT2, which computes the parameter estimates  $\hat{a} = .963$  and  $\hat{b} = 1.128$ . Subroutine

HYPCDF is then used to compute 10 values of the hypothesized Pareto CDF: .097, .185, .201, .288, .354, .539, .608, .750, .962, and .976.

The values of  $n$ ,  $c$ , and the hypothesized Pareto CDF are input into subroutine TESTAT, which computes the test statistics  $K-S = .162$ ,  $A-D = .416$ , and  $C-VM = .042$ . From Table VI, the  $K-S$  critical value for  $\alpha = .05$ ,  $n = 10$ , and  $c = 2.5$  is .265. Since the test statistic does not exceed the critical value, there is insufficient evidence to reject the null hypothesis. The same conclusion is reached from the  $A-D$  and  $C-VM$  critical values (Tables VII and VIII).

Now suppose the unit wants to test the null hypothesis that a set of  $n = 25$  observed service times follows the Pareto distribution of shape  $c = 3.35$ . The analyst computes the  $K-S$  or  $C-VM$  test statistic as before, but the critical values are not listed for  $c = 3.35$ . Therefore, the next step is to determine the appropriate regression coefficients from Table XI or XII. For  $n = 25$  and  $\alpha = .05$  the  $K-S$  coefficients are  $b_0 = .1623$  and  $b_1 = .0075$ . The  $K-S$  critical value is  $Y = b_0 + b_1 c = .1623 + .0075 (3.35) = .1874$ .

#### Chapter Summary

This chapter presented the results of the research conducted in response to the three objectives listed in Chapter I. Tables of critical values for the modified  $K-S$ ,  $A-D$ , and  $C-VM$  tests were presented. Also included were

tables comparing powers of the K-S, A-D, and C-VM statistics against the Chi-square. Tables of regression coefficients were presented as well. The use of the tables was explained, and an example was described.

The research results are further analysed and discussed in the next chapter.

## VI. ANALYSIS AND DISCUSSION

### Chapter Overview

This chapter discusses the results presented in Chapter V. Observations are made concerning the tables of critical values, power comparisons, and regression coefficients, including an explanation as to how the computer programs were verified and validated.

### Critical Values

The critical value tables generated for this thesis are located in Chapter V. For the K-S test (Table VI), the critical values for a given level of significance and shape parameter decrease as the sample size increases. Further, the size of the decrease becomes smaller at larger values of  $n$ . This trend suggests that the K-S critical values may converge to a lower limit as the sample size increases. However, the use of sample sizes larger than 30 would have required much more computer processing time, and thus was beyond the scope of this thesis. The A-D critical values (Table VII) exhibit a different pattern. The values for each combination of significance level and shape parameter generally decrease from  $n = 5$  to 20 and increase from  $n = 20$  to 30, suggesting a convergence between 15 and 20. Similarly, the C-VM critical values (Table VIII) appear to converge between  $n = 25$  and 30,

since the values consistently decrease until  $n = 30$ , then begin to increase.

An important observation is made when the table of modified K-S values is compared to a standard (unmodified) K-S table (13:462). For each value of  $n$  in Table VI, the critical values for shape 1 or 2 at a .05 significance level are nearly the same as the critical values for a .20 significance level using the standard table. Thus the result of using the standard K-S table when location and scale parameters are estimated would be to obtain an extremely conservative test in the sense that the actual significance level would be much lower than that given by the standard table.

#### Power Comparison

The power comparison tables generated for this thesis are located in Chapter V. Values in Table IX pertain to a null hypothesis for which the Pareto shape parameter is 1.0, whereas in Table X the hypothesized shape parameter is 3.5. Both tables are divided into two sections based on a level of significance of .05 or .01. It is obvious from the tables that none of the three tests developed in this thesis is very powerful when the sample size is only five. Nevertheless, they at least provide some means of testing goodness-of-fit for sample sizes which are too small to use the Chi-square test. For sample sizes of 15 or 25, the powers improve dramatically.

For each alternative distribution the three tests tended to be more powerful than the Chi-square. Two sets of Chi-square critical values were examined. The first set of values was taken from a standard table of Chi-square critical values corresponding to 2 degrees of freedom (13:432). After completing 5000 Monte Carlo repetitions, it was discovered that the tabled Chi-square value for a level of significance of .05 displayed a probability of a Type I error (i.e., rejecting  $H_0$  when true) of .10, which was twice the claimed level of significance of .05. Similarly, the probability of Type I error for a claimed level of significance .01 was, in fact, .02. This discrepancy was due to the fact that the tabled Chi-square values represent only an approximation of the actual asymptotic distribution of the Chi-square, so that the actual value lies somewhere between Chi-square with 2 degrees of freedom and Chi-square with 4 degrees of freedom (34:401-402). Since the Type I errors were twice their expected value, a second set of Chi-square critical values was generated using Monte Carlo simulation in the same manner as was used to generate critical values for the K-S, A-D, and C-VM tests. As apparent from Tables IX and X, the second set of Chi-square values display Type I error rates which is much closer to the claimed level of significance of .05 or .01. Therefore, these values were used in the power comparison tables rather than the less accurate values stemming from the standard Chi-square table.

The modified K-S, A-D, and C-VM tests are especially powerful when the sample data are taken from the Weibull, the Beta, or the normal distribution. On the other hand, the three tests display relatively low power in their ability to distinguish against the exponential or the Pareto with different shape parameters. In general, the K-S test has higher power than the others when the hypothesized shape parameter is 1.0. When the shape parameter is 3.5, the C-VM test tends to be more powerful. Next to the Chi-square, the A-D test appears to have the lowest power in most cases.

#### Regression Analysis

The regression tables generated for this thesis are also located in Chapter V. Table XI contains regression coefficients and correlation values for the modified Kolmogorov-Smirnov test, while Table XII contains regression values for the Cramer-von Mises test.

It is apparent from Tables VI and VIII that for a given significance level and sample size except  $n = 5$ , the K-S and C-VM critical values decrease from shape parameter 0.5 to 1.5, then steadily increase for shapes 1.5 to 4.0. Using the method of least squares, a simple linear regression analysis was performed on the critical values. The correlation of regression on the shape parameter interval 0.5 to 1.5 was in most cases less than .80. However, the regression relationships on the shape interval 1.5 to 4.0 showed very

strong correlation (.97 or higher in most cases). Therefore, regression coefficients corresponding to the interval  $1.5 \leq c \leq 4.0$  were included in Tables XI and XII.

No consistent linear trend could be identified for the Anderson-Darling critical values. In general the values seem to decrease on the interval  $0.5 \leq c \leq 2.5$  and then increase on  $2.5 \leq c \leq 4.0$ . However, when least squares regression was applied to the two intervals, the correlation values tended to be less than .80 in most cases. Therefore, it was decided not to include a regression table for the A-D test.

#### Verification and Validation

The critical values were computed by the CRITVAL program and associated subroutines contained in Appendix A. The power study was conducted using the POWER program and subroutines in Appendix B. The purpose of verification was to ensure that the concepts and equations developed in this thesis were reflected accurately in the computer code. The five verification techniques suggested by Banks and Carson (5:379) were implemented as follows:

1. Have the code checked. The code was checked by two individuals knowledgeable of FORTRAN programming. One of the individuals, Charek, was also very familiar with the logic required for computing parameter estimates for the Pareto distribution, since he too has conducted extensive research in this area (12).

2. Make a flow diagram. Flow diagrams illustrating the logic involved in generating critical values served as the basis of the program and were closely followed during the actual writing of the program. The diagrams are included in Appendices A and B.

3. Examine a wide variety of output. The output of each subroutine and the results of each individual computation was checked through extensive use of print statements. Each computational stage was checked at least once against manual calculations to ensure the expected values were produced. A pre-production run involving 50 replications was thoroughly examined for reasonableness prior to the final production run of 5000 replications.

4. Print the input parameters. During the test runs, input parameters were printed before and after each calculation to ensure against any inadvertant alteration of parameters.

5. Make the code self-documenting. Extensive comments have been incorporated into the programs and subroutines to allow easy interpretation of the logic. At the beginning of each program component, every variable is defined and the purpose explained.

Validation of the computer programs was provided in the results of the power study. For each hypothesized shape parameter and sample size, the K-S, A-D, and C-VM tests

displayed a Type I error rate equal to or very near the claimed level of significance. This fact validates the critical values as well as the power comparison values.

#### Chapter Summary

The results of this thesis are presented in Tables VI-XII. The results of the power study show that the three tests developed in this thesis offer tests which can be used with small sample sizes and are more powerful than the Chi-square at larger sample sizes. The programs used to generate the tables were thoroughly verified and validated.

Conclusions and recommendations for further study are presented in the next chapter.

## VII. CONCLUSIONS AND RECOMMENDATIONS

### Conclusions

The following conclusions are based on the results contained in this thesis:

1. The first research objective listed in Chapter I has been successfully fulfilled. Tables VI-VIII contain critical values of the modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) tests. The validity of these critical values has been verified by a Monte Carlo power study which has shown that all three tests achieve the claimed level of significance when the null hypothesis is true. Therefore, each table of critical values can be used to test whether a random sample of data follows the three-parameter Pareto distribution with specified shape parameter.

2. The second research objective has also been completed successfully. The results of the power study are contained in Tables IX and X. It appears that none of the three tests developed in this thesis is very powerful when the sample size is only five. For sample sizes of 15 or 25, however, the powers improve dramatically. For each of the alternative distributions considered, the three tests tended to be more powerful than the Chi-square, as expected. The three tests are especially powerful when the sample data are

taken from the Weibull, the Beta, or the normal distribution. In general, the K-S test has higher power than the others when the hypothesized shape parameter is 1.0. When the shape parameter is 3.5, the C-VM test tends to be more powerful. Next to the Chi-square, the A-D test appears to have the lowest power in most cases.

3. Successful completion of the third research objective has revealed a strong linear relationship between shape parameters and critical values for the K-S and C-VM tests. Linear coefficients and correlation values are contained in Tables XI and XII. However, no consistent functional relationship could be identified for the A-D test.

#### Recommendations

Based on observations made during the investigation for this thesis, the following research areas are proposed for further study:

1. Apply the techniques used in this thesis to generate modified K-S, A-D, and C-VM tests for other distribution functions.

2. Investigate whether other types of goodness-of-fit tests can be modified through Monte Carlo techniques. For example, if the S statistic of Mann, Scheuer, and Fertig (38) can be modified for the Pareto distribution, a power study can be conducted to determine whether the S statistic is more powerful than the K-S, A-D, or C-VM tests.

3. Derive the maximum likelihood estimators of location and scale for the three-parameter Pareto.

4. Compute critical values for sample sizes and Pareto shape parameters not specifically included in Tables VI-VIII. For example, the tables can be expanded to include all sample sizes from 3 to 100 and shape parameters from 0.25 to 10.

5. Increase the accuracy of the critical values by using various techniques (5:406-442) of experimental design (e.g., increased repetitions, multiple batch runs, replications, antithetic random number seeds, analysis of variance, etc.) to reduce the inherent uncertainty and to determine the amount of variance involved.

6. Apply more sophisticated regression techniques to determine the functional relationship between Pareto shape parameters and Anderson-Darling critical values.

7. Apply the results of this thesis to earlier studies (Chapter III) involving the Pareto distribution. For example, Berger and Mandelbrot's (7) conclusion that the Pareto can be used to model errors in communications circuits can now be tested for goodness-of-fit.

8. Further investigate potential applications of the Pareto distribution as an accurate model of actual phenomena. The tests developed in this thesis contribute to the usefulness of the Pareto distribution which, in many situations, should be considered as a viable model when simulating or testing the underlying distribution of a given population.

**APPENDIX A**  
**Computer Program and Subroutines**  
**for Generating Critical Values**

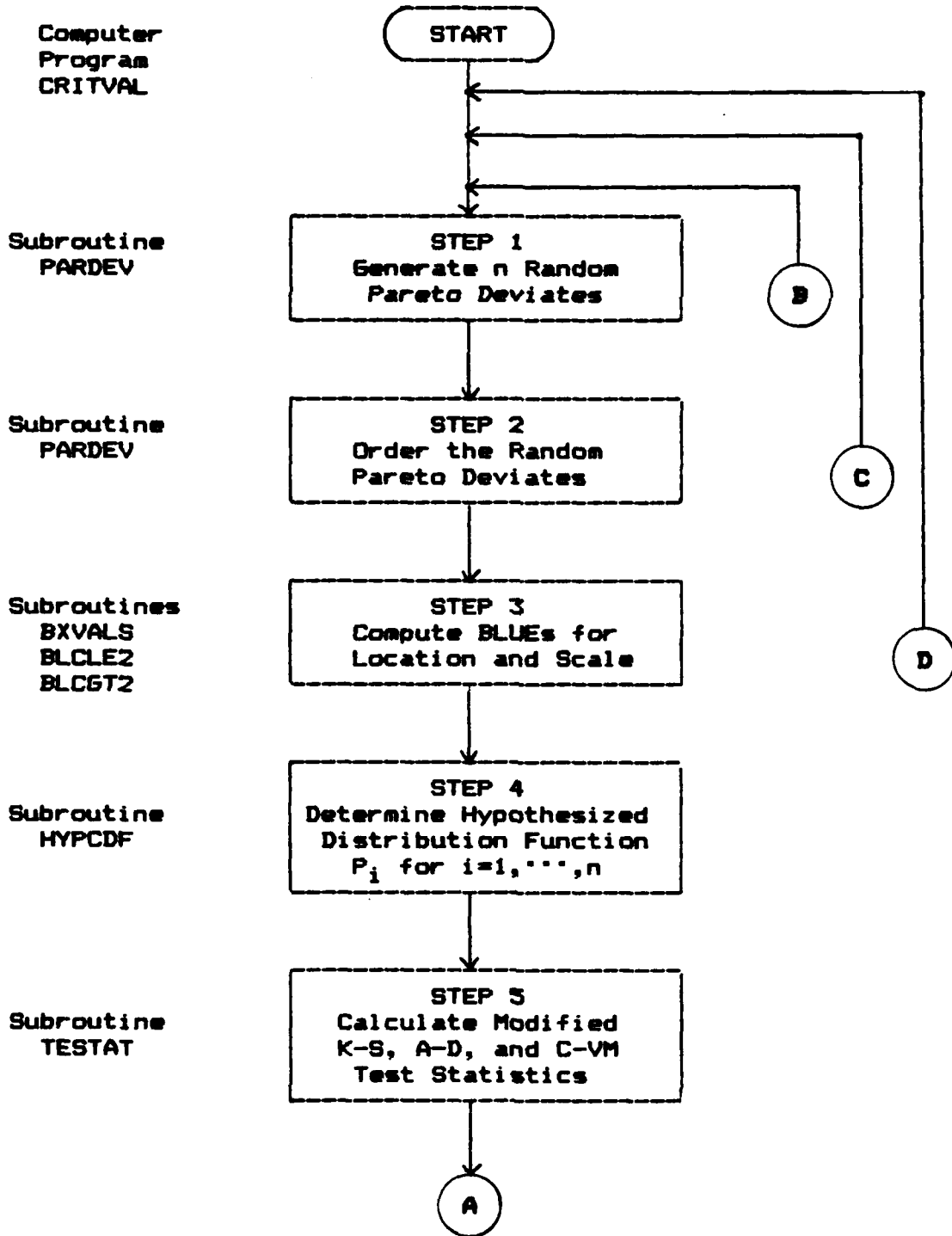
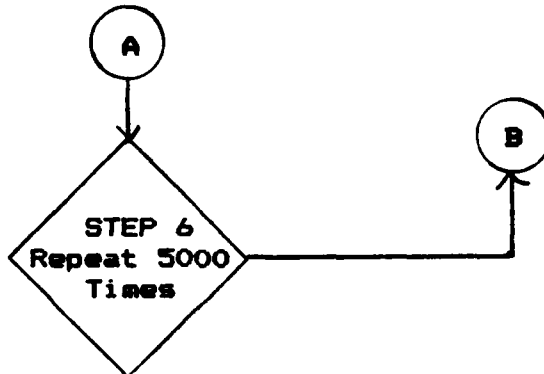
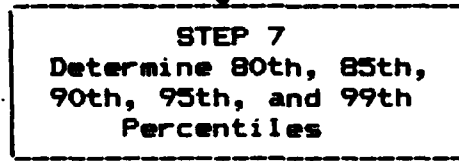


Fig 6. Procedure for Generating Critical Values

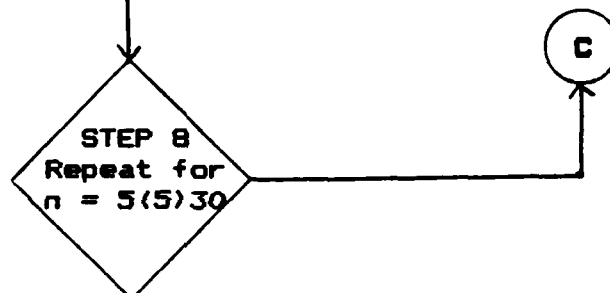
Main Program  
DO Loop 60



Subroutine  
CRTVAL



Main Program  
DO Loop 80



Main Program  
DO Loop 90

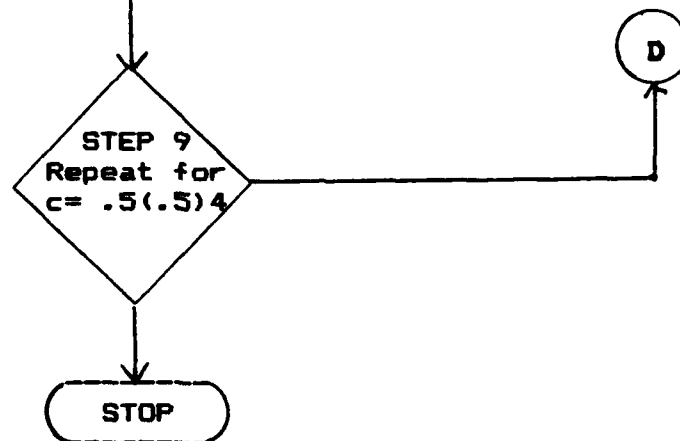


Fig 6 (Continued). Procedure for Generating Critical Values

```

1      c***** Classroom Support Computer (CSC) - VAX 11/785 - VMS 4.1 ****
2      c
3      c***** CRITVAL PROGRAM FOR PARETO GOODNESS-OF-FIT TESTS *****
4      c
5      c*****
6      c*****
7      c**
8      c** BEGIN CRITVAL MAIN PROGRAM **
9      c**
10     c*****
11     c
12     c Ref: Appendix A, Figure 6.
13     c
14     c=====
15     c
16     c Purpose:
17     c     1. Generate critical value tables for the modified
18     c     Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D),
19     c     and Cramer-von Mises (C-VM) tests for the three-
20     c     parameter Pareto distribution when location and
21     c     scale parameters must be estimated from sample data.
22     c     2. Provide extensive commentary to help novice prog-
23     c     rammers develop similar goodness-of-fit programs.
24     c     Thus, diagnostic print routines have been retained as
25     c     part of the commentary rather than deleted.
26     c
27     c=====
28     c
29     c Variables:
30     c     dseed = random number seed
31     c     c = shape parameter
32     c     n = sample size
33     c     nshp = shape parameter counter (8 different values)
34     c     nsiz = sample size counter (6 different values of n)
35     c     npct = percentile counter (5 different percentiles)
36     c     nst = number of test statistics to be used
37     c     it = iteration counter (5000 repetitions required)
38     c     KS = array of values of modified K-S test statistic
39     c     CVM = array of values of modified C-VM test statistic
40     c     AD = array of values of modified A-D test statistic
41     c     alpha = level of significance
42     c
43     c=====
44     c
45     c Input:
46     c     nst = number of repetitions (input at computer terminal)
47     c     dseed = random number seed (input at computer terminal)
48     c
49     c=====

