AD-A112 469 UNCLASSIFIED	ILLINOIS UNIV A MULTI-SAMPLE CL JAN 82 H BOZDO TR-82-1	T CHICAGO CIRCLE USTER ANALYSIS U BAN, S L SCLOVE	DEPT OF QUANT	ITATIVEETC INFORMATION CR N00014-80-C	F/6 12/1 NITERETC(U) 0408 NL-	
						i.





### MULTI-SAMPLE CLUSTER ANALYSIS USING AKAIKE'S INFORMATION CRITERION\*

by

HAMPARSUM BOZDOGAN Quantitative Methods Department University of Illinois at Chicago Circle

and

STANLEY L. SCLOVE Departments of Mathematics and Quantitative Methods University of Illinois at Chicago Circle

> TECHNICAL REPORT NO. 82-1 January 30, 1982

PREPARED FOR THE OFFICE OF NAVAL RESEARCH UNDER CONTRACT NOOO14-80-C-0408, TASK NR042-443 with the University of Illinois at Chicago Circle

Principal Investigator: Stanley L. Sclove

Reproduction in whole or in part is permitted for any purpose of the United States Government.

Approved for public release; distribution unlimited



QUANTITATIVE METHODS DEPARTMENT UNIVERSITY OF ILLINOIS AT CHICAGO CIRCLE CHICAGO, ILLINOIS 60680

VV45

<sup>\*</sup>Presented by the first author as an Invited Paper, Special Session on Cluster Analysis, 789th Meeting, American Mathematical Society, University of Massachusetts, Amherst, MA, October 16-18, 1981.

Hamparsum Bozdogan and Stanley L. Sclove University of Illinois at Chicago Circle

#### CONTENTS

Abstract; Key Words and Phrases

- 1. Introduction
- 2. The Multi-Sample Cluster Problem
- 3. The Number of Clustering Alternatives for a Given K Samples into k Nonempty Clusters
- 4. AIC for the Univariate Model
- 5. AIC for the Multivariate Model
- 6. Numerical Examples of Multi-Sample Cluster Analysis on Real Data Sets
  - 6.1. Univariate Examples
  - 6.2. A Multivariate Example

Manufactor and the stable of most store and

Acknowledgement

References

# MULTI-SAMPLE CLUSTER ANALYSIS USING AKAIKE'S INFORMATION CRITERION\*

Hamparsum Bozdogan and Stanley L. Sclove University of Illinois at Chicago Circle

#### ABSTRACT

Multi-sample cluster analysis, the problem of grouping samples, is studied from an <u>information-theoretic</u> viewpoint via Akaike's Information Criterion (AIC). This criterion combines the maximum value of the likelihood with the number of parameters used in achieving that value. The multi-sample cluster problem is defined, and AIC is developed for this problem.

The form of AIC is derived in both univariate and multivariate analysis of variance models. Numerical examples are presented and results are shown to demonstrate the utility of AIC in identifying the best clustering alternatives.

Key Words and Phrases: Multi-sample cluster analysis; Akaike's Information Criterion (AIC); ANOVA Model, MANOVA Model; maximum likelihood.

Accession 15

-11-

<sup>\*</sup>Presented by the first author as an Invited Paper, Special Session on Cluster Analysis, 789th Meeting, American Mathematical Society, University of Massachusetts, Amherst, MA, October 16-18, 1981.

#### MULTI-SAMPLE CLUSTER ANALYSIS USING AKAIKE'S INFORMATION CRITERION\*

Hamparsum Bozdogan and Stanley L. Sclove University of Illinois at Chicago Circle

1. Introduction

In this paper, we shall develop Akaike's Information Criterion (AIC) for multi-sample cluster analysis. The problem of multi-sample cluster analysis arises when we are given a collection of samples (groups, treatments), to be clustered into homogeneous groups.

It is reasonable to provide a practically useful statistical procedure that would use some sort of statistical model to aid in comparisons of various collections of samples, identify homogeneous groups of samples, and tell us which should be clustered together and which samples should not.

Examples of multi-sample clustering situations are abundant. Here we mention a few.

Example 1.1. Botany: grouping of three types of species of iris, namely <u>Iris</u> setosa (S), <u>Iris versicolor</u> (Ve), and <u>Iris virginica</u> (Vi), given in Example 6.2 and Table 6.3 in Section 6, on the basis of each and of all the four variables. <u>Example 1.2</u>. <u>Zoology</u>: grouping of geographical locations to study the differences of populations of two types of species of <u>Crocidura</u>. Delany and Healy [7] studied variation in white-toothed shrews, that is, nocturnal mammals, in the British Isles. White-toothed shrews of genus <u>Crocidura</u> occur in the Channel and Scilly Islands of the British Isles and the French mainland. From p = 10 measurements on each of n = 399 skulls obtained from the K = 10 locations, Tresco, Bryher, St. Agnes, St. Martin's, St. Mary's, Sark, Jersey, Alderney, Guernsey, and Cap Gris Nez. The sample sizes for the data from the

<sup>\*</sup>Presented by the first author as an Invited Paper, Special Session on Cluster Analysis, 789th Meeting, American Mathematical Society, University of Massachusetts, Amherst, MA, October 16-18, 1981.

ten locations are, respectively,  $n_1 = 144$ ,  $n_2 = 16$ ,  $n_3 = 12$ ,  $n_4 = 7$ ,  $n_5 = 90$ ,  $n_6 = 25$ ,  $n_7 = 6$ ,  $n_8 = 26$ ,  $n_9 = 53$ ,  $n_{10} = 20$ . Attempts were made to analyze the pattern of variation between these ten populations to examine the belief that there may be two species of <u>Crocidura</u>, namely <u>Crocidura russula</u>, <u>Crocidura</u> <u>suaveolens</u>. The locations were geographically close, but it is assumed that only one sub-species was present in any one place. Thus the problem here is to cluster the locations, that is, "samples" into homogeneous groups to discover the origin of the two species.

Example 1.3. Air and Water Pollution: grouping of weather class types or nitrate sites to distinguish whether the source of nitrate is weather type or local. Heidorn [12] studied synoptic, that is, general weather patterns associated with nitrates in southern Ontario. In recent years, there has been growing concern over the potential hazard of particulate nitrate in the atmosphere which acts as a respiratory irritant, especially to those who have asthma problems. Nitrate is also suspected to lower the pH level in freshwater lakes.

A sample of n = 17 cities across southern Ontario from Windsor in the west to Kingston in the east was chosen as the location of nitrate sites. Nitrate concentrations for the 17 sites were measured. In order to determine the effect of weather patterns on the measurement of nitrate, eight weather class types were defined for the nitrate sites. Thus the problem here is to cluster the weather class types or the sites into homogeneous groups to determine whether the source of particulate nitrate is due to weather class type or is local.