```

```

50      c
51      c Subroutines:
52      c
53      c PARDEV - Generates n ordered Pareto deviates
54      c BXVALS - Calculates B values and summations of B and Bx
55      c BLCLE2 - Finds BLUEs for location and scale when c <= 2
56      c BLCGT2 - Finds BLUEs for location and scale when c > 2
57      c HYPCDF - Computes the Hypothesized Pareto CDF
58      c TESTAT - Calculates the K-S, A-D, and C-VM test statistics
59      c CRTVAL - Determines critical values from plotting positions
60      c
61      c=====
62      c
63      c Calculate:
64      c
65      c     nc = n * c
66      c
67      c     Plotting Positions (Eqn 51):
68      c
69      c     Y(i) = (i - 0.3)/(nst + 0.4) for i = 1,...,nst(=5000)
70      c
71      c=====
72      c
73      c Output:
74      c
75      c     KScrit = 3-D array of critical values for modified K-S test
76      c     ADcrit = 3-D array of critical values for modified A-D test
77      c     CVcrit = 3-D array of critical values for modified C-VM test
78      c
79      c=====
80      c
81      c Declare Variables:
82      c
83      c     common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
84      c     1     Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
85      c     1     KScrit,ADcrit,CVcrit,Y
86      c     integer n,nsiz,nshp,it,npct,nst
87      c     real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),
88      c     1     CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
89      c     1     KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),r(30),pct,
90      c     1     Y(5002),alpha
91      c     double precision dseed
92      c
93      c     ** Open Output Files to Store Computed Critical Values: **
94      c     open (unit=7,file='CRIT',status='new')
95      c
96      c     ** Number of Test Statistics to be Used on Each Run: **
97      c     print*, 'The Monte Carlo analysis will require'
98      c     print*, '     5000 test statistics.'
99      c     print*, 'Enter the number to be used for this run:'
100     c     read*,nst
101     c

```

```

102      c      ** Calculate 5002 Plotting Positions on the Y-axis: **
103      c
104          Y(0) = 0.0
105      do 10 i = 1,nst
106          Y(i) = (i - 0.3)/(nst + 0.4)
107      10 continue
108          Y(nst + 1) = 1.0
109      c
110          print*, ' '
111          print*, 'SELECTED MEDIAN RANKS PLOTTING POSITIONS'
112          print*, ' TO BE USED TO FIND CRITICAL VALUES:'
113          print*, ' '
114          print*, '      Y(5001) = ',Y(5001)
115          print*, '      Y(5000) = ',Y(5000)
116          print*, '99PCT: Y(4950) = ',Y(4950)
117          print*, '95PCT: Y(4750) = ',Y(4750)
118          print*, '90PCT: Y(4500) = ',Y(4500)
119          print*, '85PCT: Y(4250) = ',Y(4250)
120          print*, '80PCT: Y(4000) = ',Y(4000)
121          print*, '      Y(0001) = ',Y(1)
122          print*, '      Y(0000) = ',Y(0)
123          print*, '===== '
124      c
125      c      ** Plotting Positions Computation Complete      **
126      c
127          print*, 'Enter random number seed or "1." for default:'
128          read*,dseed
129          if (dseed .eq. 1.) dseed = 123457.d00
130          print*, ' '
131          print*, 'STANDBY . . . COMPUTATIONS IN PROGRESS'
132      c
133          nshp = 0
134      c
135      c --- Begin DO Loop 90 for Shape Parameter Values c = .5(.5)4 ---
136      c
137          do 90 shape = 0.5,4.0,.5
138              c = shape
139              nshp = nshp + 1
140      c
141      c      Write Headings for Output Data:
142          write(7,52)
143          write(7,51)
144          write(7,52)
145          write(7,54)
146          write(7,52)
147          write(7,56)
148      c
149          nsiz = 0
150      c
151      c --- Begin DO Loop 80 for Sample Sizes n = 5(5)30 ---
152      c
153          do 80 nsamp = 5,30,5

```

```

154
155         if ( (c.eq.0.5) .and. (nsamp.eq.5) ) then
156     c         the BLUEs do not exist, so we must let:
157         n = 6
158         else
159             n = nsamp
160         end if
161     c
162         nsiz = nsiz + 1
163         nc = n * c
164     c
165         write(7,58)
166     c
167     c         --- Begin DO Loop 60 for 5000 Iterations ---
168     c
169         do 60 it = 1,nsi
170     c
171     c         ** Perform Steps 1 & 2 of Fig 6: **
172     c
173             call PARDEV
174     c
175     c         ** Perform Step 3 of Figure 6: **
176     c
177             call BXVALS
178     c
179             if (c .le. 2.0) then
180                 call BLCLE2
181             else
182                 call BLCGT2
183             end if
184     c
185     c         ** Perform Step 4 of Figure 6: **
186     c
187             call HYPCDF
188     c
189     c         ** Perform Step 5 of Figure 6: **
190     c
191             call TESTAT
192     c
193     c         60 continue
194     c
195     c         --- End DO Loop 60 for 5000 Iterations ---
196     c         ** Completes Step 6 of Figure 6 **
197     c
198     c         ** Perform Step 7 of Figure 6: **
199     c
200     c         --- Begin DO Loop 70 for Percentiles ---
201     c
202         do 70 npct = 1,5
203     c
204     c
205         call CRTVAL

```

```

206      c
207      c          -- Write CRTVAL Output to File --
208      c          write(7,62),1,-pct,n,c,KScrit(nsiz,nshp,npct),
209      c          1          ADcrit(nsiz,nshp,npct),CVcrit(nsiz,nshp,npct)
210      c
211      c          print*,' '
212      c          print*,' CRITICAL VALUES FROM MAIN PROGRAM'
213      c          print*,' pct =',pct,' n=',n,' ** c=',c
214      c          print*,' K-S =',KScrit(nsiz,nshp,npct),
215      c          1          ' A-D =',ADcrit(nsiz,nshp,npct),
216      c          1          ' CVM =',CVcrit(nsiz,nshp,npct)
217      c          print*,' '
218
219      c          70          continue
220      c
221      c          ---      End DO Loop 70 for Percentiles      ---
222      c
223      c          80          continue
224      c
225      c          ---      End DO Loop 80 for Sample Sizes n = 5(5)30      ---
226      c          ** Completes Step 8 of Figure 6 **
227      c
228      c          90          continue
229      c
230      c          ---      End DO Loop 90 for Shape Parameter Values c = .5(.5)4      ---
231      c          ** Completes Step 9 of Figure 6 **
232      c
233      c*****
234      c
235      c  OUTPUT INSTRUCTIONS:  The remainder of the main program
236      c  consists of commands to format the output data and write
237      c  the data and headers to a file which can be printed out
238      c  in hardcopy.
239      c
240      c*****
241      c
242      c *** Write KS Critical Value Tables to File by Alpha Level: ***
243      c
244      c          write(7,52)
245      c          write(7,130)
246      c          write(7,52)
247      c          write(7,132)
248      c          write(7,52)
249      c          write(7,200)
250      c          write(7,201)
251      c          write(7,52)
252      c
253      c          npct = 0
254      c
255      c ---Begin DO Loop 105 to Sort Critical Values by Alpha Level---
256      c
257      c          do 105 npct = 1,5

```

```

258      c
259      if (npct .ne. 5) alpha = .25 - (.05*npct)
260      if (npct .eq. 5) alpha = .01
261      c
262      nsiz = 0
263      n = 0
264      c
265      --- Begin DO Loop 107 to Sort Output by Sample Size ---
266      c
267      do 107 nsiz = 1,6
268      c
269      n = 5 * nsiz
270
271      Write(7,120),alpha,n,KScrit(nsiz,1,npct),KScrit
272      1      (nsiz,2,npct),KScrit(nsiz,3,npct),KScrit(nsiz,
273      1      4,npct),KScrit(nsiz,5,npct),KScrit(nsiz,6,npct),
274      1      KScrit(nsiz,7,npct),KScrit(nsiz,8,npct)
275      c
276      107      continue
277      c
278      --- End DO Loop 107 After Sorting by Sample Size ---
279      c
280      write(7,201)
281      c
282      105 continue
283      c
284      --- End DO Loop 105 After Sorting Output by Alpha Level ---
285      c
286      c
287      c *** Write AD Critical Value Tables to File by Alpha Level: ***
288      c
289      write(7,52)
290      write(7,140)
291      write(7,52)
292      write(7,142)
293      write(7,52)
294      write(7,200)
295      write(7,201)
296      write(7,52)
297      c
298      npct = 0
299      c
300      ---Begin DO Loop 115 to Sort Critical Values by Alpha Level---
301      c
302      do 115 npct = 1,5
303      c
304      if (npct .ne. 5) alpha = .25 - (.05*npct)
305      if (npct .eq. 5) alpha = .01
306      c
307      nsiz = 0
308      n = 0
309      c

```

```

310      c      --- Begin DO Loop 117 to Sort Output by Sample Size ---
311      c
312      c      do 117 nsiz = 1,6
313      c
314      c      n = 5 * nsiz
315      c
316      c      Write(7,120),alpha,n,ADcrit(nsiz,1,npct),ADcrit
317      1      (nsiz,2,npct),ADcrit(nsiz,3,npct),ADcrit(nsiz,
318      1      4,npct),ADcrit(nsiz,5,npct),ADcrit(nsiz,6,npct),
319      1      ADcrit(nsiz,7,npct),ADcrit(nsiz,8,npct)
320      c
321      c      117 continue
322      c
323      c      --- End DO Loop 117 After Sorting by Sample Size ---
324      c
325      c      write(7,201)
326      c
327      c      115 continue
328      c
329      c      --- End DO Loop 115 After Sorting Output by Alpha Level ---
330      c
331      c
332      c      *** Write CVM Critical Value Tables to File by Alpha Level ***
333      c
334      c      write(7,52)
335      c      write(7,150)
336      c      write(7,52)
337      c      write(7,152)
338      c      write(7,52)
339      c      write(7,200)
340      c      write(7,201)
341      c      write(7,52)
342      c
343      c      npct = 0
344      c
345      c      ---Begin DO Loop 125 to Sort Critical Values by Alpha Level---
346      c
347      c      do 125 npct = 1,5
348      c
349      c      if (npct .ne. 5) alpha = .25 - (.05*npct)
350      c      if (npct .eq. 5) alpha = .01
351      c
352      c      nsiz = 0
353      c      n = 0
354      c
355      c      --- Begin DO Loop 127 to Sort Output by Sample Size ---
356      c
357      c      do 127 nsiz = 1,6
358      c
359      c      n = 5 * nsiz
360

```

```

361          Write(7,120),alpha,n,CVcrit(nsiz,1,npct),CVcrit
362          1      (nsiz,2,npct),CVcrit(nsiz,3,npct),CVcrit(nsiz,
363          1      4,npct),CVcrit(nsiz,5,npct),CVcrit(nsiz,6,npct),
364          1      CVcrit(nsiz,7,npct),CVcrit(nsiz,8,npct)
365      c
366      127      continue
367      c
368      c      --- End DO Loop 127 After Sorting by Sample Size ---
369      c
370          write(7,201)
371      c
372      125      continue
373      c
374      c      --- End DO Loop 125 After Sorting Output by Alpha Level ---
375      c
376      c      Specify Format for Hardcopy Output Data and Headers:
377      c
378          51      format(' *****')
379          52      format(' ')
380          54      format(' ** PARETO CRITICAL VALUES FOR SHAPE C = **')
381          56      format(' alpha',3X,'n',4X,'c',7X,'KS',8X,'AD',8X,'CVM')
382          58      format(' -----')
383          62      format(' ',T3,F3.2,I5,F6.1,3F10.4)
384          120     format(' ',T3,F3.2,I5,F8.3,7F9.3)
385          130     format('1',36X,'Table VI')
386          132     format(20X,'CRITICAL VALUES FOR THE MODIFIED K-S TEST')
387          140     format('1',36X,'Table VII')
388          142     format(20X,'CRITICAL VALUES FOR THE MODIFIED A-D TEST')
389          150     format('1',35X,'Table VIII')
390          152     format(19X,'CRITICAL VALUES FOR THE MODIFIED C-VM TEST')
391          200     format(' alpha',3X,'n',4X,'c=.5',5X,'1.0',6X,'1.5',6X,
392          1      '2.0',6X,'2.5',6X,'3.0',6X,'3.5',6X,'4.0')
393          201     format(81(' '))
394      c
395          close(7)
396      c
397          end
398      c
399      c=====
400      c      END MAIN PROGRAM
401      c*****

```

```

402           Subroutine PARDEV
403 c*****
404 c**
405 c**           B E G I N   S U B R O U T I N E   P A R D E V           **
406 c**
407 c*****
408 c
409 c Ref: Appendix A, Fig 6, Steps 1 & 2.
410 c
411 c=====
412 c
413 c Purpose: For a specified sample size n, generate n random
414 c           deviates from a Pareto distribution with location and
415 c           scale parameters set to one (a = b = 1) and the shape
416 c           parameter c set to some specified positive value.
417 c
418 c=====
419 c
420 c Variables:
421 c           r = array containing n random numbers
422 c           c = shape parameter
423 c           x = array containing n Pareto deviates
424 c           n = sample size
425 c           dseed = random number seed
426 c
427 c=====
428 c
429 c Input:   dseed = random number seed (from MAIN program)
430 c           c = shape parameter = .5(.5)4 (MAIN DO Loop 90)
431 c           n = sample size = 5(5)30 (MAIN DO Loop 80)
432 c
433 c=====
434 c
435 c IMSL Subroutines:
436 c
437 c GGUBS - generates random numbers uniformly distributed on (0,1)
438 c VSRTA - arranges a set of numbers in ascending order
439 c
440 c=====
441 c
442 c Calculate:
443 c
444 c   x(j) = (1/r(j)) ** (1/c)   for j = 1,2,...,n (from eqn 48)
445 c
446 c=====
447 c
448 c Output:   x = array of n ordered Pareto deviates
449 c
450 c=====
451 c
452 c Declare Variables:
453 c

```

```

454      common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
455      1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
456      1      KScrit,ADcrit,CVcrit,Y
457      real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),
458      1      CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
459      1      r(30),KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),
460      1      Y(5002)
461      integer n,npct
462      double precision dseed
463      c
464      c--- Begin DO Loop 10 to Generate n Random Pareto Deviates ---
465      c
466      do 10 j = 1,n
467      c
468      c      Use IMSL subroutine to generate random numbers:
469      c      call GGUBS(dseed,n,r)
470      c
471      c      Use eqn 48 to transform them to Pareto deviates:
472      c      x(j) = (1.0/r(j))**(1.0/c)
473      c
474      10  continue
475      c
476      c--- End DO Loop 10 after Generating n Random Deviates ---
477      c      ** (Completes Step 1 of Figure 6) **
478      c
479      c Use IMSL subroutine to place the deviates in ascending order:
480      c call vsrta(x,n)
481      c      ** (Completes Step 2 of Figure 6) **
482      c
483      c      return
484      c      end
485      c
486      c=====
487      c      END SUBROUTINE PARDEV
488      c*****

```

```

489           Subroutine BXVALS
490 c*****
491 c**
492 c**      BEGIN SUBROUTINE BXVALS      **
493 c**
494 c*****
495 c
496 c Ref: Appendix A, Fig. 6, Step 3.
497 c
498 c=====
499 c
500 c Purpose: For a given sample size n, calculate the B values
501 c          used to find the BLUEs of location and scale. Also
502 c          find the sum of the first n-1 values of B(i). Then,
503 c          compute the three values equal to the sums of the
504 c          first n-1, the first n-2, and (for c = .5, 1, or 2)
505 c          the first n -2/c values of B(i)*x(i).
506 c
507 c=====
508 c
509 c Variables:  c = shape parameter
510 c             n = sample size
511 c             x = array containing n ordered Pareto deviates
512 c             B = array containing n values of B
513 c             Bsum1 = sum of B(i) values for i = 1,2,...,(n-1)
514 c             Bxsum1 = sum of B(i)*x(i) for i = 1,2,...,(n-1)
515 c             Bxsum2 = sum of B(i)*x(i) for i = 1,2,...,(n-2)
516 c             Bxsm2c = sum of B(i)*x(i) for i = 1,2,...,(n-2/c)
517 c
518 c=====
519 c
520 c Input:     c = shape parameter = .5(.5)4 (from MAIN DO Loop 90)
521 c           n = sample size = 5(5)30 (from MAIN DO Loop 80)
522 c           nc = n*c (from MAIN program)
523 c           x = ordered Pareto deviates (from PARDEV)
524 c
525 c=====
526 c
527 c Calculate:
528 c
529 c           B(i) = [1 - 2/c(n-i+1)] * B(i-1) (eqn 29)
530 c
531 c           Bsum1 = B(1) + B(2) + ... + B(n-1)
532 c
533 c           Bxsum1 = B(1)*x(1) + ... + B(n-1)*x(n-1)
534 c
535 c           Bxsum2 = B(1)*x(1) + ... + B(n-2)*x(n-2)
536 c
537 c           Bxsm2c = B(1)*x(1) + ... + B(n-2/c)*x(n-2/c)
538 c
539 c=====
540 c

```

```

541 c Output:
542 c      B = array containing n values of B
543 c      Bsum1 = sum of first (n-1) B values
544 c      Bxsum1 = sum of first (n-1) B*x values
545 c      Bxsum2 = sum of first (n-2) B*x values
546 c      Bxsm2c = sum of first (n-2/c) B*x (if 2/c is integer)
547 c
548 c=====
549 c
550 c Declare Variables:
551 c
552 c      common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
553 c      1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
554 c      1      KScrit,ADcrit,CVcrit,Y
555 c      real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),
556 c      1      CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
557 c      1      KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),Y(5002)
558 c      integer n
559 c      double precision dseed
560 c
561 c Calculate the first B value (eqn 25):
562 c
563 c      B(1) = 1.0 - 2.0/nc
564 c
565 c --- Begin DO Loop 10 to Find the 2nd thru nth B values ---
566 c
567 c      do 10 j = 2,n
568 c          B(j) = B(j-1) * (1.0 - (2.0/(c*(n-j+1))))
569 c      10 continue
570 c
571 c --- End DO Loop 10 ---
572 c
573 c      Bsum1 = 0
574 c
575 c --- Begin DO Loop 20 to Sum the First n-1 Values of B ---
576 c
577 c      do 20 k=1,(n-1)
578 c          Bsum1 = Bsum1 + B(k)
579 c      20 continue
580 c
581 c --- End DO Loop 20 ---
582 c
583 c      Bxsum1 = 0
584 c
585 c --- Begin DO Loop 30 to Sum the First n-1 Values of Bx ---
586 c
587 c      do 30 l=1,(n-1)
588 c          Bxsum1 = Bxsum1 + (B(l))*x(l)
589 c      30 continue
590 c
591 c --- End DO Loop 30 ---
592 c

```

```

593         Bxsum2 = Bxsum1 - (B(n-1)*x(n-1))
594     c
595     c --- Find Bxsm2c When 2/c is an Integer (c=.5, 1, or 2) ---
596     c
597         Bxsm2c = 0
598     c
599         if (c .eq. 1.0) then
600             Bxsm2c = Bxsum2
601         else if (c .eq. 2.0) then
602             Bxsm2c = Bxsum1
603         else if (c .eq. 0.5) then
604             Bxsm2c = Bxsum2 - (B(n-3)*x(n-3)) - (B(n-2)*x(n-2))
605         end if
606     c
607         return
608     end
609     c
610     c=====
611     c                      END SUBROUTINE BXVALS
612     c*****

```

```

613           Subroutine BLCLE2
614 c*****
615 c**
616           BEGIN SUBROUTINE BLCLE2
617 c**
618 c*****
619 c
620 c Ref: Appendix A, Figure 6, Step 3.
621 c
622 c=====
623 c
624 c Purpose: Given an ordered sample of size n and specified shape
625 c           c<=2, calculate the BLUEs of location a and scale b.
626 c
627 c=====
628 c
629 c Variables:
630 c           x = array containing n ordered Pareto deviates
631 c           c = shape parameter
632 c           n = sample size
633 c           B = array of B values used to calculate the BLUEs
634 c           nc = product of n and c
635 c           Coef1 = coefficient used to compute BLUE of location a
636 c           Coef2 = coefficient used to compute BLUE of location a
637 c           Coef3 = coefficient used to compute BLUE of scale b
638 c           Bxsum2 = sum of B(i)*x(i) terms for i = 1,...,n-2
639 c           Bxsm2c = sum of B(i)*x(i) terms for i = 1,...,n-2/c
640 c           ablu = BLUE of the location parameter a
641 c           bblu = BLUE of the scale parameter b
642 c           U = value used to compute BLUEs when c = 1.5
643 c           Termi = terms used to compute U (i=1,2,3)
644 c
645 c=====
646 c
647 c Input:   x = array of n ordered Pareto deviates (from PARDEV)
648 c           c = shape parameter = 0.5, 1.0, 1.5, or 2.0
649 c           n = sample size = 5(5)30 (from MAIN DO Loop 80)
650 c           nc = n*c (from MAIN program)
651 c           B = array containing n values of B (from BXVALS)
652 c           Bxsum2 = sum of first n-2 values of B (from BXVALS)
653 c           Bxsm2c = sum of first n-2/c values of B (from BXVALS)
654 c
655 c=====
656 c
657 c Calculate (if c = 0.5, 1, or 2):
658 c
659 c           Coef1 = [(c+1)*(c+2)] / [(nc-2)*(nc-c-2)]
660 c           Coef2 = (nc-2) / (c+2)
661 c
662 c           ablu = x(1) - Coef1 * [Bxsm2c - (Coef2*x(1))] (eqn 34)
663 c           bblu = (nc-1) * [x(1) - ablu] (eqn 35)
664 c

```

```

665 c -----
666 c Calculate (if c = 1.5):
667 c
668 c      Term1 = (nc-2) * (nc-c-2)
669 c      Term2 = nc * (c-2) * B(n-1)
670 c      Term3 = (nc-1) * (c+2)
671 c      Coef3 = [(nc-1)/nc] * (nc-2-U)
672 c      U = (Term1 - Term2) / Term3      (eqn 39)
673 c
674 c
675 c      ablu = x(1) - bblu / (nc-1)      (eqn 37)
676 c      bblu = (1/U) * [(c+1)*(Bxsum2) + (2c-1)*B(n-1)*x(n-1)
677 c             - Coef3 * x(1)]      (eqn 38)
678 c
679 c =====
680 c Output:
681 c      ablu = BLUE of location parameter a
682 c      bblu = BLUE of scale parameter b
683 c
684 c =====
685 c Declare Variables:
686 c
687 c      common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
688 c             1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
689 c             1      KScrit,ADcrit,CVcrit,Y
690 c
691 c      integer n
692 c      real x(30),ablu,bblu,B(30),D,c,nc,Bsum1,Bxsum1,Bxsum2,
693 c             1      Bxsm2c,P(30),Term1,Term2,Term3,Coef1,Coef2,Coef3,U,
694 c             1      KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),Y(5002)
695 c      double precision dseed
696 c
697 c
698 c      if ((c.eq.0.5) .or. (c.eq.1.0) .or. (c.eq.2.0)) then
699 c          Coef1 = ((c+1.0)*(c+2.0)) / ((nc-2.0)*(nc-c-2.0))
700 c          Coef2 = (nc-2.0) / (c+2.0)
701 c          ablu = x(1) - Coef1 * (Bxsm2c - (Coef2*x(1)))
702 c          bblu = (nc-1.0) * (x(1) - ablu)
703 c
704 c      else if (c .eq. 1.5) then
705 c          Term1 = (nc-2.0) * (nc-c-2.0)
706 c          Term2 = nc * (c-2.0) * B(n-1)
707 c          Term3 = (nc-1.0) * (c+2.0)
708 c          U = (Term1 - Term2) / Term3
709 c          Coef3 = ((nc-1.0)/nc) * (nc-2.0-U)
710 c          bblu = (1.0/U) * ((c+1.0) * (Bxsum2)
711 c             1      + (2.0*c-1.0)*B(n-1)*x(n-1) - Coef3 * x(1) )
712 c          ablu = x(1) - (bblu / (nc-1.0))
713 c
714 c      end if
715 c
716 c      return
717 c      end
718 c
719 c =====
720 c      END SUBROUTINE BLCLE2
721 c *****

```

```

722           Subroutine BLCGT2
723 c*****
724 c**
725 c**           B E G I N   S U B R O U T I N E   B L C G T 2           **
726 c**
727 c*****
728 c
729 c Ref: Appendix A, Figure 6, Step 3.
730 c
731 c=====
732 c
733 c Purpose: Given an ordered sample of size n and a specified
734 c           shape c > 2, calculate the best linear unbiased
735 c           estimates (BLUEs) of location and scale.
736 c
737 c=====
738 c
739 c Variables: x = array containing n ordered Pareto deviates
740 c             c = shape parameter
741 c             n = sample size
742 c             nc = product of n and c
743 c             B = array of B values used to calculate the BLUEs
744 c             Bsum1 = sum of B(i) terms for i = 1,...,n-1
745 c             Bxsum1 = sum of B(i)*x(i) terms for i = 1,...,n-1
746 c             D = value used to calculate the BLUEs
747 c             YV = value used to calculate the BLUEs
748 c             ablu = BLUE for location parameter a
749 c             bblu = BLUE for scale parameter b
750 c
751 c=====
752 c
753 c Input:      x = array of ordered Pareto deviates (from PARDEV)
754 c             c = shape parameter = 2.5, 3.0, 3.5, or 4.0
755 c             n = sample size = 5(5)30 (from MAIN DO Loop 80)
756 c             nc = n*c (from MAIN Program)
757 c             B = array of B values (from BXVALS)
758 c             Bsum1 = sum of first (n-1) B values (from BXVALS)
759 c             Bxsum1 = sum of first n-1 B*x values (from BXVALS)
760 c
761 c=====
762 c
763 c Calculate:
764 c
765 c           D = [(c+1) * Bsum1] + [(c-1) * B(n)]           (eqn 21)
766 c
767 c           YV = (c+1)*Bxsum1 + (c-1)*B(n)*x(n) - D*x(1) (eqn 22)
768 c
769 c           ablu = x(1) - YV/[(nc-1)*(nc-2) - D*nc]       (eqn 17)
770 c
771 c           bblu = (nc-1) * [ x(1) - ablu ]               (eqn 18)
772 c
773 c=====

```

```

774      c
775      c  Output:   ablu = BLUE for location a
776      c           bblu = BLUE for scale b
777      c
778      c=====
779      c
780      c  Declare Variables:
781      c
782      c      common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
783      1          Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
784      1          KScrit,ADcrit,CVcrit,Y
785      c      integer n
786      c      real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),YV,
787      1          CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
788      1          KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),r(30),
789      1          Y(5002)
790      c      double precision dseed
791      c
792      c      D = ((c+1.0) * Bsum1) + ((c-1.0) * B(n))
793      c      YV = ((c+1.0)*Bxsum1) + ((c-1.0)*B(n)*x(n)) - (D*x(1))
794      c      ablu = x(1) - YV/((nc-1.0)*(nc-2.0) - (D*nc))
795      c      bblu = (nc-1.0) * (x(1) - ablu)
796      c
797      c      return
798      c      end
799      c
800      c=====
801      c                      END SUBROUTINE BLCGT2
802      c*****

```

```

803           Subroutine HYPcdf
804 c*****
805 c**
806 c**       B E G I N   S U B R O U T I N E   H Y P C D F       **
807 c**
808 c*****
809 c
810 c   Ref:  Appendix A. Figure 6, Step 4.
811 c
812 c=====
813 c
814 c   Purpose:  Given an ordered sample of size n, a specified
815 c             shape c, and the BLUEs of location a and scale b,
816 c             compute the hypothesized Pareto distribution
817 c             function P(i) for i = 1,2,...,n.
818 c
819 c=====
820 c
821 c   Variables:
822 c             x = array containing n ordered Pareto deviates
823 c             n = sample size
824 c             c = shape parameter
825 c             ablu = BLUE of location a
826 c             bblu = BLUE of scale b
827 c             P = array containing n points of the
828 c             hypothesized Pareto CDF
829 c
830 c=====
831 c
832 c   Input:
833 c             x = array of n ordered Pareto deviates (from PARDEV)
834 c             c = shape parameter = .5(.5)4 (from MAIN DO Loop 90)
835 c             n = sample size = 5(5)30 (from MAIN DO Loop 80)
836 c             ablu = BLUE of location a (from BLCLE2 or BLCGT2)
837 c             bblu = BLUE of scale b (from BLCLE2 or BLCGT2)
838 c
839 c=====
840 c
841 c   Calculate:
842 c
843 c       P(i) = 1 - [1 / [1 + (x(i) - ablu)/bblu] ]**c      (eqn 40)
844 c
845 c=====
846 c
847 c   Output:  P = array of n points of the hypothesized CDF
848 c
849 c=====
850 c
851 c   Declare Variables:

```

```

852      c
853      common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
854      1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
855      1      KScrit,ADcrit,CVcrit,v
856      integer n
857      real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),
858      1      CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
859      1      KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),r(30),
860      1      Y(5002)
861      double precision dseed
862      c
863      do 10 i = 1,n
864      P(i) = 1.0 - (1.0 + (x(i) - ablu)/bblu) ** (-c)
865      10 continue
866      c
867      return
868      end
869      c
870      c=====
871      c                      END SUBROUTINE HYPcdf
872      c*****

```

```

873           Subroutine TESTAT
874 c*****
875 c**
876 c**           B E G I N   S U B R O U T I N E   T E S T A T           **
877 c**
878 c*****
879 c
880 c  Ref:  Appendix A, Figure 6, Step 5.
881 c
882 c=====
883 c  Purpose:  Given a sample size n, and the hypothesized Pareto
884 c            distribution function P(i), compute values of the
885 c            test statistics of the modified K-S, A-D, and CVM
886 c            goodness-of-fit tests.
887 c=====
888 c  Variables:
889 c            n = sample size
890 c            nshp = shape parameter counter (8 values, 1-8)
891 c            nsiz = sample size counter (6 values, 1-6)
892 c            it = iteration counter (1-5000)
893 c            P = array of n values of the hypothesized Pareto CDF
894 c-----
895 c            DP = positive differences between EDF and CDF points
896 c            DM = negative differences between EDF and CDF points
897 c            DPLUS = maximum positive difference (largest DP value)
898 c            DMINUS = maximum negative difference (largest DM value)
899 c            KS = values of the modified K-S test statistic
900 c-----
901 c            AL = value used to calculate the A-D test statistic
902 c            AM = value used to calculate the A-D test statistic
903 c            AN = AL + AM
904 c            AAA = values to be summed for A-D test statistic
905 c            SAAA = sum of AAA values
906 c            AD = values of the modified A-D test statistic
907 c-----
908 c            ACV = squared quantities in the C-VM formula
909 c            SACV = sum of the ACV values
910 c            CVM = values of the modified C-VM test statistic
911 c=====
912 c  Input:
913 c            n = sample size = 5(5)30 (from MAIN DO Loop 80)
914 c            P = array of n values of hypothesized CDF (from HYPcdf)
915 c            it = iteration counter (from MAIN DO Loop 60)
916 c            nsiz = sample size counter (from MAIN DO Loop 80)
917 c            nshp = shape parameter counter (from MAIN DO Loop 90)
918 c=====
919 c  Calculations for K-S test statistic (eqns 41 & 42):
920 c
921 c            DP(i) = ABS[ (i/n) - P(i) ]
922 c            DM(i) = ABS[ P(i) - (i-1)/n ]
923 c
924 c            DPLUS = max [ DP(i) ] for i=1,2,...,n

```

```

932      c           DMINUS = max [ DM(i) ] for i=1,2,...,n
933      c
934      c           KS = max (DPLUS,DMINUS)
935      c
936      c -----
937      c
938      c   Calculations for A-D test statistic (eqn 43):
939      c
940      c           AL(j) = ln (P(j))
941      c           AM(j) = ln (1 - P(n+1-j))
942      c           AN(j) = AL(j) + AM(j)
943      c
944      c           AAA(j) = (2*j - 1) * AN(j)
945      c           SAAA = AAA(1) + AAA(2) + ... + AAA(n)
946      c
947      c           AD = -n - (1/n) * SAAA
948      c
949      c -----
950      c
951      c   Calculations for C-VM test statistic (eqn 44):
952      c
953      c           ACV(k) = [ P(k) - (2*k - 1)/(2*n) ]**2
954      c           SACV = ACV(1) + ACV(2) + ... + ACV(n)
955      c
956      c           CVM = (1/(12*n)) + SACV
957      c
958      c
959      c =====
960      c
961      c   Declare Variables:
962      c
963      c           common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
964      c           1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
965      c           1      KScrit,ADcrit,CVcrit,Y
966      c           integer n,nsiz,nshp,it
967      c           real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),
968      c           1      CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
969      c           1      KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),r(30),
970      c           1      DP(30),DM(30),DPLUS,DMINUS,AL(30),AM(30),
971      c           1      AN(30),AAA(30),SAAA,ACV(30),SACV,Y(5002)
972      c           double precision dseed
973      c
974      c           DPLUS = 0
975      c           DMINUS = 0
976      c
977      c           do 5 ik = 1,30
978      c               DP(ik) = 0
979      c               DM(ik) = 0
980      c           5      continue
981      c
982      c -----   Compute the K-S Test Statistic (eqns 41 & 42):   -----
983      c