<u>Example 1.4</u>. <u>Business and Economics</u>: grouping of corporations by their financial characteristics. Chen <u>et al.</u> [6], Williams and Goodman [16], and others, studied the statistical methods for clustering corporations on the

-2-

basis of yearly data concerning several of their financial characteristics. Thus the general problem here is to cluster the sets of corporations in order to detect, describe and distinguish relatively homogeneous groups of companies so that the formation of the groups and organizational behavior of companies can be studied and compared.

So, as we see, multi-sample cluster analysis examples are quite rich and varied.

The analysis of variance (ANOVA) is a widely used model for comparing two or more univariate samples, where the familiar Student's t and F statistics are used for formal comparisons among two or more samples. Multivariate analysis of variance (MANOVA) is a widely used model for comparing two or more multivariate samples. In the MANOVA model, the likelihood ratio principle leads to <u>Wiiks' [17] lambda</u>, or in short <u>Wilks' A criterion</u> as the test statistic. It plays the same role in multivariate analysis that the F-ratio statistic plays in the univariate case.

Often, however, the formal analyses involved in MANOVA are not revealing or informative. Therefore, in this paper we shall propose Akaike's Information Criterion (AIC) as a new procedure for comparing the clusters, and use it to identify the best clustering alternatives.

In 1971, Akaike first introduced an information criterion, referred to as an automatic (model) identification criterion or Akaike's information criterion (AIC), for the identification and comparison of statistical models in a class of competing models with different numbers of parameters. It is defined by

(1.1) AIC = -2 log<sub>e</sub> (maximized likelihood)

+2 (number of independently adjusted parameters within the model).

-3-

It was obtained with the aid of an information theoretic interpretation of the method of maximum likelihood by Akaike ([2], [3]). It estimates minus twice the expected log likelihood of the model whose parameters are determined by the method of maximum likelihood. When several competing models are being compared or fitted, AIC is a simple procedure which measures the <u>badness of fit</u> or the <u>discrepancy</u> of the estimated model from the true model when a set of data is given.

The first term in (1.1) stands for the penalty of <u>badness of fit</u> or <u>downward bias</u> when the maximum likelihood estimators of the parameters of the model are used. The second term in the definition of AIC, on the other hand, stands for the penalty of increased <u>unreliability</u> or <u>compensation for the bias</u> in the first term as a consequence of increasing number of parameters. If more parameters are used to describe the data, it is natural to get a larger likelihood, possibly without improving the true goodness of fit by penalizing the use of additional parameters.

Thus, when there are several competing models, the parameters within the models are estimated by the method of maximum likelihood and the AIC-values are computed and compared to find a model with the minimum value of AIC. This procedure is called the <u>minimum AIC procedure</u>. The model with the minimum AIC is called the <u>minimum AIC estimate</u> (MAICE) and is designated as the <u>best model</u>.

In Section 2, we shall define the general multi-sample cluster problem, and in Section 3, we shall briefly discuss the number of clustering alternatives for a given K groups or samples into k nonempty clusters. In the subsequent sections, that is, in Section 4 and in 5, we shall derive the AIC procedure for the univariate analysis of variance (ANOVA) model, and the multivariate analysis of variance (MANOVA) model. In Section 6, we shall give

-4-

numerical examples for both univariate and multivariate multi-sample cluster analysis on real data sets to demonstrate our results of AIC and minimum AIC procedures obtained from different computer analyses.

# 2. The Multi-Sample Cluster Problem

Suppose each individual, object, or case, has been measured on p response or outcome measures (dependent variables) simultaneously in K independent groups or samples (factor levels). Let

(2.1) 
$$\underline{X} (n \times p) = \begin{bmatrix} \underline{X}_1 \\ \underline{X}_2 \\ \vdots \\ \vdots \\ \underline{X}_K \end{bmatrix}$$

be a single data matrix of K groups or samples, where  $\underline{X}_{g}(n_{g}xp)$  represents the observations from the g-th group or sample, g=1,2,...,K, and n  $\sum_{g=1}^{K} n_{g}$ . The goal of cluster analysis is to put the K groups or samples into k homogeneous groups, samples, or classes where k is unknown, but k<K.

Often individuals or objects have been sampled from K>1 populations. For multi-samples or multiple groups of individuals or objects the data matrix may be represented in partitioned form as above. Let  $n_g$  represent the number of individuals in the g-th (random) sample, g=1,2,...,K. The  $n_g$  are not restricted to being equal or proportional to other  $n_g$ 's. The total number of observations is  $n = \sum_{g=1}^{K} n_g$ . Let  $X_{gi}$  be the px1 vector of observations in group g=1,2,...,K, and for individual i=1,2,...,n\_g.

-5-

#### 3. The Number of Clustering Alternatives for a Given K Samples into k Nonempty Clusters

In this section, we shall briefly discuss how to obtain the total number of clustering alternatives for a given K, the number of <u>groups</u> or <u>samples</u>. For this, we shall recall some established results.

<u>Theorem 3.1</u>. The number of ways of clustering K groups or samples into k clusters such that none of the k clusters is empty is given by

(3.1) 
$$\sum_{q=0}^{k} {\binom{k}{g}} {(-1)}^{g} {(k-g)}^{K}$$

where the order of groups or samples within each cluster is irrelevant.

# Proof. Duran and Odell [9].

Artice And Andrewson, Station, and

In this theorem the k clusters are assumed to be distinct. However, in clustering K groups or samples into k clusters, none of which is empty, the order of the k clusters is irrelevant. Consequently, from this fact and Theorem 3.1, it follows that the <u>total number of ways of clustering K groups or</u> samples into k clusters is given by

(3.2) 
$$S(K,k) = \frac{1}{k!} \sum_{g=0}^{k} {\binom{k}{g}} (-1)^{g} (k-g)^{K}$$

which is known as the Stirling Number of the Second Kind (see, e.g., Abramowitz and Stegun [1]) and also called the <u>number of clustering alternatives</u>.

If k, the number of clusters of groups or samples is known in advance, then the total number of clustering alternatives is given by S(K,k). However, if k is not specified a priori and varies, then the total number of clustering

-6-

alternatives for a given K, the number of groups or samples, is given by

(3.3) 
$$\sum_{k=1}^{K} S(K,k)$$
.