```

```

984      do 10 i = 1,n
985          DP(i) = ABS( (i/real(n)) - P(i) )
986          DM(i) = ABS( P(i) - (i-1)/real(n) )
987      10  continue
988      c
989          DPLUS = MAX( DP(1),DP(2),DP(3),DP(4),DP(5),DP(6),DP(7),
990      1          DP(8),DP(9),DP(10),DP(11),DP(12),DP(13),DP(14),
991      1          DP(15),DP(16),DP(17),DP(18),DP(19),DP(20),
992      1          DP(21),DP(22),DP(23),DP(24),DP(25),DP(26),
993      1          DP(27),DP(28),DP(29),DP(30) )
994      c
995          DMINUS = MAX( DM(1),DM(2),DM(3),DM(4),DM(5),DM(6),DM(7),
996      1          DM(8),DM(9),DM(10),DM(11),DM(12),DM(13),DM(14),
997      1          DM(15),DM(16),DM(17),DM(18),DM(19),DM(20),
998      1          DM(21),DM(22),DM(23),DM(24),DM(25),DM(26),
999      1          DM(27),DM(28),DM(29),DM(30) )
1000      c
1001          KS(it,nsiz,nshp) = MAX(DPLUS,DMINUS)
1002      c
1003      c ----- Compute the A-D Test Statistic (eqn 43): -----
1004      c
1005          SAAA = 0
1006      c
1007          do 20 j = 1,n
1008      c
1009              AL(j) = log (P(j))
1010              AM(j) = log (1.0 - P(n+1-j))
1011              AN(j) = AL(j) + AM(j)
1012              AAA(j) = (2.0*j - 1.0) * AN(j)
1013              SAAA = SAAA + AAA(j)
1014      c
1015      20  continue
1016      c
1017          AD(it,nsiz,nshp) = -n - (1.0/real(n)) * SAAA
1018      c
1019      c ----- Compute the C-VM Test Statistic (eqn 44): -----
1020      c
1021          SACV = 0
1022      c
1023          do 30 k = 1,n
1024              ACV(k) = ( P(k) - (2.0*k - 1.0)/(2.0*real(n)) )**2
1025              SACV = SACV + ACV(k)
1026      30  continue
1027      c
1028          CVM(it,nsiz,nshp) = SACV + (1.0/(12.0*real(n)))
1029      c
1030          return
1031          end
1032      c
1033      c=====
1034      c          END SUBROUTINE TESTAT
1035      c*****

```

```

1036           Subroutine CRTVAL
1037 c*****
1038 c**
1039 c**           B E G I N   S U B R O U T I N E   C R T V A L           **
1040 c**
1041 c*****
1042 c
1043 c   Ref: Appendix A, Figure 6, Step 7.
1044 c
1045 c=====
1046 c
1047 c   Purpose:
1048 c
1049 c       Given a set of 5000 values of test statistics from the
1050 c       modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D),
1051 c       or Cramer-von Mises (C-VM) test, select critical values
1052 c       by using median ranks plotting positions to compute
1053 c       specified percentile levels.
1054 c
1055 c=====
1056 c
1057 c   Variables:
1058 c           c = shape parameter
1059 c           n = sample size
1060 c           pct = percentile value
1061 c           nshp = shape parameter counter (1: c=.5; 2: c=1.0;
1062 c           3: c=1.5; 4: c=2.0; 5: c=2.5; 6: c=3.0;
1063 c           7: c=3.5; 8: c=4.0)
1064 c           nsiz = sample size counter (1: n=5 or 6; 2: n=10;
1065 c           3: n=15; 4: n=20; 5: n=25; 6: n=30)
1066 c           npct = percentile counter (0: pct=0; 1: pct=.80;
1067 c           2: pct=.85; 3: pct=.9; 4: pct=.95; 5: pct=.99)
1068 c           nst = total number of statistics used
1069 c           it = iteration counter (5000 repetitions required)
1070 c           KS = 3D array of 5000 modified K-S test statistics
1071 c           KS1 = 1D array of 5000 K-S test statistics
1072 c           CVM = 3D array of 5000 modified C-VM test statistics
1073 c           CV1 = 1D array of 5000 C-VM test statistics
1074 c           AD = 3D array of 5000 modified A-D test statistics
1075 c           AD1 = 1D array of 5000 A-D statistics
1076 c           STAT = 1D array of test stats (KS, AD, or CVM)
1077 c           KScrit = array of critical values for the K-S test
1078 c           CVMcrit = array of critical values for the C-VM test
1079 c           ADCrit = array of critical values for the A-D test
1080 c           CRIT = either the KS, AD, or CVM critical value array
1081 c           Y = array containing 5002 plotting positions
1082 c           slpm = array of slopes used to find critical values
1083 c           bi = array of intercepts used to find critical vals
1084 c
1085 c=====
1086 c
1087 c   Input:

```

```

1088      c          Y = array of plotting positions (MAIN DO Loop 10)
1089      c          c = shape parameter (from MAIN DO Loop 90)
1090      c          n = sample size (from MAIN DO Loop 80)
1091      c          nshp = shape parameter counter (from MAIN DO Loop 90)
1092      c          nsiz = sample size counter (from MAIN DO Loop 80)
1093      c          npct = percentile counter (from MAIN DO Loop 70)
1094      c          nst = number of test statistics used (from MAIN Prog)
1095      c          KS = array of 5000 K-S test statistics (from TESTAT)
1096      c          CVM = array of 5000 C-VM test stats (from TESTAT)
1097      c          AD = array of 5000 A-D test statistics (from TESTAT)
1098      c
1099      c=====
1100      c
1101      c  IMSL Subroutine:  VSRTA - orders the test statistic values
1102      c
1103      c=====
1104      c
1105      c  Calculate Endpoints of Test Statistics (Eqns 52 - 57):
1106      c
1107      c          slpm(0) = ( Y(2) - Y(1) ) / ( STAT(2) - STAT(1) )
1108      c          bi(0) = Y(1) - slpm(0) * STAT(1)
1109      c          STAT(0) = max ( 0, - bi(0)/slpm(0) )
1110      c
1111      c          slpm(6) = (Y(5000) - Y(4999))/(STAT(5000) - STAT(4999))
1112      c          bi(6) = Y(4999) - slpm(6) * STAT(4999)
1113      c          STAT(6) = (1.0 - bi(6)) / slpm(6)
1114      c
1115      c-----
1116      c
1117      c  Calculate Critical Values (Eqns 58 - 60):
1118      c
1119      c          slpm(npct) = ( Y(j+1) - Y(j) ) / ( STAT(j+1) - STAT(j) )
1120      c          bi(npct) = Y(j) - slpm(npct) * STAT(j)
1121      c          CRIT(npct) = (pct - bi(npct)) / slpm(npct)
1122      c
1123      c=====
1124      c
1125      c  Output:
1126      c
1127      c          KScrit - array of critical values for modified K-S test
1128      c          ADcrit - array of critical values for modified A-D test
1129      c          CVcrit - array of critical values for modified C-VM test
1130      c
1131      c=====
1132      c
1133      c  Declare Variables:
1134      c
1135      c          common dseed,x,n,c,nc,B,D,ablu,bblu,P,pct,Bsum1,Bxsum1,
1136      c          1          Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,npct,nst,
1137      c          1          KScrit,ADcrit,CVcrit,Y
1138      c          integer n,nsiz,nshp,it,npct,nst,nstest
1139      c          real x(30),ablu,bblu,B(30),D,KS(5000,6,8),AD(5000,6,8),

```

```

1140      1      CVM(5000,6,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,P(30),
1141      1      KScrit(6,8,5),ADcrit(6,8,5),CVcrit(6,8,5),r(30),
1142      1      Y(5002),STAT(5002),CRIT(6,8,7),slpm(7),bi(7),pct,
1143      1      KS1(5000),CV1(5000),AD1(5000)
1144      double precision dseed
1145      c
1146      if (npct .eq. 1) pct = .80
1147      if (npct .eq. 2) pct = .85
1148      if (npct .eq. 3) pct = .90
1149      if (npct .eq. 4) pct = .95
1150      if (npct .eq. 5) pct = .99
1151      c
1152      c      ** Store the 3 Sets of 5000 Test Stats into 1D Arrays: **
1153      c
1154      do 16 ncnt = 1,nst
1155      KS1(ncnt) = KS(ncnt,nsiz,nshp)
1156      AD1(ncnt) = AD(ncnt,nsiz,nshp)
1157      CV1(ncnt) = CVM(ncnt,nsiz,nshp)
1158      16  continue
1159      c
1160      c      ** Use IMSL Subroutine to Order the Test Statistics: **
1161      c
1162      Call VSRTA(KS1,nst)
1163      c      print*, 'ORDERED KS STATISTICS FROM CRTVAL:'
1164      c      print*, 'n=',n, ' c=',c
1165      c      do 2 jks = 1,nst
1166      c      print*, 'KS STAT =',KS1(jks)
1167      c  2  continue
1168      c
1169      Call VSRTA(AD1,nst)
1170      c      print*, 'ORDERED AD STATISTICS FROM CRTVAL:'
1171      c      print*, 'n=',n, ' c=',c
1172      c      do 4 jad = 1,nst
1173      c      print*, 'AD STAT =',AD1(jad)
1174      c  4  continue
1175      c
1176      Call VSRTA(CV1,nst)
1177      c      print*, 'ORDERED CVM STATISTICS FROM CRTVAL:'
1178      c      print*, 'n=',n, ' c=',c
1179      c      do 6 jcv = 1,nst
1180      c      print*, 'CV STAT =',CV1(jcv)
1181      c  6  continue
1182      c
1183      c      --- Begin DO Loop 20 to Rotate Through KS, AD, and CVM ---
1184      c
1185      do 20 ntest = 1,3
1186      c
1187      c      --- Begin DO Loop 30 for 5000 Data Points ---
1188      c
1189      do 30 j = 1,nst
1190      c
1191      if (ntest .eq. 1) then

```

```

1192          STAT(j) = KS1(j)
1193          else if (ntest .eq. 2) then
1194              STAT(j) = AD1(j)
1195          else if (ntest .eq. 3) then
1196              STAT(j) = CV1(j)
1197          end if
1198      c
1199      30      continue
1200      c
1201      ---      End DO Loop 30 for 5000 Data Points      ---
1202      c
1203      ** Extrapolate Left Endpoint of the Test Statistics: **
1204      c
1205          if (STAT(1) .eq. STAT(2)) then
1206      c
1207          print*, '*****'
1208      c          print*, 'TWO LEFT ENDPOINT STATS EQUAL'
1209      c              if (ntest .eq. 1) print*, 'FOR KS TEST'
1210      c              if (ntest .eq. 2) print*, 'FOR AD TEST'
1211      c              if (ntest .eq. 3) print*, 'FOR CVM TEST'
1212      c          print*, 'n=', n, ' c=', c, ' pct=', pct
1213      c          print*, 'STAT(1)=', STAT(1)
1214      c          print*, 'STAT(2)=', STAT(2)
1215      c          print*, '%%%%%%%%%'
1216      c          print*, ' '
1217      c
1218          dif0 = STAT(3) - STAT(1)
1219          if (dif0 .eq. 0.0) dif0 = .00001
1220          slpm(0) = (Y(3) - Y(1)) / dif0
1221      else
1222          dif0 = STAT(2) - STAT(1)
1223          slpm(0) = (Y(2) - Y(1)) / dif0
1224      end if
1225      c
1226          bi(0) = Y(1) - slpm(0) * STAT(1)
1227          STAT(0) = max( 0.0, - bi(0)/slpm(0) )
1228      c          print*, ' '
1229      c          print*, '=====
1230      c              if (ntest .eq. 1) print*, 'FOR KS TEST STATISTICS'
1231      c              if (ntest .eq. 2) print*, 'FOR AD TEST STATISTICS'
1232      c              if (ntest .eq. 3) print*, 'FOR CVM TEST STATISTICS'
1233      c          print*, 'LEFT ENDPT X(0000) =', STAT(0)
1234      c          print*, '-----FIRST X(0001) =', STAT(1)
1235      c          print*, '80PCT STAT X(4000) =', STAT(4000)
1236      c          print*, '85PCT STAT X(4250) =', STAT(4250)
1237      c          print*, '90PCT STAT X(4500) =', STAT(4500)
1238      c          print*, '95PCT STAT X(4750) =', STAT(4750)
1239      c          print*, '99PCT STAT X(4950) =', STAT(4950)
1240      c          print*, '----- LAST X(5000) =', STAT(5000)
1241      c
1242      c          ** Extrapolate Right Endpoint of the Test Statistic: **
1243      c

```

```

1244         if (STAT(nst-1) .eq. STAT(nst)) then
1245     c
1246         print*, '*****'
1247         print*, 'TWO RIGHT ENDPOINT STATS EQUAL:'
1248         if (ntest .eq. 1) print*, 'FOR KS TEST'
1249         if (ntest .eq. 2) print*, 'FOR AD TEST'
1250         if (ntest .eq. 3) print*, 'FOR CVM TEST'
1251         print*, 'n=', n, ' c=', c, ' pct=', pct
1252         print*, 'STAT(4999)=', STAT(nst-1)
1253         print*, 'STAT(5000)=', STAT(nst)
1254         print*, 'XXXXXXXXXXXXXXXXXXXXXXXXXXXX'
1255         print*, ' '
1256     c
1257         dif6 = STAT(nst) - STAT(nst-2)
1258         if (dif6 .eq. 0.0) dif6 = .00001
1259         slpm(6) = (Y(nst)-Y(nst-2)) / dif6
1260     else
1261         dif6 = STAT(nst) - STAT(nst-1)
1262         slpm(6) = (Y(nst)-Y(nst-1)) / dif6
1263     end if
1264     c
1265     bi(6) = Y(nst-1) - slpm(6)*STAT(nst-1)
1266     STAT(nst+1) = (1.0 - bi(6)) / slpm(6)
1267     print*, 'RGHT ENDPT X(5001) =', STAT(nst+1)
1268     c
1269     ** Interpolate Critical Values Between Test Stats: **
1270     c
1271     -- Begin DO Loop 50 to Find Max Y(k) < pct: --
1272     c
1273     do 50 k, 1, nst
1274         k = nst+1 - k
1275     c
1276         if (Y(k) .le. pct) then
1277     c
1278         if (STAT(k) .eq. STAT(k+1)) then
1279     c
1280         print*, '*****'
1281         print*, 'TWO ADJACENT STATS EQUAL:'
1282         if (ntest .eq. 1) print*, 'FOR KS TEST'
1283         if (ntest .eq. 2) print*, 'FOR AD TEST'
1284         if (ntest .eq. 3) print*, 'FOR CVM TEST'
1285         print*, 'n=', n, ' c=', c, ' pct=', pct
1286         print*, 'STAT(k)=', STAT(k)
1287         print*, 'STAT(k+1)=', STAT(k+1)
1288         print*, 'XXXXXXXXXXXXXXXXXXXXXXXXXXXX'
1289         print*, ' '
1290     c
1291         dif = STAT(k+1) - STAT(k-1)
1292         if (dif .eq. 0.0) dif = .00001
1293         slpm(npct) = (Y(k+1)-Y(k-1)) / dif
1294     else
1295         dif = STAT(k+1) - STAT(k)

```

```

1296         slpm(npct) = (Y(k+1)-Y(k)) / dif
1297     end if
1298 c
1299         bi(npct) = Y(k) - slpm(npct) * STAT(k)
1300         CRIT(nsiz,nshp,npct)
1301     1         = (pct-bi(npct))/slpm(npct)
1302             GOTO 75
1303 c
1304         end if
1305 c
1306     50     continue
1307 c
1308         -- End DO Loop 50 Upon Finding Crit Val --
1309 c
1310     ** Associate the Critical Values with Test Type: **
1311 c
1312     75     if (ntest .eq. 1) then
1313             KScrit(nsiz,nshp,npct) = CRIT(nsiz,nshp,npct)
1314     c         print*,'n=',n,' ** c=',c,' pct=',pct
1315     c         print*,'CRTVAL KS Crit Val =',KScrit(nsiz,nshp,npct)
1316     else if (ntest .eq. 2) then
1317             ADcrit(nsiz,nshp,npct) = CRIT(nsiz,nshp,npct)
1318     c         print*,'CRTVAL AD Crit Val =',ADcrit(nsiz,nshp,npct)
1319     else if (ntest .eq. 3) then
1320             CVcrit(nsiz,nshp,npct) = CRIT(nsiz,nshp,npct)
1321     c         print*,'CRTVAL CV Crit Val =',CVcrit(nsiz,nshp,npct)
1322     c         print*,' '
1323     end if
1324 c
1325     20 continue
1326 c
1327     --- End DO Loop 20 After Rotating Through KS, AD, and CVM ---
1328 c
1329     return
1330     end
1331 c
1332 c=====
1333 c         END SUBROUTINE CRTVAL
1334 c*****

```

**APPENDIX B**  
**Computer Program and Subroutines**  
**for Determining Power Values**

Computer Program  
POWER

Subroutines  
PARETO GGWIB  
GGAMR GGBTR  
GGEXN GGNML

Subroutines  
VSRTA BLCGT2  
BXVALS HYPCDF  
BLCLE2 TESTAT

Subroutine  
COMPAR

Main Program  
DO Loop 40

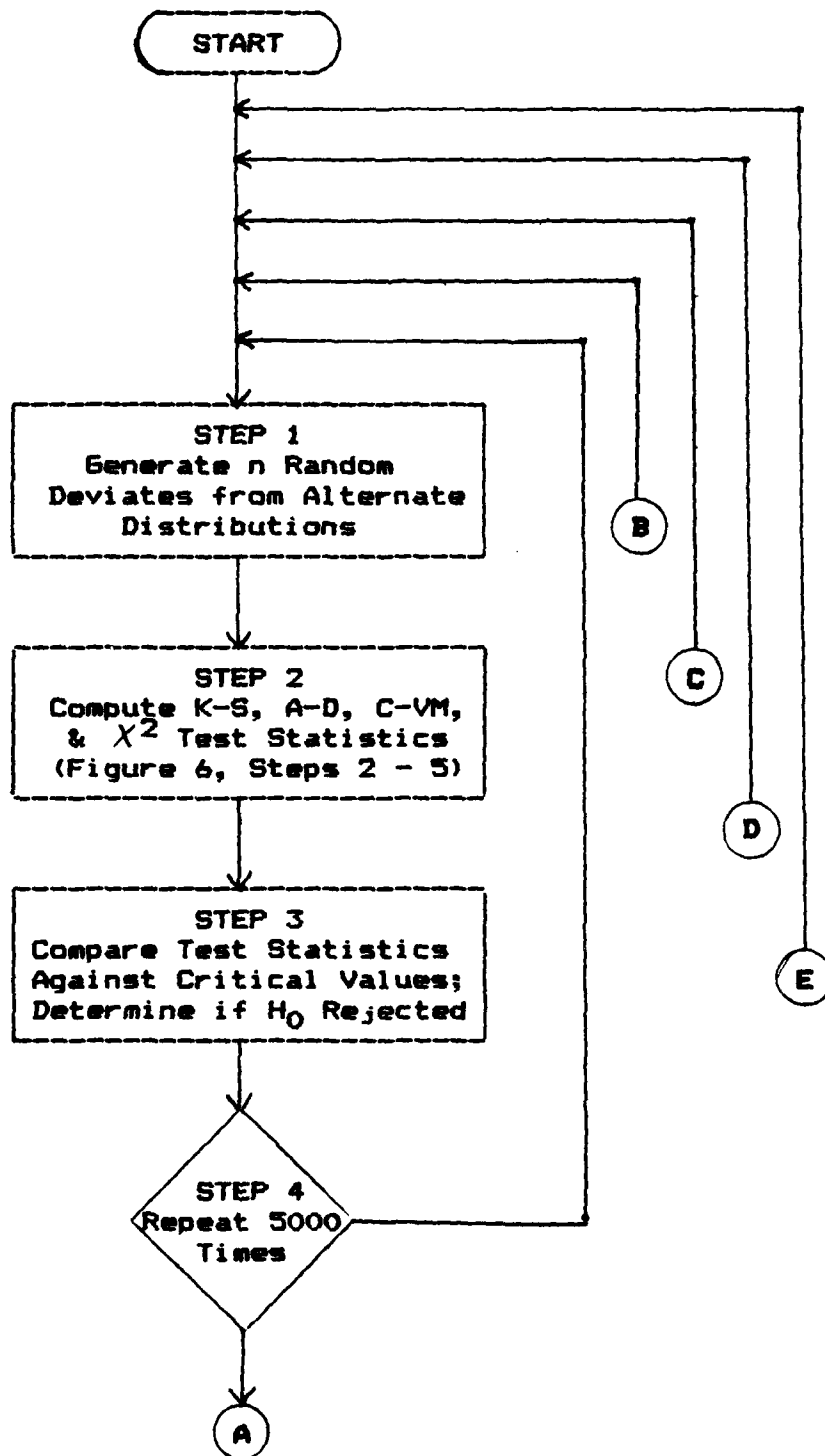


Fig 7. Procedure for Determining Power Values

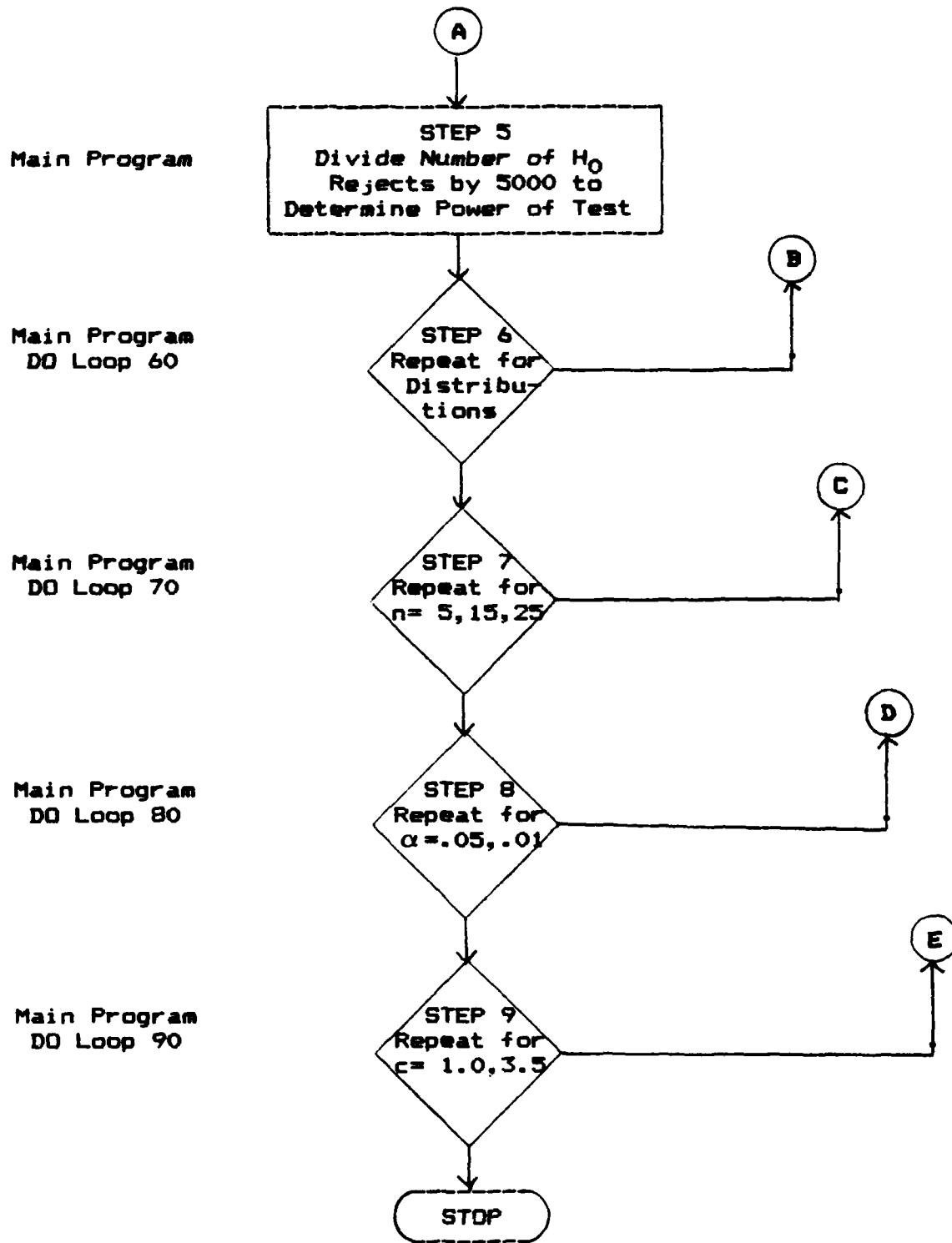


Fig 7 (Continued). Procedure for Determining Power Values

```

1      c***** Classroom Support Computer (CSC) - VAX 11/785 - VMS 4.1 ****
2      c
3      c*****      POWER PROGRAM FOR PARETO GOODNESS-OF-FIT TESTS      *****
4      c
5      c*****
6      c*****
7      c**
8      c**      BEGIN      POWER      MAIN      PROGRAM      **
9      c**
10     c*****
11     c
12     c Ref: Appendix B, Figure 7.
13     c
14     c=====
15     c
16     c Purpose: Test the null hypothesis that a set of sample data
17     c follows the Pareto distribution with hypothesized shape c
18     c against the alternate hypothesis that the data follow some
19     c other distribution. The goals are to:
20     c
21     c 1. Compare powers of the modified Kolmogorov-Smirnov (K-S),
22     c Anderson-Darling (A-D), and Cramer-von Mises (C-VM) tests
23     c against the Chi-Square test to determine which test can
24     c best detect a false Pareto distribution hypothesis.
25     c
26     c 2. When the Pareto null hypothesis is true, confirm that
27     c the hypothesis rejection rates under the modified K-S, A-D,
28     c and C-VM statistics are low enough to satisfy a claimed level
29     c of significance.
30     c
31     c 3. Provide extensive commentary to assist novice programmers
32     c to conduct similar power studies in statistical analysis.
33     c Diagnostic print statements have been retained as commentary
34     c to contribute to this goal.
35     c
36     c=====
37     c
38     c Variables:
39     c      dseed = random number seed
40     c      alpha = level of significance (.01 or .05 used here)
41     c      n = sample size
42     c      c = null-hypothesis Pareto shape parameter
43     c      nshp = null-hyp Pareto shape counter (1:c=1.0, 2:c=3.5)
44     c      nalf = significance level counter (1:  $\alpha$  = .05, 2:  $\alpha$  = .01)
45     c      nsiz = sample size counter (1:n=5, 2:n=15, 3:n=25)
46     c      nalt = alternative distribution counter (8 in all)
47     c      nrep = number of repetitions to be used
48     c      it = iteration counter (5000 repetitions required)
49     c      KS = array of values of modified K-S test statistic
50     c      CVM = array of values of modified C-VM test statistic
51     c      AD = array of values of modified A-D test statistic
52     c      X2 = array of values of Chi-square test statistic

```

```

53      c          nrKS = number of hypothesis rejects under the K-S test
54      c          nrAD = number of hypothesis rejects under the A-D test
55      c          nrCV = number of hypothesis rejects under the CVM test
56      c          nrX2 = number of hypothesis rejects under Chi-square
57      c
58      c=====
59      c
60      c  Input:
61      c      nrep = number of repetitions (input at computer terminal)
62      c      dseed = random number seed (input at computer terminal)
63      c
64      c=====
65      c
66      c  Subroutines:
67      c
68      c  PARETO - Generates n random Pareto deviates
69      c  BXVALS - Calculates B values and summations of B and Bx
70      c  BLCLE2 - Finds BLUEs for location and scale when c <= 2
71      c  BLCGT2 - Finds BLUEs for location and scale when c > 2
72      c  HYPCDF - Computes the Hypothesized Pareto CDF
73      c  TESTAT - Calculates the K-S, A-D, and C-VM test statistics
74      c  COMPAR - Compares test stats vs. crit vals and counts rejects
75      c
76      c-----
77      c
78      c  IMSL Subroutines:
79      c
80      c  GGWIB - Generates random Weibull deviates
81      c  GGAMR - Generates random Gamma deviates
82      c  GGBTR - Generates random Beta deviates
83      c  GGEXN - Generates random Exponential Deviates
84      c  GGNML - Generates random Normal Deviates
85      c  VSRTA - Arranges data in ascending order
86      c
87      c=====
88      c
89      c  Output:
90      c
91      c  KSpwr (nshp,nalf,nsiz,nalt) = power values for K-S test
92      c  ADpwr (nshp,nalf,nsiz,nalt) = power values for A-D test
93      c  CVpwr (nshp,nalf,nsiz,nalt) = power values for C-VM test
94      c  X2pwr (nshp,nalf,nsiz,nalt) = power values for Chi-square
95      c
96      c=====
97      c
98      c  Declare Variables:
99      c
100     c      common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
101     1          Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
102     1          nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
103     integer  n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
104     1          nrCV(2,2,3,8),nrX2(2,2,3,8)

```

```

105      real      x(25), ablu, bblu, B(25), D, KS(2, 2, 3, 8), AD(2, 2, 3, 8),
106      1          CVM(2, 2, 3, 8), c, nc, Bsum1, Bxsum1, Bxsum2, Bxsm2c,
107      1          P(25), r(25), alpha, KSpwr(2, 2, 3, 8), ADpwr(2, 2, 3, 8),
108      1          CVpwr(2, 2, 3, 8), X2crit(2, 2, 3), X2(2, 2, 3, 8),
109      1          X2pwr(2, 2, 3, 8)
110      character test(4)*3, altcdf(8)*12
111      double precision dseed
112      c
113          test(1) = 'K-S'
114          test(2) = 'A-D'
115          test(3) = 'CVM'
116          test(4) = 'CHI'
117      c
118          altcdf(1) = 'Pareto c=1.0'
119          altcdf(2) = 'Pareto c=3.5'
120          altcdf(3) = 'Pareto c=2.0'
121          altcdf(4) = 'Weibull'
122          altcdf(5) = 'Gamma'
123          altcdf(6) = 'Beta'
124          altcdf(7) = 'Exponential'
125          altcdf(8) = 'Normal'
126      c
127      c      ** Open Output File to Store Computed Power Values: **
128          open (unit=7, file='X2ALL', status='new')
129      c
130      c      ** Number of Repetitions to be Used on Each Run: **
131          print*, 'The Monte Carlo power analysis will require'
132          print*, '      5000 repetitions.'
133          print*, 'Enter the number to be used for this run:'
134          read*, nrep
135      c
136          print*, 'Enter random number seed or "1." for default:'
137          read*, dseed
138          if (dseed .eq. 1.) dseed = 123457.d00
139          print*, ' '
140          print*, 'STANDBY . . . COMPUTATIONS IN PROGRESS'
141      c
142      c --- Begin DO Loop 90 for Null-Hypothesis Pareto Shape c ---
143      c
144          do 90 nshp = 1, 2
145      c
146          if (nshp .eq. 1) then
147              c = 1.0
148              write(7, 51)
149              write(7, 56)
150              write(7, 58)
151              write(7, 62)
152          else if (nshp .eq. 2) then
153              c = 3.5
154              write(7, 52)
155              write(7, 56)
156              write(7, 59)

```

```

157         write(7,62)
158     end if
159     c
160     c     --- Begin DO Loop 80 for Alpha Significance Levels ---
161     c
162     do 80 nalf =1,2
163     c
164         if (nalf .eq. 1) then
165             alpha = .05
166             write(7,64)
167         else if (nalf .eq. 2) then
168             alpha = .01
169             write(7,66)
170         end if
171     c
172         write(7,54)
173         write(7,74)
174         write(7,68)
175         write(7,72)
176         write(7,76)
177         write(7,72)
178     c
179         nsiz = 0
180     c
181     c     print*, '=====
182     c     print*, 'Numbers of Rejects After do 80/Before do 70'
183     c     print*, 'c =', c, 'alpha =', alpha, 'n=', n, 'CDF: ', altcdf(nalt)
184     c     print*, 'KS Rejects = ', nrKS(nshp, nalf, nsiz, nalt)
185     c     print*, 'AD Rejects = ', nrAD(nshp, nalf, nsiz, nalt)
186     c     print*, 'CV Rejects = ', nrCV(nshp, nalf, nsiz, nalt)
187     c     print*, '=====
188     c
189     c     --- Begin DO Loop 70 for Sample Sizes ---
190     c
191     do 70 n = 5,25,10
192     c
193         nsiz = nsiz + 1
194     c
195         nc = n * c
196     c
197     c     -- Begin DO Loop 60 for Alternate CDFs --
198     c
199     do 60 nalt = 1,8
200     c
201     c
202         nrKS(nshp, nalf, nsiz, nalt) = 0
203         nrAD(nshp, nalf, nsiz, nalt) = 0
204         nrCV(nshp, nalf, nsiz, nalt) = 0
205         nrX2(nshp, nalf, nsiz, nalt) = 0
206     c
207     c     -- Begin DO Loop 40 for Repetitions --
208     c

```

```

209          do 40 it = 1,nrep
210          c
211          c      ** Perform Step 1 of Figure 7: **
212          c
213          if (nalt .eq. 1) call PARETO
214          if (nalt .eq. 2) call PARETO
215          if (nalt .eq. 3) call PARETO
216          if (nalt .eq. 4) call GGWIB(dseed,3.5,n,x)
217          if (nalt .eq. 5) call GGAMR(dseed,2.,n,1,x)
218          if (nalt .eq. 6) call GGBTR(dseed,2.,3.,n,x)
219          if (nalt .eq. 7) call GGEXN(dseed,2.,n,x)
220          if (nalt .eq. 8) call GGNML(dseed,n,x)
221          c
222          c      ** Perform Step 2 of Figure 7: **
223          c
224          call VSRTA(x,n)
225          call BXVALS
226          c
227          if (c .eq. 1.0) call BLCLE2
228          if (c .eq. 3.5) call BLCGT2
229          c
230          call HYPCDF
231          call TESTAT
232          c
233          c      ** Perform Step 3 of Figure 7: **
234          c
235          call COMPAR
236          c
237          40      continue
238          c
239          c      -- End DO Loop 40 for Repetitions --
240          c      ** Completes Step 4 of Figure 7 **
241          c
242          c      ** Perform Step 5 of Figure 7: **
243          c
244          c      print*, '======'
245          c      print*, 'Numbers of Rejects Prior to Power Calculation'
246          c      print*, 'c =', c, 'alpha =', alpha, 'n =', n, 'nalt =', nalt
247          c      print*, 'KS Rejects = ', nrKS(nshp,nalf,nsiz,nalt)
248          c      print*, 'AD Rejects = ', nrAD(nshp,nalf,nsiz,nalt)
249          c      print*, 'CV Rejects = ', nrCV(nshp,nalf,nsiz,nalt)
250          c      print*, 'X2 Rejects = ', nrX2(nshp,nalf,nsiz,nalt)
251          c      print*, '======'
252          c
253          c      KSpwr(nshp,nalf,nsiz,nalt) =
254          1      nrKS(nshp,nalf,nsiz,nalt)/real(nrep)
255          c
256          c      ADpwr(nshp,nalf,nsiz,nalt) =
257          1      nrAD(nshp,nalf,nsiz,nalt)/real(nrep)
258          c
259          c      CVpwr(nshp,nalf,nsiz,nalt) =
260          1      nrCV(nshp,nalf,nsiz,nalt)/real(nrep)

```

```

261      c
262      X2pwr (nshp,nalf,nsiz,nalt) =
263      1      nrX2(nshp,nalf,nsiz,nalt)/real (nrep)
264      c
265      print*, '*****'
266      print*, ' POWER VALUES FROM MAIN PROGRAM'
267      print*, ' Null-hyp c =',c,'alpha =',alpha
268      print*, ' n=',n,' Alternate CDF: ',altcdf(nalt)
269      c      print*, '*****'
270      c      print*, ' KS Rejects = ',nrKS(nshp,nalf,nsiz,nalt)
271      c      print*, ' AD Rejects = ',nrAD(nshp,nalf,nsiz,nalt)
272      c      print*, ' CV Rejects = ',nrCV(nshp,nalf,nsiz,nalt)
273      c      print*, ' X2 Rejects = ',nrX2(nshp,nalf,nsiz,nalt)
274      c      print*, '*****'
275      print*, ' KS Power = ',KSpwr (nshp,nalf,nsiz,nalt)
276      print*, ' AD Power = ',ADpwr (nshp,nalf,nsiz,nalt)
277      print*, ' CV Power = ',CVpwr (nshp,nalf,nsiz,nalt)
278      print*, ' X2 Power = ',X2pwr (nshp,nalf,nsiz,nalt)
279      print*, '*****'
280      print*, ' '
281      c
282      60      continue
283      c
284      c      -- End DO Loop 60 for Alternate CDFs --
285      c      ** Completes Step 6 of Figure 7 **
286      c
287      c      -- Write Power Results to File --
288      c
289      write(7,110),n,test(1),KSpwr (nshp,nalf,nsiz,1),
290      1      KSpwr (nshp,nalf,nsiz,2),KSpwr (nshp,nalf,nsiz,3),
291      1      KSpwr (nshp,nalf,nsiz,4),KSpwr (nshp,nalf,nsiz,5),
292      1      KSpwr (nshp,nalf,nsiz,6),KSpwr (nshp,nalf,nsiz,7),
293      1      KSpwr (nshp,nalf,nsiz,8)
294      c
295      write(7,110),n,test(2),ADpwr (nshp,nalf,nsiz,1),
296      1      ADpwr (nshp,nalf,nsiz,2),ADpwr (nshp,nalf,nsiz,3),
297      1      ADpwr (nshp,nalf,nsiz,4),ADpwr (nshp,nalf,nsiz,5),
298      1      ADpwr (nshp,nalf,nsiz,6),ADpwr (nshp,nalf,nsiz,7),
299      1      ADpwr (nshp,nalf,nsiz,8)
300      c
301      write(7,110),n,test(3),CVpwr (nshp,nalf,nsiz,1),
302      1      CVpwr (nshp,nalf,nsiz,2),CVpwr (nshp,nalf,nsiz,3),
303      1      CVpwr (nshp,nalf,nsiz,4),CVpwr (nshp,nalf,nsiz,5),
304      1      CVpwr (nshp,nalf,nsiz,6),CVpwr (nshp,nalf,nsiz,7),
305      1      CVpwr (nshp,nalf,nsiz,8)
306      c
307      write(7,110),n,test(4),X2pwr (nshp,nalf,nsiz,1),
308      1      X2pwr (nshp,nalf,nsiz,2),X2pwr (nshp,nalf,nsiz,3),
309      1      X2pwr (nshp,nalf,nsiz,4),X2pwr (nshp,nalf,nsiz,5),
310      1      X2pwr (nshp,nalf,nsiz,6),X2pwr (nshp,nalf,nsiz,7),
311      1      X2pwr (nshp,nalf,nsiz,8)
312      c