Table 3.1 gives S(K,k) for values of K and k up to 10.

k	1	2	3	4	5	6	7	8	9	10	Total
К											
1 2 3 4 5 6 7 8 9 10	1 1 1 1 1 1 1 1 1	1 3 7 15 31 63 123 255 511	1 6 25 90 301 966 3021 9318	1 10 65 350 1701 7770 34101	1 15 140 1050 6951 42525	1 21 266 2645 22821	1 28 462 5879	1 36 750	1 45	1	1 2 5 203 877 4136 21142 115952

TABLE 3.1. NUMBER OF CLUSTERING ALTERNATIVES FOR VARIOUS VALUES OF K AND K

Consider, for example, K=3 samples. We now wish to cluster K=3 groups or samples first into k=3 groups or samples, then into k=2 groups or samples, and k=1 group or sample in a hierarchical fashion.

From Table 3.1, we have the total number of ways of clustering K=3 groups or samples into k=3 homogeneous groups or samples is 1. The total number of ways of clustering K=3 groups or samples into k=2 homogeneous groups or samples is 3. The total number of ways of clustering k=3 groups or samples into k=1 homogeneous group or sample is 1. Thus adding up these results, we obtain, in total 5 clustering alternatives as the total for K=3 groups or samples into k=1,2, and 3 homogeneous groups. We note that 5 is nothing but the sum of the values of row 3 in Table 3.1.

The 5 clustering alternatives can be classified according to their <u>representation forms</u> to make it easy to list all 5 possible clustering alternatives. The representation forms in this case are denoted by

- (1) (1) (1) (1),
- (**ii**) {2} {1},
- (111) {3},

where each of the components in a representation {g} denotes the number, g, of groups or samples in the corresponding cluster. The components of a representation form will always be written in a hierarchical order to depict the patterns of clustering alternatives. In our example there are 5 clustering alternatives but only 3 representation forms. In general the number of representation forms is much smaller then the number of clustering alternatives.

We now list the clustering alternatives corresponding to their representation forms in Table 3.2 as follows:

Alternatives	Clustering	Number of Parameters m
1 2 3 4 5	(1) (2) (3) (1 2) (3) (1 3) (1) (2 3) (1) (1 2 3)	3 2 2 2 1

TABLE 3.2. A SIMPLE PATTERN OF CLUSTERING ALTERNATIVES WHEN K=3 AND k=3, 2, and 1

For example, in alternative one, the group or sample 1, 2, and 3 are clustered as singletons. In terms of a hypothesis on means, this corresponds

-8-

to  $\mu_1$ ,  $\mu_2$ , and  $\mu_3$  all being different, and therefore, the number of parameters, m, is equal to 3. Hence, indicating that group 1, 2, and 3 are all heterogeneous. In alternative two, groups or samples 1 and 2 are clustered together forming a homogeneous subset, and group or sample 3 is clustered alone forming a heterogeneous subset. In terms of a hypothesis on means, this corresponds to  $\mu_1 = \mu_2$ , and  $\mu_3$  is different from both  $\mu_1$  and  $\mu_2$  with the total number of parameters m being equal to 2. In a similar fashion, we interpret the other clustering alternatives continuing down the line of Table 3.2.

As a last example, we shall just list the results of the total number of possible clustering alternatives when K=4 groups or samples in Table 3.3 as follows.

Alternatives	Clustering	Number of Parameters, m
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15	$ \begin{array}{c} (1) & (2) & (3) & (4) \\ (1 & 2) & (3) & (4) \\ (1 & 3) & (2) & (4) \\ (1 & 4) & (2) & (3) \\ (2 & 3) & (1) & (4) \\ (2 & 4) & (1) & (3) \\ (3 & 4) & (1) & (2) \\ (1 & 2) & (3 & 4) \\ (1 & 3) & (2 & 4) \\ (1 & 4) & (2 & 3) \\ (1 & 2 & 4) & (3) \\ (1 & 2 & 4) & (3) \\ (1 & 2 & 3 & 4) \\ (1 & 2 & 3 & 4) \\ (1 & 2 & 3 & 4) \\ (1 & 2 & 3 & 4) \\ \end{array} $	4 3 3 3 3 3 3 2 2 2 2 2 2 2 2 2 2 2 2 2

TABLE 3.3. A SIMPLE PATTERN OF CLUSTERING ALTERNATIVES WHEN K=4 AND k=4, 3, 2, AND 1

In concluding this section, we see that in general the total number of

-9-

ways of clustering K groups or samples into k homogeneous groups or samples is given by equation (3.2), and the total number of possible clustering alternatives is given by the expression (3.3).

4. AIC For The Univariate Model

We now turn our attention to consider situations with several univariate normal samples. The general layout for such data (one-way ANOVA) is represented in the following tabular form.

	Groups	
	1 2 K	
	$z_{11}$ $z_{21}$ $\cdots$ $z_{K_2}$	
	$z_{12} z_{22} \cdots z_{K_2}$	
Observations	• • • • • •	
	• • • • • •	
	z <sub>1n1</sub> z <sub>2n2</sub> • • • <sup>z</sup> Kn <sub>K</sub>	
TOTALS	$T_1 T_2 \cdots T_K$	т
SAMPLE SIZES	n <sub>1</sub> n <sub>2</sub> n <sub>K</sub>	n = ∑ng g=1
SAMPLE MEANS	$\overline{z}_1, \overline{z}_2, \ldots, \overline{z}_{K}$	= Z
VARIANCES	$s_1^2$ $s_2^2$ $s_K^2$	s <sup>2</sup>

TABLE 4.1. GENERAL DATA REPRESENTATION FOR ONE-WAY ANOVA

-10-

٦

For example, we may have multi-sample data with samples of sizes  $n_1, n_2, \ldots, n_K$  which are assumed to have come from K populations, the first with mean  $\mu_1$  and variance  $\sigma^2$ , the second with mean  $\mu_2$  and variance  $\sigma^2$ ,..., the Kth with mean  $\mu_K$  and variance  $\sigma^2$ . We may want to compare these K group or sample means  $\mu_1, \mu_2, \ldots, \mu_K$  given that all have a common  $\sigma^2$ . Hence, this is the well known analysis of variance (ANOVA) model. In terms of the parameters the ANOVA model is  $\theta = (\mu_1, \mu_2, \ldots, \mu_K, \sigma^2)$  with m=k+1 parameters, where k is the number of groups.

We shall derive the form of AIC for this model. Recall the definition of AIC from Section 1,

AIC = -2 
$$\log_e L(\theta) + 2m$$
  
= -2  $\log_e (maximized likelihood) + 2m$ ,

where m denotes the number of independently adjusted parameters within the model.