```

```

313         write(7,72)
314     c
315     70     continue
316     c
317     c     --- End DO Loop 70 for Sample Sizes ---
318     c     ** Completes Step 7 of Figure 7 **
319     c
320     80     continue
321     c
322     c     --- End DO Loop 80 for Alpha Significance Levels ---
323     c     ** Completes Step 8 of Figure 7 **
324     c
325     write(7,74)
326     c
327     90     continue
328     c
329     c --- End DO Loop 90 for Null-Hypothesis Pareto Shape Parameter ---
330     c
331     c*****
332     c
333     c     Specify Format for Hardcopy Output Data and Headers:
334     c
335     51 format('1',36X,'Table XVII')
336     52 format('1',35X,'Table XVIII')
337     54 format(' ')
338     56 format('0',22X,'POWER TEST FOR THE PARETO DISTRIBUTION')
339     58 format(22X,'Ho: Pareto Distribution at Shape c = 1.0')
340     59 format(22X,'Ho: Pareto Distribution at Shape c = 3.5')
341     62 format(22X,'Ha: The data follow another distribution')
342     64 format('0',28X,'Level of Significance = .05')
343     66 format('0',28X,'Level of Significance = .01')
344     68 format(35X,'Alternate Distributions')
345     72 format(80(' '))
346     74 format(80('='))
347     76 format(2X,' n',3X,'Test',4X,'Par.1',3X,'Par.2',3X,'Par.3',3X,
348     1      'Weibl',3X,'Gamma',3X,'Beta',4X,
349     1      'Expon',3X,'Norm1')
350     110 format(' ',I3,A7,F9.3,7F8.3)
351     c
352     close(7)
353     c
354     end
355     c
356     c=====
357     c     END MAIN PROGRAM
358     c*****

```

```

359          Subroutine PARETO
360  c*****
361  c**
362  c**      BEGIN  SUBROUTINE  PARETO      **
363  c**
364  c*****
365  c
366  c  Ref:  Appendix B, Fig 7, Step 1.
367  c
368  c=====
369  c
370  c  Purpose:  For a specified sample size n, generate n random
371  c            deviates from a Pareto distribution with parameters of
372  c            location, scale, and shape set to specified positive
373  c            values.
374  c
375  c=====
376  c
377  c  Variables:
378  c            r = array containing n random numbers
379  c            ac = actual shape parameter of Pareto deviates
380  c            x = array containing n Pareto deviates
381  c            n = sample size
382  c            dseed = random number seed
383  c
384  c=====
385  c
386  c  Input:    dseed = random number seed (from MAIN program)
387  c            n = sample size = 5,15, or 25 (MAIN DO Loop 70)
388  c            nalt = alternate CDF counter (MAIN DO Loop 60)
389  c
390  c=====
391  c
392  c  IMSL Subroutines:
393  c
394  c  GGUBFS - generates random numbers distrib uniformly on (0,1)
395  c  VSRTA  - arranges a set of numbers in ascending order
396  c
397  c=====
398  c
399  c  Calculate:
400  c
401  c       $x(j) = (1/r(j)) ** (1/ac)$  for j = 1,2,...,n (from eqn 48)
402  c
403  c       $x(a',b') = b' * ((x(a,b) - a) / b) + a'$  (from eqn 50)
404  c
405  c=====
406  c
407  c  Output:   x = array of n random Pareto deviates
408  c
409  c=====
410  c

```

```

411 c Declare Variables:
412 c
413 common dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
414 1 Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nsho,nrep,
415 1 nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
416 integer n,nsiz,nsho,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
417 1 nrCV(2,2,3,8)
418 real x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
419 1 CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
420 1 P(25),r(25),alpha,KSpowr(2,2,3,8),ADpowr(2,2,3,8),
421 1 CVpowr(2,2,3,8),ac
422 double precision dseed
423 c
424 if (nalt .eq. 1) ac = 1.0
425 if (nalt .eq. 2) ac = 3.5
426 if (nalt .eq. 3) ac = 2.0
427 c
428 c--- Begin DO Loop 10 to Generate n Random Pareto Deviates ---
429 c
430 do 10 j = 1,n
431 c
432 c Use IMSL subroutine to generate random numbers:
433 r(j) = GGUBFS(dseed)
434 c
435 c Use eqn 48 to transform them to Pareto deviates
436 c with location a = 1 and scale b = 1:
437 x(j) = (1.0/r(j))**(1.0/ac)
438 c
439 c Use eqn 50 to transform to Pareto deviates with
440 c a = 2, b = 3 for the second alternate CDF:
441 if (nalt .eq. 2) x(j) = 3. * x(j) - 1.
442 c
443 c Use eqn 50 to transform to Pareto deviates with
444 c a = 10, b = 5 for the third alternate CDF:
445 if (nalt .eq. 3) x(j) = 5. * x(j) + 5.
446 c
447 10 continue
448 c
449 c--- End DO Loop 10 after Generating n Random Deviates ---
450 c
451 return
452 end
453 c
454 c=====
455 c END SUBROUTINE PARETO
456 c*****

```

```

457           Subroutine BXVALS
458 c*****
459 c**
460 c**           B E G I N   S U B R O U T I N E   B X V A L S           **
461 c**
462 c*****
463 c
464 c Ref: Appendix B, Fig. 7, Step 2.
465 c
466 c=====
467 c
468 c Purpose: For a given sample size n, calculate the B values
469 c           used to find the BLUEs of location and scale. Also
470 c           find the sum of the first n-1 values of B(i). Then,
471 c           compute the three values equal to the sums of the
472 c           first n-1, the first n-2, and (for hypothesized
473 c           c = .5, 1, or 2) the first n -2/c values of B(i)x(i).
474 c
475 c=====
476 c
477 c Variables:  c = null-hypothesis shape parameter
478 c              n = sample size
479 c              x = array containing n ordered deviates
480 c                 from an alternate distribution
481 c              B = array containing n values of B
482 c              Bsum1 = sum of B(i) values for i = 1,2,...,(n-1)
483 c              Bxsum1 = sum of B(i)x(i) for i = 1,2,...,(n-1)
484 c              Bxsum2 = sum of B(i)x(i) for i = 1,2,...,(n-2)
485 c              Bxsm2c = sum of B(i)x(i) for i = 1,2,...,(n-2/c)
486 c
487 c=====
488 c
489 c Input:      c = null-hyp shape parameter (from MAIN DO Loop 90)
490 c              n = sample size = 5, 15, or 25 (from MAIN DO Loop 70)
491 c              nc = n*c (from MAIN program)
492 c              x = ordered deviates of alternate CDF MAIN)
493 c
494 c=====
495 c
496 c Calculate:
497 c
498 c           B(i) = [1 - 2/c(n-i+1)] * B(i-1)           (eqn 29)
499 c
500 c           Bsum1 = B(1) + B(2) + ... + B(n-1)
501 c
502 c           Bxsum1 = B(1)*x(1) + ... + B(n-1)*x(n-1)
503 c
504 c           Bxsum2 = B(1)*x(1) + ... + B(n-2)*x(n-2)
505 c
506 c           Bxsm2c = B(1)*x(1) + ... + B(n-2/c)*x(n-2/c)
507 c
508 c=====

```

```

509      c
510      c   Output:
511      c           B = array containing n values of B
512      c           Bsum1 = sum of first (n-1) B values
513      c           Bxsum1 = sum of first (n-1) B*x values
514      c           Bxsum2 = sum of first (n-2) B*x values
515      c           Bxsm2c = sum of first (n-2/c) B*x (if 2/c is integer)
516      c
517      c=====
518      c
519      c   Declare Variables:
520      c
521      c           common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
522      c           1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
523      c           1      nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
524      c           integer n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
525      c           1      nrCV(2,2,3,8)
526      c           real    x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
527      c           1      CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
528      c           1      P(25),r(25),alpha,KSpr(2,2,3,8),ADpwr(2,2,3,8),
529      c           1      CVpwr(2,2,3,8)
530      c           double precision dseed
531      c
532      c   Calculate the first B value (eqn 25):
533      c
534      c           B(1) = 1.0 - 2.0/nc
535      c
536      c   --- Begin DO Loop 10 to Find the 2nd thru nth B values ---
537      c
538      c           do 10 j = 2,n
539      c               B(j) = B(j-1) * (1.0 - (2.0/(c*(n-j+1))))
540      c           10  continue
541      c
542      c   --- End DO Loop 10 ---
543      c
544      c           Bsum1 = 0
545      c
546      c   --- Begin DO Loop 20 to Sum the First n-1 Values of B ---
547      c
548      c           do 20 k=1,(n-1)
549      c               Bsum1 = Bsum1 + B(k)
550      c           20  continue
551      c
552      c   --- End DO Loop 20 ---
553      c
554      c           Bxsum1 = 0
555      c
556      c   --- Begin DO Loop 30 to Sum the First n-1 Values of Bx ---
557      c
558      c           do 30 l=1,(n-1)
559      c               Bxsum1 = Bxsum1 + (B(l))*x(l)
560      c           30  continue

```

```

561      c
562      c      --- End DO Loop 30 ---
563      c
564      c      Bxsum2 = Bxsum1 - (B(n-1)*x(n-1))
565      c
566      c      --- Find Bxsm2c When 2/c is an Integer (c=.5, 1, or 2) ---
567      c
568      c      Bxsm2c = 0
569      c
570      c      if (c .eq. 1.0) then
571      c          Bxsm2c = Bxsum2
572      c      else if (c .eq. 2.0) then
573      c          Bxsm2c = Bxsum1
574      c      else if (c .eq. 0.5) then
575      c          Bxsm2c = Bxsum2 - (B(n-3)*x(n-3)) - (B(n-2)*x(n-2))
576      c      end if
577      c
578      c      return
579      c      end
580      c
581      c=====
582      c          END SUBROUTINE BXVALS
583      c*****

```

```

584           Subroutine BLCLE2
585 c*****
586 c**
587 c**           BEGIN   SUBROUTINE   BLCLE2           **
588 c**
589 c*****
590 c
591 c Ref: Appendix B, Figure 7, Step 2 (continued).
592 c
593 c-----
594 c
595 c Purpose: Given an ordered sample of size n and null-hypothesis
596 c           c<=2, calculate the BLUEs of location a and scale b.
597 c
598 c-----
600 c Variables:
601 c           x = array containing n ordered deviates from a CDF
602 c           c = null-hypothesis Pareto shape parameter
603 c           n = sample size
604 c           B = array of B values used to calculate the BLUEs
605 c           nc = product of n and c
606 c           Coef1 = coefficient used to compute BLUE of location a
607 c           Coef2 = coefficient used to compute BLUE of location a
608 c           Coef3 = coefficient used to compute BLUE of scale b
609 c           Bxsum2 = sum of B(i)*x(i) terms for i = 1,...,n-2
610 c           Bxsm2c = sum of B(i)*x(i) terms for i = 1,...,n-2/c
611 c           ablu = BLUE of the location parameter a
612 c           bblu = BLUE of the scale parameter b
613 c           U = value used to compute BLUEs when c = 1.5
614 c           Termi = terms used to compute U (i=1,2,3)
615 c-----
616 c Input:   x = array of n ordered deviates (from MAIN Program)
617 c           c = null-hyp shape = 1.0 (from MAIN DO Loop 90)
618 c           n = sample size = 5, 15, or 25 (from MAIN DO Loop 70)
619 c           nc = n*c (from MAIN program)
620 c           B = array containing n values of B (from BXVALS)
621 c           Bxsum2 = sum of first n-2 values of B (from BXVALS)
622 c           Bxsm2c = sum of first n-2/c values of B (from BXVALS)
623 c-----
624 c Calculate (if c = 0.5, 1.0, or 2.0):
625 c
626 c           Coef1 = [(c+1)*(c+2)] / [(nc-2)*(nc-c-2)]
627 c           Coef2 = (nc-2) / (c+2)
628 c           ablu = x(1) - Coef1 * [Bxsm2c - (Coef2*x(1))] (eqn 34)
629 c           bblu = (nc-1) * [x(1) - ablu] (eqn 35)
630 c-----
631 c Calculate (if c = 1.5):
632 c
633 c           Term1 = (nc-2) * (nc-c-2)
634 c           Term2 = nc * (c-2) * B(n-1)
635 c           Term3 = (nc-1) * (c+2)

```

```

643      c          Coef3 = [(nc-1)/nc] * (nc-2-U)
644      c          U = (Term1 - Term2) / Term3      (eqn 39)
645      c
646      c          ablu = x(1) - bblu / (nc-1)      (eqn 37)
647      c          bblu = (1/U) * [(c+1)*(Bxsum2) + (2c-1)*B(n-1)*x(n-1)
648      c          - Coef3 * x(1)]      (eqn 38)
650      c=====
652      c  Output:
653      c          ablu = BLUE of location parameter a
654      c          bblu = BLUE of scale parameter b
656      c=====
657      c
658      c  Declare Variables:
659      c
660      c          common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
661      1          Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
662      1          nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
663      c          integer n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
664      1          nrCV(2,2,3,8)
665      c          real    x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
666      1          CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
667      1          P(25),r(25),alpha,KSpwr(2,2,3,8),ADpwr(2,2,3,8),
668      1          CVpwr(2,2,3,8),Term1,Term2,Term3,Coef1,Coef2,
669      1          Coef3,U
670      c          double precision dseed
671      c
672      c          if ((c.eq.0.5) .or. (c.eq.1.0) .or. (c.eq.2.0)) then
673      c              Coef1 = ((c+1.0)*(c+2.0)) / ((nc-2.0)*(nc-c-2.0))
674      c              Coef2 = (nc-2.0) / (c+2.0)
675      c              ablu = x(1) - Coef1 * (Bxsm2c - (Coef2*x(1)))
676      c              bblu = (nc-1.0) * (x(1) - ablu)
677      c
678      c          else if (c .eq. 1.5) then
679      c              Term1 = (nc-2.0) * (nc-c-2.0)
680      c              Term2 = nc * (c-2.0) * B(n-1)
681      c              Term3 = (nc-1.0) * (c+2.0)
682      c              U = (Term1 - Term2) / Term3
683      c              Coef3 = ((nc-1.0)/nc) * (nc-2.0-U)
684      c              bblu = (1.0/U) * ( (c+1.0) * (Bxsum2)
685      1              + (2.0*c-1.0)*B(n-1)*x(n-1) - Coef3 * x(1) )
686      c              ablu = x(1) - (bblu / (nc-1.0))
687      c
688      c          end if
689      c
690      c          return
691      c          end
692      c
693      c=====
694      c          END SUBROUTINE BLCLE2
695      c*****

```

```

696           Subroutine BLCGT2
697 c*****
698 c**
699 c**       B E G I N   S U B R O U T I N E   B L C G T 2       **
700 c**
701 c*****
702 c
703 c   Ref:  Appendix B, Figure 7, Step 2 (continued).
704 c
705 c=====
706 c
707 c   Purpose:  Given an ordered sample of size n and a Pareto null
708 c             hypothesis with shape c > 2, calculate the best
709 c             linear unbiased estimates (BLUEs) of location and
710 c             scale.
711 c
712 c=====
713 c
714 c   Variables:  x = array containing n ordered deviates
715 c              c = null-hypothesis Pareto shape parameter
716 c              n = sample size
717 c              nc = product of n and c
718 c              B = array of B values used to calculate the BLUEs
719 c              Bsum1 = sum of B(i) terms for i = 1,...,n-1
720 c              Bxsum1 = sum of B(i)*x(i) terms for i = 1,...,n-1
721 c              D = value used to calculate the BLUEs
722 c              YV = value used to calculate the BLUEs
723 c              ablu = BLUE for location parameter a
724 c              bblu = BLUE for scale parameter b
725 c
726 c=====
727 c
728 c   Input:      x = array of ordered deviates (from MAIN Program)
729 c              c = shape parameter = 3.5 (from MAIN DO Loop 90)
730 c              n = sample size = 5, 15, or 25 (MAIN DO Loop 70)
731 c              nc = n*c (from MAIN Program)
732 c              B = array of B values (from BXVALS)
733 c              Bsum1 = sum of first (n-1) B values (from BXVALS)
734 c              Bxsum1 = sum of first n-1 B*x values (from BXVALS)
735 c
736 c=====
737 c
738 c   Calculate:
739 c
740 c           D = [(c+1) * Bsum1] + [(c-1) * B(n)]           (eqn 21)
741 c
742 c           YV = (c+1)*Bxsum1 + (c-1)*B(n)*x(n) - D*x(1)  (eqn 22)
743 c
744 c           ablu = x(1) - YV/[(nc-1)*(nc-2) - D*nc]        (eqn 17)
745 c
746 c           bblu = (nc-1) * [ x(1) - ablu ]                (eqn 18)
747 c