Suppose there are K independent samples of independent observations, with ng, g=1,2,...,K, observations in the g-th group and n =  $\sum_{\substack{g=1\\g=1}}^{K} n_g$ . Denote the unknown means of the groups by  $\mu_1, \mu_2, \dots, \mu_K$ . Assume that the samples  $(z_{11}, z_{12}, \dots, z_{1n_1}; \dots; z_{K^1}, \dots, z_{Kn_K})$  are drawn randomly from K populations which are N( $\mu_g, \sigma^2$ ). If the groups can differ only in their means, we may express this as

(4.1)  $z_{gi} = \mu g + \epsilon_{gi}, g=1,2,...,K; i=1,2,...,n_g,$ 

where  $z_{gi}$  is the value of the response or outcome variable in the g-th group for the i-th individual or object,  $\mu_g$  are parameters,  $\epsilon_{gi}$  are independent N(0, $\sigma^2$ ) error variables.

This equation is called the one-way ANOVA model.

Thus, the basic null hypothesis of interest in this case is given by

(4.2)  $H_0: \mu_1 = \mu_2 = \dots = \mu_{\kappa}$ 

The alternative hypothesis is given by

H, : the K population means are not all equal.

Every analysis of variance involves a partitioning of the total sum of squares of deviations, SST, into the within-group sum of squares of deviations, SSW, and the between-group sum of squares of deviations, SSB. For more details on this, we refer the reader to any basic text on statistics, e.g., Anderson and Sclove [4].

We now derive the form of Akaike's Information Criterion (AIC) for the one-way ANOVA model given in (4.1).

The likelihood function is given by

(4.3) 
$$L({u_g}, \sigma^2; z) = (2\pi\sigma^2)^{-n/2} \exp\left[-\sum_{g=1}^{K} \sum_{i=1}^{n_g} (z_{gi} - u_g)^2/(2\sigma^2)\right].$$

The log likelihood function is

(4.4) 
$$l(\{\mu_g\},\sigma^2;z) \equiv \log L(\{\mu_g\},\sigma^2;z)$$
  
=  $-\frac{n}{2}\log(2\pi) - \frac{n}{2}\log(\sigma^2) - \sum_{\substack{j=1 \ j=1}}^{K} \sum_{\substack{j=1 \ j=1}}^{ng} (z_{gj} - \mu_g)/(2\sigma^2).$ 

As is well known, the MLE's are

(4.5) 
$$\hat{\mu}_{g} = \frac{1}{n_{g}} \sum_{i=1}^{n_{g}} z_{gi} = \bar{z}_{g_{s}}, g=1,2,\ldots,K,$$

-12-

-13-

and

(4.6) 
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{g=1}^{K} \sum_{i=1}^{ng} (z_{gi} - \bar{z}_{g_i})^2 = \frac{SSW}{n}$$
,

where SSW =  $\sum_{g=1}^{K} (z_{g1} - \overline{z}_{g_*})^2$ , the Within Group Sum of Squares.

Substituting these back into (4.4), we have

$$l(\{\hat{\mu}_{g}\}, \hat{\sigma}^{2}; z) \equiv \log L(\{\hat{\mu}\}, \hat{\sigma}^{2}; z)$$
$$= -\frac{n}{2}[\log(2\pi) + \log \frac{SSW}{n}] - \frac{n}{2}.$$

Since

(4.7) AIC = 
$$-2 \log_e L(\frac{\theta}{2}) + 2m$$
,

where m is the number of parameters, and since

(4.8) -2 log L({
$$\hat{\mu}_g$$
}, $\hat{\sigma}^2$ ) = n log(2 $\pi$ ) + n log  $\frac{SSW}{n}$  + n ,

then AIC becomes

(4.9) AIC = 
$$n \log(2\pi) + n \log \frac{SSW}{n} + n + 2(k+1)$$
.

Since the constants do not affect the result of comparison of models, we could ignore them and use the simplified version

(4.10) AIC\* = nlog<sub>e</sub> SSW + 2(k+1)  
where 
$$n = \sum_{g=1}^{K} n_g$$
 = the total sample size,

SSW = Within Group Sum of Squares, and

k = number of groups or samples compared, or the number of independently adjusted parameters within the model.

However, for purposes of comparison we retain the constants and use AIC.

### 5. AIC For the Multivariate Model

In this section we shall study the natural extension of the univariate model we considered in Section 4 to its multivariate analogue. Therefore, throughout this section we shall suppose that we may have independent data matrices  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_K$ , where the rows of  $\underline{X}_g$  (ngxp) are independent and identically distributed (i.i.d.) Np( $\mu g, \underline{\Sigma}$ ), g=1,2,...,K. In terms of the parameters  $\theta = (\mu_1, \mu_2, \dots, \mu_K, \underline{\Sigma})$  the model we shall consider here is

 $\boldsymbol{\theta} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma})$ 

with m = kp + p(p+1)/2 parameters, where k is the number of groups, and p is the number of variables.

As in the univariate case, consider K normal populations with different mean vectors  $\mu_g$ , g=1,2,...,k,...,K. Let  $z_{gi}$ , g=1,2,...,K; i=1,2,...,n<sub>g</sub>, be a random sample of observations from the g-th population N<sub>p</sub>( $\mu_g, \underline{\Sigma}$ ). If the groups or samples can differ only in their mean vectors, we can write the multivariate one-way analysis variance (MANOVA) model as

(5.1)  $z_{gi} = \mu g + \varepsilon_{gi}$ , g=1,...,K;  $i=1,2,...,n_g$ ,

where  $z_{gi}$  is the (p x 1) response or outcome vector in the g-th group for i-th individual or object,  $\mu_g$  are vector parameters, and  $\epsilon_{gi}$  are independent  $N_p(0, \underline{\Sigma})$  random vector errors.

-14-

Thus, the basic <u>null hypothesis</u> we usually are interested in testing is given by

(5.2) 
$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

The alternative hypothesis is given by

$$H_1$$
 : Not all  $\mu$  are equal.

Wilks' lambda is a <u>general</u> statistic for handling this problem. Although, there are several other conventional statistics for this purpose, they all can be viewed as special cases of Wilks'  $\Lambda$  which we shall not discuss here.

For notational purposes, we shall denote  $\underline{T}$  to be the "total" sum of squares and products (SSP) matrix,  $\underline{W}$  to be the "within-group" or "withinsample" SSP matrix, and  $\underline{B}$  to be the "between-group" SSP matrix. Hence, it can be shown that

$$(5.3) \quad \underline{T} = \underline{W} + \underline{B} ,$$

where

ومردادة فالمستحكرتهم ومستوحة المتقام والكلاحص ومكالتكم وتكرا كالت

(5.4) 
$$\underline{T} = \sum_{g=1}^{K} \sum_{i=1}^{n_g} (z_{gi} - \overline{z})(z_{gi} - \overline{z})',$$

(5.5) 
$$\underline{W} = \sum_{g=1}^{K} \sum_{i=1}^{n_g} (z_{gi} - \overline{z}_g)(z_{gi} - \overline{z}_g)',$$

and

(5.6) 
$$\frac{B}{B} = \sum_{q=1}^{K} n_q (\overline{z}_q - \overline{z})(\overline{z}_q - \overline{z})',$$

-15-

with

$$\bar{z}_{g} = \frac{1}{n_{g}} \sum_{i=1}^{n_{g}} z_{gi} , g=1,2,...,K ,$$

$$\bar{z} = \frac{1}{n} \sum_{g=1}^{K} \sum_{i=1}^{n_{g}} z_{gi} , n = \sum_{g=1}^{K} n_{g} .$$

Therefore, we can present multivariate one-way analysis of variance (MANOVA) table as follows.