```

```

748 c=====
749 c
750 c Output:   ablu = BLUE for location a
751 c          bblu = BLUE for scale b
752 c
753 c=====
754 c
755 c Declare Variables:
756 c
757 c      common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
758 1            Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
759 1            nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
760 integer  n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
761 1            nrCV(2,2,3,8)
762 real     x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
763 1            CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
764 1            P(25),r(25),alpha,KSpwr(2,2,3,8),ADpwr(2,2,3,8),
765 1            CVpwr(2,2,3,8),YV
766 double precision dseed
767 c
768 c      D = ((c+1.0) * Bsum1) + ((c-1.0) * B(n))
769 c      YV = ((c+1.0)*Bxsum1) + ((c-1.0)*B(n)*x(n)) - (D*x(1))
770 c      ablu = x(1) - YV/((nc-1.0)*(nc-2.0) - (D*nc))
771 c      bblu = (nc-1.0) * (x(1) - ablu)
772 c
773 c      return
774 c      end
775 c
776 c=====
777 c      END SUBROUTINE BLCGT2
778 c*****

```

```

779           Subroutine HYPCDF
780 c*****
781 c**
782 c**      B E G I N      S U B R O U T I N E      H Y P C D F      **
783 c**
784 c*****
785 c
786 c  Ref:  Appendix B, Figure 7, Step 2 (continued).
787 c
788 c-----
789 c
790 c  Purpose:  Given an ordered sample of size n, a Pareto null-hyp
791 c            of shape c, and the BLUEs of location a and scale b,
792 c            compute the hypothesized Pareto distribution
793 c            function P(i) for i = 1,2,...,n.
794 c
795 c-----
796 c
797 c  Variables:
798 c            x = array containing n ordered deviates
799 c            n = sample size
800 c            c = null hypothesized Pareto shape parameter
801 c            ablu = BLUE of location a
802 c            bblu = BLUE of scale b
803 c            P = array containing n points of the
804 c                hypothesized Pareto CDF
805 c
806 c-----
807 c
808 c  Input:
809 c            x = array of n ordered deviates (from MAIN Program)
810 c            c = null hyp shape = 1.0 or 3.5 (MAIN DO Loop 90)
811 c            n = sample size = 5, 15, or 25 (from MAIN DO Loop 70)
812 c            ablu = BLUE of location a (from BLCLE2 or BLCGT2)
813 c            bblu = BLUE of scale b (from BLCLE2 or BLCGT2)
814 c
815 c-----
816 c
817 c  Calculate:
818 c
819 c      P(i) = 1 - [1 + (x(i) - ablu)/bblu] ** (-c)      (eqn 40)
820 c
821 c-----
822 c
823 c  Output:  P = array of n points of the hypothesized CDF
824 c
825 c-----
826 c
827 c  Declare Variables:
828 c
829 c      common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
830 c      1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nsnp,nrep,

```

```

831      1      nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
832      integer n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
833      1      nrCV(2,2,3,8)
834      real    x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
835      1      CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
836      1      P(25),r(25),alpha,KSpwr(2,2,3,8),ADpwr(2,2,3,8),
837      1      CVpwr(2,2,3,8)
838      double precision dseed
839      c
840      do 10 i = 1,n
841          P(i) = 1.0 - (1.0 + (x(i) - ablu)/bblu) ** (-c)
842      10  continue
843      c
844      return
845      end
846      c
847      c=====
848      c                      END SUBROUTINE HYPcdf
849      c*****

```

```

850           Subroutine TESTAT
851 c*****
852 c**
853 c**           B E G I N   S U B R O U T I N E   T E S T A T           **
854 c**
855 c*****
856 c
857 c  Ref:  Appendix B, Figure 7, Step 2.
858 c
859 c=====
860 c
861 c  Purpose:  Given a sample size n, and the hypothesized Pareto
862 c            distribution function P(i), compute values of the
863 c            test statistics of the Chi-square and the modified
864 c            K-S, A-D, and CVM goodness-of-fit tests.
865 c
866 c=====
867 c
868 c  Variables:
869 c            n = sample size
870 c            nshp = null-hyp shape counter (1: c=1.0, 2: c=3.5)
871 c            nalf = alpha level counter (1:  $\alpha$  =.05, 2:  $\alpha$  =.01)
872 c            nsiz = sample size counter (1: n=5, 2: n=15, 3: n=25)
873 c            nalt = alternate distribution counter
874 c            P = array of n values of the hypothesized Pareto CDF
875 c
876 c            DP = positive differences between EDF and CDF points
877 c            DM = negative differences between EDF and CDF points
878 c            DPLUS = maximum positive difference (largest DP value)
879 c            DMINUS = maximum negative difference (largest DM value)
880 c            KS = values of the modified K-S test statistic
881 c
882 c            AL = value used to calculate the A-D test statistic
883 c            AM = value used to calculate the A-D test statistic
884 c            AN = AL + AM
885 c            AAA = values to be summed for A-D test statistic
886 c            SAAA = sum of AAA values
887 c            AD = values of the modified A-D test statistic
888 c
889 c            ACV = squared quantities in the C-VM formula
890 c            SACV = sum of the ACV values
891 c            CVM = values of the modified C-VM test statistic
892 c
893 c            ablu = BLUE of location parameter a
894 c            bblu = BLUE of scale parameter b
895 c            c = null-hypothesized Pareto shape parameter
896 c            obs = number of observations in each of 5 cells
897 c            rtend = right endpoint of a cell
898 c            X2 = array of values of the Chi-square test statistic
899 c
900 c=====
901 c

```

```

902 c Input:
903 c     n = sample size = 5, 15, or 25 (from MAIN DO Loop 70)
904 c     P = array of n values of hypothesized CDF (from HYPCDF)
905 c     nshp = null-hyp shape counter (from MAIN DO Loop 90)
906 c     nalf = significance level counter (from MAIN DO Loop 80)
907 c     nsiz = sample size counter (from MAIN DO Loop 70)
908 c     nalt = alternate CDF counter (from MAIN DO Loop 60)
909 c     ablu = BLUE of location a (from BLCLE2 or BLCGT2)
910 c     bblu = BLUE of scale b (from BLCLE2 or BLCGT2)
911 c     c = hypothesized Pareto shape (from MAIN DO Loop 90)
912 c
913 c -----
914 c
915 c Calculations for K-S test statistic (eqns 41 & 42):
916 c
917 c     DP(i) = ABS[ (i/n) - P(i) ]
918 c     DM(i) = ABS[ P(i) - (i-1)/n ]
919 c
920 c     DPLUS = max [ DP(i) ] for i=1,2,...,n
921 c     DMINUS = max [ DM(i) ] for i=1,2,...,n
922 c
923 c     KS = max (DPLUS,DMINUS)
924 c
925 c -----
926 c
927 c Calculations for A-D test statistic (eqn 43):
928 c
929 c     AL(j) = ln (P(j))
930 c     AM(j) = ln (1 - P(n+1-j))
931 c     AN(j) = AL(j) + AM(j)
932 c
933 c     AAA(j) = (2*j - 1) * AN(j)
934 c     SAAA = AAA(1) + AAA(2) + ... + AAA(n)
935 c
936 c     AD = -n - (1/n) * SAAA
937 c
938 c -----
939 c
940 c Calculations for C-VM test statistic (eqn 44):
941 c
942 c     ACV(k) = [ P(k) - (2*k - 1)/(2*n) ]**2
943 c     SACV = ACV(1) + ACV(2) + ... + ACV(n)
944 c
945 c     CVM = (1/(12*n)) + SACV
946 c
947 c -----
948 c
949 c Calculations for Chi-square test statistic (eqn 62):
950 c
951 c     rtend(i) = ablu - bblu + bblu * (1 - .2*i) ** (-1/c)
952 c     ex = n / 5.
953 c

```

```

954      c      X2 = [(obs(1)-ex)**2] / ex + [(obs(2)-ex)**2] / ex
955      c      + ... + [(obs(5)-ex)**2] / ex
956      c
957      c-----
958      c
959      c  Declare Variables:
960      c
961      c      common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
962      1          Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
963      1          nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
964      c      integer n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
965      1          nrCV(2,2,3,8),obs(5),nrX2(2,2,3,8)
966      c      real    x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
967      1          CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
968      1          P(25),r(25),alpha,KSpwr(2,2,3,8),ADpwr(2,2,3,8),
969      1          CVpwr(2,2,3,8),DP(25),DM(25),DPLUS,DMINUS,AL(25),
970      1          AM(25),AN(25),AAA(25),SAAA,ACV(25),SACV,rtend(4),
971      1          X2crit(2,2,3),X2(2,2,3,8),ex
972      c      double precision dseed
973      c
974      c -----  Compute the K-S Test Statistic (eqns 41 & 42):  -----
975      c
976      c          DPLUS = 0
977      c          DMINUS = 0
978      c          do 5 ik = 1,25
979      c              DP(ik) = 0
980      c              DM(ik) = 0
981      5          continue
982      c
983      c          do 10 i = 1,n
984      c
985      c              DP(i) = ABS( (i/real(n)) - P(i) )
986      c              DM(i) = ABS( P(i) - (i-1)/real(n) )
987      c
988      c              if (nshp.eq.1 .and. nalf.eq.2 .and. n.eq.5 .and.
989      c 1          nalt .lt. 3) then
990      c              print*, 'P(i)=',P(i), 'DP(i)=',DP(i), 'DM(i)=',DM(i)
991      c              end if
992      c
993      10          continue
994      c
995      c          DPLUS = MAX( DP(1),DP(2),DP(3),DP(4),DP(5),DP(6),DP(7),
996      1          DP(8),DP(9),DP(10),DP(11),DP(12),DP(13),DP(14),
997      1          DP(15),DP(16),DP(17),DP(18),DP(19),DP(20),
998      1          DP(21),DP(22),DP(23),DP(24),DP(25) )
999      c
1000      c          DMINUS = MAX( DM(1),DM(2),DM(3),DM(4),DM(5),DM(6),DM(7),
1001      1          DM(8),DM(9),DM(10),DM(11),DM(12),DM(13),DM(14),
1002      1          DM(15),DM(16),DM(17),DM(18),DM(19),DM(20),
1003      1          DM(21),DM(22),DM(23),DM(24),DM(25) )
1004      c
1005      c          KS(nshp,nalf,nsiz,nalt) = MAX(DPLUS,DMINUS)

```





```

1095           Subroutine COMPAR
1096 c*****
1097 c**
1098 c**           B E G I N   S U B R O U T I N E   C O M P A R           **
1099 c**
1100 c*****
1101 c
1102 c   Ref: Appendix B, Figure 7, Step 3.
1103 c
1104 c=====
1105 c
1106 c   Purpose:
1107 c
1108 c       Compare a test statistic, calculated from Chi-square or the
1109 c       modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D),
1110 c       or Cramer-von Mises (C-VM) test, against the appropriate
1111 c       critical value. From a series of test statistics, count the
1112 c       number of times the null hypothesis is rejected, i.e., the
1113 c       number of test statistic values that exceed the critical
1114 c       value. The K-S, A-D, and C-VM critical values were taken
1115 c       from Tables VI- VIII of the thesis.
1116 c
1117 c=====
1118 c
1119 c   Variables:
1120 c       c = null-hypothesis Pareto shape parameter
1121 c       alpha = significance level
1122 c       n = sample size
1123 c       nshp = shape parameter counter (1: c=1.0; 2: c=3.5)
1124 c       nalf = significance level counter (1:  $\alpha=.05$ ; 2:  $\alpha=.01$ )
1125 c       nsiz = sample size counter (1: n=5; 2: n=15; 3: n=25)
1126 c       KS = array of modified K-S test statistics
1127 c       CVM = array of modified C-VM test statistics
1128 c       AD = array of modified A-D test statistics
1129 c       X2 = array of Chi-square test statistics
1130 c
1131 c=====
1132 c
1133 c   Input:
1134 c       c = null-hyp shape parameter (from MAIN DO Loop 90)
1135 c       alpha = significance level (from MAIN DO Loop 80)
1136 c       n = sample size (from MAIN DO Loop 80)
1137 c       nshp = shape parameter counter (from MAIN DO Loop 90)
1138 c       nalf = significance level counter (MAIN DO Loop 80)
1139 c       nsiz = sample size counter (from MAIN DO Loop 70)
1140 c       nalt = alternate CDF counter (from MAIN DO Loop 60)
1141 c       KS = array of K-S test statistics (from TESTAT)
1142 c       CVM = array of C-VM test stats (from TESTAT)
1143 c       AD = array of A-D test statistics (from TESTAT)
1144 c       KScrit(nshp,nalf,nsiz) = K-S critical values (Table VI)
1145 c       ADcrit(nshp,nalf,nsiz) = A-D critical values (Table VII)
1146 c       CVcrit(nshp,nalf,nsiz) = CVM critical values (Table VIII)

```

```

1147 c      X2crit(nshp,nalf,nsiz) = Chi-square critical values
1148 c
1149 c=====
1150 c
1151 c      Calculations: none
1152 c
1153 c=====
1154 c
1155 c      Output:
1156 c
1157 c      nrKS = number of times null hypothesis is rejected under K-S
1158 c      nrAD = number of times null hypothesis is rejected under A-D
1159 c      nrCV = number of times null hypothesis is rejected under CVM
1160 c      nrX2 = number of times null hyp is rejected under Chi-square
1161 c
1162 c=====
1163 c
1164 c      Declare Variables:
1165 c
1166 c      common  dseed,x,n,c,nc,B,D,ablu,bblu,P,Bsum1,Bxsum1,
1167 c      1      Bxsum2,Bxsm2c,KS,AD,CVM,it,nsiz,nshp,nrep,
1168 c      1      nalt,nalf,nrKS,nrAD,nrCV,nrX2,X2
1169 c      integer n,nsiz,nshp,it,nrep,nrKS(2,2,3,8),nrAD(2,2,3,8),
1170 c      1      nrCV(2,2,3,8),nrX2(2,2,3,8)
1171 c      real    x(25),ablu,bblu,B(25),D,KS(2,2,3,8),AD(2,2,3,8),
1172 c      1      CVM(2,2,3,8),c,nc,Bsum1,Bxsum1,Bxsum2,Bxsm2c,
1173 c      1      P(25),r(25),alpha,KSowr(2,2,3,8),ADpwr(2,2,3,8),
1174 c      1      CVpwr(2,2,3,8),KScrit(2,2,3),ADcrit(2,2,3),
1175 c      1      CVcrit(2,2,3),X2crit(2,2,3),X2(2,2,3,8)
1176 c      double precision dseed
1177 c
1178 c      print*, '*****'
1179 c      print*, 'Numbers of Rejects at COMPAR Entrance'
1180 c      print*, 'c =',c,'nalf =',nalf,'n =',n,'nalt =',nalt
1181 c      print*, 'KS Rejects = ',nrKS(nshp,nalf,nsiz,nalt)
1182 c      print*, 'AD Rejects = ',nrAD(nshp,nalf,nsiz,nalt)
1183 c      print*, 'CV Rejects = ',nrCV(nshp,nalf,nsiz,nalt)
1184 c      print*, '=====
1185 c
1186 c      --- Input K-S Critical Values from Table VI: ---
1187 c
1188 c      KScrit(1,1,1) = .3676251
1189 c      KScrit(1,1,2) = .2157919
1190 c      KScrit(1,1,3) = .1698559
1191 c      KScrit(1,2,1) = .4074441
1192 c      KScrit(1,2,2) = .2468265
1193 c      KScrit(1,2,3) = .2007451
1194 c      KScrit(2,1,1) = .3493998
1195 c      KScrit(2,1,2) = .2376525
1196 c      KScrit(2,1,3) = .1886063
1197 c      KScrit(2,2,1) = .3815996
1198 c      KScrit(2,2,2) = .2743093

```

```

1199          KScrit(2,2,3) = .2182668
1200          c
1201          c --- Input A-D Critical Values from Table VII: ---
1202          c
1203          ADcrit(1,1,1) = 1.236920
1204          ADcrit(1,1,2) = .8907447
1205          ADcrit(1,1,3) = .9147376
1206          ADcrit(1,2,1) = 2.076011
1207          ADcrit(1,2,2) = 1.250242
1208          ADcrit(1,2,3) = 1.311781
1209          ADcrit(2,1,1) = .6840515
1210          ADcrit(2,1,2) = .8985860
1211          ADcrit(2,1,3) = .9520599
1212          ADcrit(2,2,1) = .9126385
1213          ADcrit(2,2,2) = 1.268849
1214          ADcrit(2,2,3) = 1.449695
1215          c
1216          c --- Input C-VM Critical Values from Table VIII: ---
1217          c
1218          CVcrit(1,1,1) = .1389776
1219          CVcrit(1,1,2) = .1312229
1220          CVcrit(1,1,3) = .1386932
1221          CVcrit(1,2,1) = .1738497
1222          CVcrit(1,2,2) = .1923594
1223          CVcrit(1,2,3) = .1988135
1224          CVcrit(2,1,1) = .1186844
1225          CVcrit(2,1,2) = .1561372
1226          CVcrit(2,1,3) = .1618638
1227          CVcrit(2,2,1) = .1574178
1228          CVcrit(2,2,2) = .2217665
1229          CVcrit(2,2,3) = .2403474
1230          c
1231          c --- Input Chi-square Critical Values : ---
1232          c
1233          X2crit(1,1,1) = 6.000003
1234          X2crit(1,1,2) = 7.333337
1235          X2crit(1,1,3) = 7.600005
1236          X2crit(1,2,1) = 12.00000
1237          X2crit(1,2,2) = 10.66667
1238          X2crit(1,2,3) = 10.80000
1239          X2crit(2,1,1) = 6.000003
1240          X2crit(2,1,2) = 7.333337
1241          X2crit(2,1,3) = 7.600005
1242          X2crit(2,2,1) = 6.000003
1243          X2crit(2,2,2) = 10.46378
1244          X2crit(2,2,3) = 10.80000
1245          c
1246          c --- Compare Test Statistics vs Critical Values: ---
1247          c
1248          c print*, '*****'
1249          c print*, 'BEFORE REJ COUNTER IS INCREMENTED:'
1250          c print*, 'c =',c,'nalf =',nalf,' ** n=',n,' ** nalt=',nalt