Source	d.f.	SSP matrix	Wilks' criterion
Between samples	K-1	<u>B</u>	<u>IM</u>
			III
Within samples	n-K	W	~A(p;n-K;K-1)
Total	n-1	Ţ	

TABLE 5.1. MANOVA TABLE

Now, we derive the form of Akaike's Information Criterion (AIC) for the MANOVA model given in (5.1), subject to the constraint given in (5.2). The likelihood function of all the sample observations is given by

(5.7) 
$$L(\underline{\mu}_{g},\underline{\Sigma}_{g};\underline{Z}) = \prod_{g=1}^{K} L_{g}(\underline{\mu}_{g},\underline{\Sigma}_{g};\underline{Z}_{g}),$$

or by

(5.8) 
$$L = (2\pi)^{-np/2} \frac{K}{\pi} \frac{|\underline{\Sigma}g|}{g=1} - \frac{ng/2}{\pi} x$$

exp 
$$\{-\frac{1}{2} \operatorname{tr} \sum_{g=1}^{K} \frac{A_g}{A_g} - \frac{1}{2} \operatorname{tr} \sum_{g=1}^{K} \operatorname{ng} \Sigma_g (\overline{Z}_g - \mu_g)(\overline{Z}_g - \mu_g)'\},$$

where 
$$n = \sum_{g=1}^{K} n_g$$
 and  $\underline{A}_g = \sum_{i=1}^{n_g} (\underline{z}_{gi} - \overline{\underline{z}}_g) (\underline{z}_{gi} - \overline{\underline{z}}_g)'$ .

The log likelihood function is

(5.9) 
$$l(\underline{\mu}g,\underline{\Sigma};\underline{Z}) \equiv \log_{e}L$$
  
$$= -\frac{np}{2}\log(2\pi) - 1/2\sum_{g=1}^{K} n_{g}\log|\underline{\Sigma}g| - 1/2tr\sum_{g=1}^{K} \underline{\Sigma}g^{-1}\underline{A}g$$
$$- 1/2tr\sum_{g=1}^{K} n_{g}\underline{\Sigma}g^{-1}(\underline{Z}g - \underline{\mu})(\underline{Z}g - \underline{\mu}g)'.$$

Since the common covariance matrix is  $\underline{\Sigma}$ , then the log likelihood function becomes

$$(5.10) \quad l\{\mu_{g}\}, \underline{\Sigma}; \underline{Z}\} \equiv \log_{e}L(\{\mu_{g}\}, \underline{\Sigma}; \underline{Z})$$

$$= -\frac{np}{2} \log(2\pi) - \frac{n}{2} \log |\underline{\Sigma}| - \frac{1}{2} \operatorname{tr} \underline{\Sigma}^{-1} \sum_{g=1}^{K} \underline{A}_{g}$$
$$- \frac{1}{2} \operatorname{tr} \underline{\Sigma}^{-1} \sum_{g=1}^{K} n_{g} (\overline{\underline{Z}}_{g} - \underline{\mu}_{g}) (\overline{\underline{Z}}_{g} - \underline{\mu}_{g})',$$

and the maximum-likelihood estimates (MLE's) of  $\mu_g,$  and  $\underline{\Sigma}$  are

(5.11) 
$$\hat{\mu}_{g} = \bar{z}_{g}, g=1,2,\ldots,K,$$

and

$$(5.12) \quad \hat{\underline{\Sigma}} = n^{-1} \underline{W}$$

where  $\underline{W} = \sum_{g=1}^{K} \underline{A}_{g}$ .

Substituting these back into (5.10) and simplifying, the maximized log likelihood becomes

(5.13) 
$$1(\{\hat{\mu}_{g}\},\hat{\Sigma};Z) \equiv \log L(\{\hat{\mu}_{g}\},\hat{\Sigma};Z)$$
  
=  $-\frac{np}{2}\log(2\pi) - \frac{n}{2}\log|n^{-1}\underline{W}| - \frac{np}{2}$ ,

where  $\underline{W}$  is the "within-group" SSP matrix.

Since

(5.14) AIC = 
$$-2 \log_{eL}(\theta) + 2m$$
,

where  $m = kp + \frac{p(p+1)}{2}$  is the number of parameters, then AIC becomes

(5.15) AIC = 
$$np \log(2\pi) + n\log |n^{-1}\underline{W}| + np + 2[kp + \frac{p(p+1)}{2}]$$

Since the constants do not affect the result of comparison of models, we could ignore them and reduce the form of AIC to a much simpler form

(5.15) AIC\* = 
$$n \log_{e} |\underline{W}| + 2[kp + \frac{p(p+1)}{2}]$$

-18-

where 
$$n = \sum_{g=1}^{n} n_g = the total sample size,$$

k = number of groups or samples compared,

p = number of variables.

However, for purposes of comparison we retain the constants and use AIC.

### 6. Numerical Examples of Multi-Sample Cluster Analysis on Real Data Sets

In this section we shall give numerical examples of both univariate and multivariate multi-sample data, and cluster the groups or samples, and choose the best clusterings by using Akaike's Information Criterion (AIC) as derived in Sections 4 and 5.

Our computations were carried out for all the examples we shall present here on an IBM 370, using various statistical software packages such as MINITAB, SPSS, and SPEAKEASY (VM/CMS version).

6.1. Univariate Examples

For the univariate numerical examples we shall illustrate our results on two data sets, a biomedical data set of Dolkart, Halpern, and Perlman [8] and Fisher [10] iris data. Here we shall take 150 iris specimens on each of the four morphological variables: sepal length and width and petal length and width and demonstrate our results on these variables individually rather than considering all of them together.

<u>Example 6.1</u>. (Brown and Hollander [5]) <u>Antibody Responses in Three Groups of</u> <u>Mice</u>: "Dolkart, Halpern, and Perlman [8] compared antibody responses in normal and alloxan diabetic mice. Their investigation was designed to study the circulating antibody response in alloxan diabetic, insulin-treated

-19-

diabetic and normal CF-1 mice injected with serum albumin.