```

```

1251      c
1252      c      print*, 'KS Stat =', KS(nshp, nalf, nsiz, nalt),
1253      c      1      ' Crit =', KScrit(nshp, nalf, nsiz)
1254      c
1255      c      print*, 'AD Stat =', AD(nshp, nalf, nsiz, nalt),
1256      c      1      ' Crit =', ADcrit(nshp, nalf, nsiz)
1257      c
1258      c      print*, 'CV Stat =', CVM(nshp, nalf, nsiz, nalt),
1259      c      1      ' Crit =', CVcrit(nshp, nalf, nsiz)
1260      c
1261      c      print*, 'X2 Stat =', X2(nshp, nalf, nsiz, nalt),
1262      c      1      ' Crit =', X2crit(nshp, nalf, nsiz)
1263      c      print*, '*****'
1264      c
1265      c      if ( KS(nshp, nalf, nsiz, nalt) .gt. KScrit(nshp, nalf, nsiz) )
1266      c      1      nrKS(nshp, nalf, nsiz, nalt) = nrKS(nshp, nalf, nsiz, nalt) + 1
1267      c
1268      c      if ( AD(nshp, nalf, nsiz, nalt) .gt. ADcrit(nshp, nalf, nsiz) )
1269      c      1      nrAD(nshp, nalf, nsiz, nalt) = nrAD(nshp, nalf, nsiz, nalt) + 1
1270      c
1271      c      if ( CVM(nshp, nalf, nsiz, nalt) .gt. CVcrit(nshp, nalf, nsiz) )
1272      c      1      nrCV(nshp, nalf, nsiz, nalt) = nrCV(nshp, nalf, nsiz, nalt) + 1
1273      c
1274      c      if ( X2(nshp, nalf, nsiz, nalt) .gt. X2crit(nshp, nalf, nsiz) )
1275      c      1      nrX2(nshp, nalf, nsiz, nalt) = nrX2(nshp, nalf, nsiz, nalt) + 1
1276      c
1277      c      print*, '======'
1278      c      print*, 'Numbers of Rejects at COMPAR Exit'
1279      c      print*, 'c =', c, 'nalf =', nalf, ' n =', n, ' nalt =', nalt
1280      c      print*, 'KS Rejects = ', nrKS(nshp, nalf, nsiz, nalt)
1281      c      print*, 'AD Rejects = ', nrAD(nshp, nalf, nsiz, nalt)
1282      c      print*, 'CV Rejects = ', nrCV(nshp, nalf, nsiz, nalt)
1283      c      print*, 'X2 Rejects = ', nrX2(nshp, nalf, nsiz, nalt)
1284      c      print*, '======'
1285      c
1286      c      return
1287      c      end
1288      c
1289      c=====
1290      c      END SUBROUTINE COMPAR
1291      c*****

```

## BIBLIOGRAPHY

1. Amstatter, B. Reliability Mathematics. New York: McGraw-Hill Book Company, 1971.
2. Anderson, T. W. and D. A. Darling. "Asymptotic Theory of Goodness of Fit Criteria Based on Stochastic Processes," Annals of Mathematical Statistics, 23: 193-212 (1952).
3. Anderson, T. W. and D. A. Darling. "A Test of Goodness of Fit," Journal of the American Statistical Association, 49: 765-769 (Dec 1954).
4. Andrews, D. F. and others. Robust Estimates of Location. Princeton University Press, 1972.
5. Banks, Jerry and John S. Carson. Discrete-Event System Simulation. Englewood Cliffs: Prentice-Hall, 1984.
6. Bell, C. B. and others. Signal Detection for Pareto Renewal Processes. Technical Report No. 8-82 for the Office of Naval Research. Contract N00014-80-C-0208. San Diego State University, San Diego CA, Oct 1982 (AD-A120 972).
7. Berger, J. M. and B. Mandelbrot. "A New Model for Error Clustering in Telephone Circuits," IBM Journal of Research and Development, 7: 224-236 (July 1963).
8. Brownlee, K. A. Statistical Theory and Methodology in Science and Engineering (Second Edition). New York: John Wiley and Sons, 1965.
9. Bush, J. G. and others. "Modified Cramer-von Mises and Anderson-Darling Tests for Weibull Distributions with Unknown Location and Scale Parameters," Communications in Statistics, Part A - Theory and Methods, 12: 240-245 (1983).
10. Buslenko, N. P., and others. The Monte Carlo Method. New York: Pergamon Press, 1966.
11. Champernowne, D. G. "The Graduation of Income Distributions," Econometrica, 20: 591-615 (1952).
12. Charek, Dennis J. A Comparison of Estimation Techniques for the Three-Parameter Pareto Distribution. MS Thesis, GSO/MA/85D-3. School of Engineering, Air Force Institute of Technology (AU), Wright Patterson AFB OH, December 1985.

13. Conover, W. J. Practical Nonparametric Statistics (Second Edition). New York: John Wiley and Sons, 1980.
14. David, F. N. and N. L. Johnson. "The Probability Integral Transformation When Parameters are Estimated from the Sample," Biometrika, 35: 182-190 (1948).
15. David, Herbert A. Order Statistics (Second Edition). New York: John Wiley and Sons, 1981.
16. Davis, Henry T. and Michael L. Feldstein. "The Generalized Pareto Law as a Model for Progressively Censored Survival Data," Biometrika, 66: 299-306 (1979).
17. Fisk, P. R. "The Graduation of Income Distributions," Econometrica, 29: 171-185 (1961).
18. Freiling, E. C. A Comparison of the Fallout Mass-Size Distributions Calculated by Lognormal and Power-Law Models. Report No. USNRDL-TR-1105 for the U.S. Naval Radiological Defense Laboratory, San Francisco CA, Nov 1966 (AD-646019).
19. Green, J. and Y. Hegazy. "Powerful Modified EDF Goodness-of-Fit Tests," Journal of the American Statistical Association, 71: 204-209 (1976).
20. Hajek, Jaroslav. A Course in Non-Parametric Statistics. San Francisco: Holden-Day, Inc., 1969.
21. Hammersley, J. M. and D. C. Handscomb. Monte Carlo Methods. London: Methuen and Co., 1967.
22. Harris, Carl M. "The Pareto Distribution as a Queue Service Discipline," Operations Research, 16: 307-313 (Jan-Feb 1968).
23. Harter, H. L. Order Statistics and Their Use in Testing and Estimation, Vol 2. Aerospace Research Laboratories, Wright-Patterson AFB OH, 1969.
24. Harter, H. L. "Another Look at Plotting Positions," Communications in Statistics, A13(13): 1613-1633 (1984).
25. Harter, H. L. "A Monte Carlo Study of Plotting Positions," Communications in Statistics, B14(2): 317-343 (1985).
26. Hastings, N. A. J. and J. B. Peacock. Statistical Distributions. London: Butterworth & Co. Ltd., 1974.

27. Hines, William W. and Douglas C. Montgomery. Probability and Statistics in Engineering and Management Science. New York: The Ronald Press Co., 1972.
28. Johnson, Norman L. and Samuel Kotz. Continuous Univariate Distributions-1. Boston: Houghton Mifflin Co., 1970.
29. Kaminsky, Kenneth S. Best Linear Unbiased Prediction of Order Statistics in Exponential and Pareto Populations. Contract F33615-71-C-1463. Technical Report No. ARL 75-0201 for Aerospace Research Laboratories, Wright-Patterson AFB OH, June 1975 (AD-A014 740).
30. Kaminsky, Kenneth S. and Paul I. Nelson. "Best Linear Unbiased Prediction of Order Statistics in Location and Scale Families," Journal of the American Statistical Association, 70: 145-150 (1975).
31. Kapur, K. C. and L. R. Lamberson. Reliability in Engineering Design. New York: John Wiley and Sons, 1977.
32. Koutrouvelis, Ioannis. Estimation of Asymptotic Pareto Laws and the Tail of a Distribution. Contract Number N00014-72-C-0508. Technical Report No. 34 for Office of Naval Research, Arlington VA, Aug 1975 (AD-A018 173).
33. Kulldorff, Gunnar and Kerstin Vannman. "Estimation of the Location and Scale Parameters of a Pareto Distribution by Linear Functions of Order Statistics", Journal of the American Statistical Association, 68: 218-227 (1973).
34. Lilliefors, H. "On the Kolmogorov-Smirnov Test for Normality with Mean and Variance Unknown", Journal of the American Statistical Association, 62: 399-402 (1967).
35. Lilliefors, H. "On the Kolmogorov-Smirnov Test for the Exponential Distribution with Mean Unknown", Journal of the American Statistical Association, 64: 387-399 (1969).
36. Littel, Ramon C., James McClave, and Walter Offen. "Goodness-of-Fit Tests for the Two Parameter Weibull Distribution", Communications in Statistics, B8(3): 257-269 (1979).
37. Little, Robert E. Probability and Statistics for Engineers. Champaign IL: Matrix Publishers, Inc., 1978.
38. Mann, N. R., E. M. Scheuer, and K. W. Fertig. "A New Goodness-of-Fit Test for the Two-Parameter Weibull or Extreme-Value Distribution with Unknown Parameters", Communications in Statistics, 2: 383-400 (1973).

39. Massey, Frank J. "The Kolmogorov-Smirnov Test for Goodness of Fit", Journal of the American Statistical Association, 46: 68-78 (1951).
40. Mood, A. M. and F. A. Graybill, Introduction to the Theory of Statistics (Second Edition). New York: McGraw Hill Inc., 1963.
41. Moore, Albert H. and H. L. Harter. "One-order-statistic Conditional Estimators of Shape Parameters of Limited and Pareto Distributions and Scale Parameters of Type II Asymptotic Distributions of Smallest and Largest Values," IEEE Transactions on Reliability, R-16: 100-103 (1967).
42. Pigou, A. C., The Economics of Welfare. London: Macmillan and Co., 1948.
43. Ream, Thomas J. A New Goodness of Fit Test for Normality with Mean and Variance Unknown. MS Thesis, GOR/MA/81D-9. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, Dec 1981.
44. Steindl, Josef. Random Processes and the Growth of Firms. New York: Hafner Publishing Co., 1965.
45. Stephens, M. A. "EDF Statistics for Goodness of Fit and Some Comparisons", Journal of the American Statistical Association, 69: 730-737 (Sep 1974).
46. Stephens, M. A. "Asymptotic Results for Goodness-of-Fit Statistics with Unknown Parameters", Annals of Statistics, 4: 357-369 (1976).
47. Stephens, M. A. The Anderson-Darling Statistic. Grant No. DAAG29-77-G-0031. Technical Report No. 39 for the U.S. Army Research Office, Dept. of Statistics, Stanford University, Stanford CA, Oct 1979 (AD-A079 807).
48. Vannman, Kerstin. "Estimators Based on Order Statistics from a Pareto Distribution", Journal of the American Statistical Association, 71: 704-708 (Sep 1976).
49. Viviano, Philip J. A Modified Kolmogorov-Smirnov, Anderson-Darling, and Cramer-von Mises Test for the Gamma Distribution with Unknown Location and Scale Parameters. MS Thesis, GOR/MA/82D-4. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, Dec 1982.

50. Wingo, Dallas R. "Estimation in a Pareto Distribution: Theory and Computation". IEEE Transactions on Reliability, R-28: 35-37 (Apr 1979).

51. Wong, Wing-Yue. On the Property of Dullness of Pareto Distribution. Contract No. N00014-75-C0455. Technical Report No. 82-16 for the Office of Naval Research, Purdue University. West Lafayette IN, May 1982 (AD-A119 631).

52. Woodbury, Larry B. A New Goodness of Fit Test for the Uniform Distribution with Unspecified Parameters. MS Thesis, GOR/MA/82D-6. School of Engineering, Air Force Institute of Technology (AU), Wright-Patterson AFB OH, Dec 1982.

53. Woodruff, Brian W. and others. "A Modified Kolmogorov-Smirnov Test for Weibull Distributions with Unknown Location and Scale Parameters," IEEE Transactions on Reliability, R-32: 209-213 (Jun 1983).

54. Yoder, John D. Modified Kolmogorov-Smirnov, Anderson-Darling, and Cramer-Von Mises Tests for the Logistic Distribution with Unknown Location and Scale Parameters. MS Thesis, GOR/ENC/83D. School of Engineering, Air Force Institute of Technology (AU), Wright Patterson AFB OH, December 1983.

## VITA

Captain James E. Porter III was born in Tokyo, Japan, on 24 September 1951. He graduated from Judson High School, Converse, Texas, in 1969. He then attended the University of Texas at Austin and in 1974 graduated Phi Beta Kappa with a Bachelor of Science degree in Mathematics.

Upon completing Officer Training School and receiving his USAF commission in April 1975, he was assigned to the Space Systems (now called Space Operations) career field, Air Force Specialty Code (AFSC) 20XX. He served as a Space Surveillance Officer at the Sea-Launched Ballistic Missile Detection and Warning radar site, Fort Fisher AFS, North Carolina, from June 1975 to May 1977; and at the Ballistic Missile Early Warning System radar site, Thule, Greenland, from May 1977 to May 1978.

From June 1978 to May 1981 Captain Porter was assigned to Headquarters North American Aerospace Defense Command, Peterson AFB, Colorado, as a Space Systems Staff Officer. He next served as Space Operations Career Management Staff Officer, Air Force Manpower and Personnel Center, Randolph AFB, Texas, until May 1984. He then entered the Graduate Space Operations Program, School of Engineering, Air Force Institute of Technology.

Address: 4026 Kirby Drive, San Antonio, Texas 78219.

*AD-A163 831*

**REPORT DOCUMENTATION PAGE**

REPORT SECURITY CLASSIFICATION <b>UNCLASSIFIED</b>			1b. RESTRICTIVE MARKINGS			
2a. SECURITY CLASSIFICATION AUTHORITY			3. DISTRIBUTION/AVAILABILITY OF REPORT  <b>Approved for public release; distribution unlimited</b>			
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE						
4. PERFORMING ORGANIZATION REPORT NUMBER(S)  <b>AFIT/GSO/MA/85D-6</b>			5. MONITORING ORGANIZATION REPORT NUMBER(S)			
6a. NAME OF PERFORMING ORGANIZATION  <b>School of Engineering</b>		6b. OFFICE SYMBOL (If applicable) <b>AFIT/ENS</b>	7a. NAME OF MONITORING ORGANIZATION			
6c. ADDRESS (City, State and ZIP Code)  <b>Air Force Institute of Technology Wright-Patterson AFB OH 45433-6583</b>			7b. ADDRESS (City, State and ZIP Code)			
8a. NAME OF FUNDING/SPONSORING ORGANIZATION		8b. OFFICE SYMBOL (If applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER			
8c. ADDRESS (City, State and ZIP Code)			10. SOURCE OF FUNDING NOS.			
11. TITLE (Include Security Classification) <b>See Box 19</b>			PROGRAM ELEMENT NO.	PROJECT NO.	TASK NO.	WORK UNIT NO.
			PERSONAL AUTHOR(S) <b>James E. Porter III, Captain, USAF</b>			
13a. TYPE OF REPORT <b>MS Thesis</b>		13b. TIME COVERED FROM _____ TO _____	14. DATE OF REPORT (Yr., Mo., Day) <b>1985 December</b>		15. PAGE COUNT <b>183</b>	
16. SUPPLEMENTARY NOTATION						
17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) <b>Monte Carlo Method; Statistical Functions; Probability Distribution Function; Statistical Analysis; Statistical Decision Theory; Statistical Tests; Order Statistics</b>			
FIELD	GROUP	SUB. GR.				
<b>12</b>	<b>01</b>					
19. ABSTRACT (Continue on reverse if necessary and identify by block number)						
TITLE: <b>MODIFIED KOLMOGOROV-SMIRNOV, ANDERSON-DARLING, AND CRAMER-VON MISES TESTS FOR THE PARETO DISTRIBUTION WITH UNKNOWN LOCATION AND SCALE PARAMETERS</b>						
Approved for public release: IAW AFR 180-17. <i>Ly. Wolaver</i> <b>LYON E. WOLAVER</b> 16 JAN 86 Dean for Research and Professional Development Air Force Institute of Technology (AFIT) Wright-Patterson AFB OH 45433						
THESIS ADVISOR: <b>Dr Albert H. Moore</b> <b>Professor of Mathematics</b>						
DISTRIBUTION/AVAILABILITY OF ABSTRACT <b>UNCLASSIFIED/UNLIMITED</b> <input checked="" type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS <input type="checkbox"/>			21. ABSTRACT SECURITY CLASSIFICATION <b>UNCLASSIFIED</b>			
22a. NAME OF RESPONSIBLE INDIVIDUAL <b>Prof. Albert H. Moore</b>			22b. TELEPHONE NUMBER (Include Area Code) <b>(513)255-3098</b>		22c. OFFICE SYMBOL <b>AFIT/ENC</b>	

**19. ABSTRACT**

Modified Kolmogorov-Smirnov (K-S), Anderson-Darling (A-D), and Cramer-von Mises (C-VM) critical values are generated for the three-parameter Pareto distribution. The values may be used to test whether a set of observations follows a Pareto distribution when the location and scale parameters are unspecified and thus must be estimated from the sample. A Monte Carlo simulation of 5000 repetitions is used to generate critical values for sample sizes 5(S)30 (i.e., 5 to 30 in increments of 5) and Pareto shape parameters .5(.5)4.0.

A 5000-repetition Monte Carlo investigation is carried out by using 5, 15, and 25 observations from eight alternate distributions to compare the powers of the K-S, A-D, C-VM, and Chi-square tests. The power values of the tests are relatively low for a sample size of five. However, the powers of the modified K-S, A-D, and C-VM tests are considerably better than the Chi-square test at larger sample sizes. Next to the Chi-square test, the A-D test has the lowest power in most cases.

A functional relationship is identified between the modified K-S and C-VM test statistics and the Pareto shape parameter. The critical values are found to be a linear function of the shape parameters between 1.5 and 4.0.

**END**

**FILMED**

3-86

**DTIC**