"Only those animals treated with alloxan who had elevated serum glucose levels (250mg/100 ml or higher) were included in the study, together with a group of normal animals. Animals were bled from the orbital sinus, and the serum analyzed for antigen binding capacity of BSA, glucose concentration, and serum proteins. BSA was iodinated with I-131, and the antigen-binding capacity of each serum sample was determined as micrograms of BSA nitrogen bound by 1 ml of undiluted serum." The data are given in Table 6.1.

Norma 1	Alloxan Diabetic	Alloxan Diabetic- Treated with Insulin
155.76	390,72	82.50
282.00	46-20	99,66
197.34	468,60	97.66
297.00	86.46	150.48
115.50	174-02	242.88
126.72	132.66	67.98
110 16	13 20	227 70
20 04	A08 06	130 68
252 70	167 64	73 26
122 10	62 04	17 92
240 14	127 28	10 90
100 00	167 • 30 275 00	100 22
142 22	2/3.00	
143.22	1/0.22	/1.54
04.02	145.80	133.32
25.54	108.24	464.64
85.80	275.88	36.96
122.10	50.16	46.20
454.85	72.60	34.32
655.38		43.56
13.86		

TABLE 6.1MICROGRAMS OF BSA NITROGEN BOUND PER m1 OF UNDILUTED<br/>MOUSE SERUM ON DAY 39, FOLLOWING INJECTION OF 5 mg<br/>BSA ANTIGEN INTO EACH ANIMAL ON DAY 0 AND 28

Source: R.E. Dolkart, B. Halpern, and J. Perlman [8].

In this example we are given K=3 groups or samples and we wish to cluster them into k=1, 2, and 3 homogeneous groups. From Table 3.1, as we know, there

are 5 total possible clustering alternatives, namely, (1) (2) (3) all separate, and (1 2) (3), (1 3) (2), (2 3) (1), and (1 2 3) all together. Let us code <u>Normal Group=1, Alloxan Diabetic Group=2</u>, and <u>Alloxan Diabetic-Treated with</u> <u>Insulin Group=3</u>. Considering the ANOVA model as our underlying model for comparisons of these groups, from a simple ANOVA run on the computer we computed the AIC's for each of the 5 clustering alternatives. The results are shown in Table 6.2.

Alternative	Clustering	nlog <sub>e</sub> (2 <del>m</del> )	nloge <sup>SSW</sup> /n	n	k	2(k+1)	AIC
1 2 3 4 5	(1) (2) (3) (1 2) (3) (1 3) (2) (2 3) (1) (1 2 3)	104.758 104.758 104.758 104.758 104.758	559.139 559.149 561.945 561.513 562.581	57 57 57 57 57 57	3 2 2 2 1	8 6 6 6 4	728.897 <sup>C</sup> 726.907 <sup>a</sup> 729.703 729.271 728.339 <sup>b</sup>

TABLE 6.2 THE AIC'S FOR ANTIBODY RESPONSES IN THREE GROUPS OF MICE

n = 20 + 18 + 19 = 57

AIC =  $nlog_e(2\pi)$  +  $nlog_e \frac{SSW}{n}$  + n + 2 (k+1)

<sup>a</sup>First Minimum AIC

<sup>b</sup>Second Minimum AIC

<sup>C</sup>Third Minimum AIC

In this example the first minimum AIC occurs at the alternative submodel 2. That is, the MAICE is submodel 2 indicating to us that in terms of clustering, Normal Group=1 and Alloxan Diabetic Group=2 should be clustered together, and Alloxan Diabetic-Treated with Insulin Group=3 should be clustered by itself. Therefore, in terms of a hypothesis on means, (1 2) (3) corresponds to  $\mu_1 = \mu_2 \neq \mu_3$  indicating that Normal and Alloxan Diabetic Groups form the best homogeneous set in terms of their nitrogen-binding capacities, and the Alloxan Diabetic-Treated with Insulin Group forms a set by itself. On the other hand, the second minimum AIC occurs at the alternative submodel 5, and the third minimum AIC is at the alternative submodel 1 indicating that either we should cluster all the groups together or treat each group separately, but if we were to compare each group separately to the Normal Group=1, then we should choose Normal Group=1 with Alloxan Diabetic Group=2 together as the best choice by the minimum AIC procedure.

<u>Example 6.2</u>. <u>Clustering of Irises by Groups</u>: As we mentioned in Example 1.2, the iris data set is composed of 150 iris species belonging to three groups or species, namely <u>Iris setosa</u> (S), <u>Iris versicolor</u> (Ve), and <u>Iris virginica</u> (Vi) measured on sepal and petal length and width. Each group is represented by 50 plants. The data set for the 150 irises are given in Table 6.3.

This data set has been quite extensively studied in classification and cluster analysis since it was published by Fisher [10], and still today, is being used as a "testing ground" for classification and clustering methods proposed by many investigators such as Friedman and Rubin [11], Kendall [13], Solomon [15], Mezzich and Solomon [14], and many others, including the present authors.

For each of the 150 plants we already know the group structure of the iris species, namely K=3 groups or samples. Even though the two species, <u>Iris</u> <u>setosa</u> and <u>Iris versicolor</u> were found growing in the same colony, and <u>Iris</u> <u>virginica</u> was found growing in a different colony, Fisher reports in his linear discriminant analysis the separation of <u>I. setosa</u> completely from <u>I.</u> <u>versicolor</u> and <u>I. virginica</u>. Since then other investigators have shown similar results in their studies such as the ones we mentioned above.

-22-

	Iris s	etosa		In	is ver	sicolor	color Iris virginica				
Sepal length	Sepal width	Pstal Tength	Peta] width	Sepal length	Sepal width	Petal length	Petal width	Sepal length	Sepal width	Petal length	Petal width
5444554604948838741714161800227842590594105401 5444554545445555555555555555545545554455544555445554455544555445554455544555445554554555455455545554555455545554555455545555	502169449174000495884763404542141212560453238 333333332333344333333333333333333344333333	$1.4 \\ 1.5 \\ 4.7 \\ 4.5 \\ 5.6 \\ 4.1 \\ 1.5 $			333222332223322233222332223322222222222	4.4.9065573695207645159809734805558791555741046032	1.4 1.5 1.5 1.5 1.6 0.3 4 0.5 0 4 3 4 5 0 5 1.8 3 5 2 3 4 4 7 5 0 102 6 5 6 5 3 3 3 2 4 2 0 3 1.2 1.2 1.2 1.2 1.2 1.2 1.2 1.2 1.2 1.2	6.813569372548784577096737221424943173409798887 6.5766747655665566777665777666576667777666576665566	370900595627058208622887328080888860410111723 3232332223322332233322332232233222333333	019686538113501357907979708968146161658451197	2121221122104382353008188816902543488143935
4.5 5.1 4.6 5.3 5.0	3.8 3.2 3.7 3.3	1.6 1.4 1.5 1.4	0.2 0.2 0.2 0.2	5.7 5.2 5.1 5.7	2.9 2.9 2.5 2.8	4.2 4.3 3.0 4.1	1.2 1.3 1.1 1.3	6.3 6.5 6.2 5.9	2.5 3.0 3.4 3.0	5.0 5.2 5.4 5.1	1.9 2.0 2.3 1.8

TABLE 6.3 . MEASUREMENTS ON THREE TYPES OF IRIS

With this in mind, let us take K=3 groups or species on each of the variables separately and cluster them into k=1, 2, and 3 homogeneous groups. Since we are dealing with K=3 groups, by now we know that there are 5 total possible clustering alternatives. Denoting <u>I. setosa</u> by S, <u>I. versicolor</u> by Ve, and <u>I. virginica</u> by Vi, we have (S) (Ve) (Vi), (S, Ve) (Vi), (S, Vi) (Ve), (Ve, Vi) (S), and (S, Ve, Vi) as all the possible clustering alternatives of three iris species. Using the ANOVA model as our underlying model for comparisons of these iris groups, from a simple ANOVA run on the computer by using SPSS MANOVA program which performs both univariate and multivariate linear estimation and tests of hypotheses, we obtained the AIC's for each of the 5 clustering alternatives of iris groups on each of the four variables separately. We report our results on each of the four variables, respectively, as follows.

Alternative	Clustering	nlog <sub>e</sub> (2π)	nloge <sup>SSW</sup> /n	n	k	2(k+1)	AIC
1 2 3 4 5	(S) (Ve) (V1) (S, Ve) (V1) (S, V1) (Ve) (Ve, V1) (S) (S, Ve, V1)	275.681 275.681 275.681 275.681 275.681 275.681	-200.295 -135.669 - 58.550 -163.740 - 56.966	150 150 150 150 150	3 2 2 2 1	8 6 6 4	233.386a 296.012 373.131 267.941b 372.715

TABLE 6.4. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE SEPAL LENGTH

TABLE 6.5. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE SEPAL WIDTH

Alternative	Clustering	nlog <sub>e</sub> (2m)	nloge <sup>SSW</sup> /n	n	k	2(k+1)	AIC
1 2 3 4 5	(S) (Ve) (V1) (S, Ve) (V1) (S, V1) (Ve) (Ve, V1) (S) (S, Ve, V1)	275.681 275.681 275.681 275.681 275.681 275.681	-326.949 -252.915 -287.157 -318.019 -250.129	150 150 150 150 150	3 2 2 2 1	8 6 6 4	106.732 <sup>a</sup> 178.766 144.524 113.662 <sup>b</sup> 179.552

-24-

Alternative	Clustering	nlog <sub>e</sub> (2π)	nloge <sup>SSW</sup> /n	n	k	2(k+1)	AIC
1 2 3 4 5	(S) (Ve) (Vi) (S, Ve) (Vi) (S, Vi) (Ve) (Ve, Vi) (S) (S, Ve, Vi)	275.681 275.681 275.681 275.681 275.681 275.681	-255.988 59.442 163.259 -116.579 169.493	150 150 150 150 150	3 2 2 2 1	8 6 6 6 4	177.693 <sup>a</sup> 491.123 594.940 315.102 <sup>b</sup> 599.174

TABLE 6.6. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE PETAL LENGTH

TABLE 6.7. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE PETAL WIDTH

Alternative	Clustering	nlog <sub>e</sub> (2π)	nlog <sub>e</sub> SSW/n	n	k	2(k+1)	AIC
1 2 3 4 5	(S) (Ve) (Vi) (S, Ve) (Vi) (S, Vi) (Ve) (Ve, Vi) (S) (S, Ve, Vi)	275.681 275.681 275.681 275.681 275.681 275.681	-478.966 -216.942 - 84.552 -314.688 - 82.452	150 150 150 150 150	3 2 2 2 1	8 6 6 6 4	-45.285 <sup>a</sup> 214.739 347.129 116.993 <sup>b</sup> 347.229

AIC =  $nlog_e(2\pi)$  +  $nlog_e$  SSW/n + n + 2(k+1)

<sup>a</sup>First Minimum AIC

**b**Second Minimum AIC

Looking at each of the tables above, we see that on each of the variables the first minimum AIC occurs at the alternative submodel 1, namely (S) (Ve) (Vi). That is, the MAICE is submodel 1 indicating that indeed there are three types of species across all the variables. But the second minimum AIC is at the alternative submodel 4 again across all the variables indicating that if we were to cluster any iris species, we should cluster <u>I. versicolor</u> and <u>I. virginica</u> together, as one homogeneous group.

Thus our minimum AIC results for each of the variables confirm other investigators' findings, including Fisher's results on the iris data. Moreover, if we were to choose among the submodels then we would choose the one with smallest minimum AIC as the best submodel. Examining the Tables 6.4, 6.5, 6.6, and 6.7, we see that the smallest minimum AIC occurs at the submodel 1 in Table 6.7 on variable petal width. This indicates that petal width alone separates the three iris species with virtual certainty, confirming again Fisher's results (see, e.g., Fisher [10]).

# 6.2. A Multivariate Example

Now we consider Fisher iris data again and this time we cluster K=3 groups or species into k=1, 2, and 3 homogeneous groups on the basis of all the four variables, assuming the MANOVA model as the underlying model for comparisons of these three iris groups. On the iris data, running SPSS MANOVA program, we obtain the following "within-group" sum of squares and products (SSP) matrices for each of the clustering alternatives. These are:

(1) (S) (VE) (VI) 
$$\underline{W}_{1} = \begin{bmatrix} 39.462 & 13.818 & 24.729 & 5.6554 \\ 13.818 & 16.962 & 8.1208 & 4.8084 \\ 24.729 & 8.1208 & 27.223 & 6.2718 \\ 5.6554 & 4.8084 & 6.2718 & 6.1566 \end{bmatrix}$$

(2) (S, VE) (VI) 
$$\underline{W}_{2} = \begin{cases} 60.714 & -1.3489 & 89.222 & 30.549 \\ -1.3489 & 27.786 & -37.906 & -12.958 \\ 89.222 & -37.906 & 222.94 & 81.818 \\ 30.549 & -12.958 & 81.818 & 35.317 \\ 150 \log_{e} |150^{-1}\underline{W}_{2}| = -1,085.9 \end{cases}$$

-26-

(3) (S, VI) (VE) 
$$\underline{W}_{3} = \begin{bmatrix} 101.52 & -4.3257 & 186.38 & 76.044 \\ -4.3257 & 22.115 & -38.301 & 15.395 \\ 186.38 & -38.301 & 445.43 & 188.28 \\ 76.044 & -15.395 & 188.28 & 85.367 \end{bmatrix}$$
  
150  $\log_{e}|150 \ \underline{W}_{3}| = -988.39$ 

$$(4) (VE, VI) (S) \qquad \underline{W}_{4} = \begin{bmatrix} 50.352 & 17.184 & 46.047 & 17.205 \\ 17.184 & 18.002 & 14.71 & 8.3784 \\ 46.047 & 14.71 & 68.954 & 28.882 \\ 17.205 & 8.3784 & 28.882 & 18.407 \\ 150 \log_{e} | 150^{-1} \underline{W}_{4} | = -1,129.6 \end{bmatrix}$$

(5) (S, VE, VI) 
$$\underline{W}_{5} = \begin{bmatrix} 102.6 & -6.0197 & 189.78 & 76.884 \\ -6.0197 & 28.307 & -49.119 & -18.124 \\ 189.78 & -49.119 & 464.33 & 193.05 \\ 76.884 & -18.124 & 193.05 & 86.57 \end{bmatrix}$$
  
150  $\log_{e} | 150^{-1} \underline{W}_{5} | = -941.73$ 

After carrying out all our computations for each of the clustering alternatives (using the Matrix Algebra Routines in SPEAKEASY interactive computer package), we obtain the AIC's from (5.15). The results are shown in Table 6.8.

-27-

3

Alternative	Clustering	nplog <sub>e</sub> (2π)	_1 nlog <sub>e</sub>  n <u>W</u>	np	k	2m	AIC
1 2 3 4 5	(S) (Ve) (V1) (S, Ve) (V1) (S, V1) (Ve) (Ve, V1) (S) (S, Ve, V1)	1,102.724 1,102.724 1,102.724 1,102.724 1,102.724 1,102.724	-1,504.2 -1,085.9 - 988.39 -1,299.6 - 941.73	600 600 600 600 600	3 2 2 2 1	44 36 36 36 28	242.524 <sup>a</sup> 652.824 750.334 439.124 <sup>b</sup> 788.994

TABLE 6.8. THE AIC'S FOR IRISES BY GROUPS ON ALL VARIABLES

n = 150 plants, p = 4 variables

m = kp + p(p+1)/2 parametersAIC = nplog<sub>e</sub>(2 $\pi$ ) + nlog<sub>e</sub>[n <u>W</u>] + np + 2m

<sup>a</sup>First Minimum AIC

**b**Second Minimum AIC

Hence, looking at the Table 6.8, we see that, using all four variables simultaneously the first minimum AIC occurs at the alternative submodel 1, that is, when (S) (Ve) (Vi) are all clustered separately. This indicates again that indeed there are three types of species. Therefore, the MAICE is submodel 1. Not surprisingly, the second minimum AIC occurs at the alternative submodel 4 telling us that if we were to cluster any one of the two iris groups, we should cluster <u>I. veriscolor</u> and <u>I. virginica</u> together as one homogeneous group, and we should cluster <u>I. setosa</u> completely separate as one heterogeneous group.

Here, it is important to note that we obtained also the same results when we used the four variables separately in our computation of AIC in the previous section, which is encouraging.

Thus, in concluding, we see from these numerical results that AIC and consequently minimum AIC procedures are very successful indeed in identifying

the best clustering alternatives when we cluster samples into homogeneous sets both in the univariate and the multivariate cases.

Moreover, the definition of MAICE gives a clear formulation of the principle of parsimony in statistical model building or comparison as the above examples demonstrate. And MAICE provides a versatile procedure for statistical model identification which is free from the ambiguities inherent in the application of conventional statistical procedures.

# Acknowledgement

This paper is a summary of some of the results in the Ph.D. thesis of the first author in the Department of Mathematics at the University of Illinois at Chicago, under the supervision of the second author. The first author is grateful to Professor Stanley L. Sclove for his valuable and useful comments.

#### REFERENCES

- [1] Ambramowitz, M., and Stegun, I.A. (1968), <u>Handbook of Mathematical</u> <u>Functions with Formulas, Graphs and Mathematical Tables</u> (Nat. Bur. of Stand. Appl. Math. Ser., No. 55), 7th printing. U.S. Govt. Printing Office, Washington, D.C.
- [2] Akaike, H. (1973), "Information Theory and an Extension of the Maximum Likelihood Principle," <u>2nd International Symposium on Information</u> <u>Theory</u>, eds. B.N. Petrov and F. Csaki, Budapest: Akademiai Kiado, 267-281.
- [3] (1974), "A New Look at the Statistical Model Identification," IEEE Transactions on Automatic Control, AC-19, 716-723.
- [4] Anderson, T.W., and Sclove, S.L. (1978), <u>An Introduction to the Statistical Analysis of Data</u>, Boston: Houghton Mifflin Company.
- [5] Brown, B.W., Jr., and Hollander, M. (1977), <u>Statistics: A Biomedical</u> Introduction, New York: John Wiley.
- [6] Chen, Hwei-Ju, Gnanadesikan, R., and Kettenring, J.R. (1974), "Statistical Methods for Grouping Corporations," <u>Sankhya B</u>, <u>36</u>, 1-28.
- [7] Delany, M.J., and Healy, M.J.R. (1966), "Variation in the White-toothed Shrews (Crocidura spp.) in the British Isles," <u>Proceedings of the</u> <u>Royal Society</u>, <u>B</u>, 164, 63-74.
- [8] Dolkart, R.E., Halpern, B., and Perlman, J. (1971), "Comparison of Antibody Responses in Normal and Alloxan Diabetic Mice," <u>Diabetes</u>, <u>20</u>, 162-167.
- [9] Duran, B.S., and Odell, P.L. (1974), <u>Cluster Analysis: A Survey</u>, New York: Springer-Verlag.
- [10] Fisher, R.A. (1936), "The Use of Multiple Measurements in Taxonomic Problems," <u>Annals of Eugenics</u>, 7, 179-188.
- [11] Friedman, H.P., and Rubin, J. (1967), "On Some Invariant Criteria for Grouping Data," <u>Journal of the American Statistical Association</u>, <u>62</u>, 1159-1178.
- [12] Heidorn, K.C. (1979), "Synoptic Weather Patterns Associated with Nitrates in Suspended Particulate in Southern Ontario," <u>Water, Air, and Soil</u> <u>Pollution, 11</u>, 225-235.
- [13] Kendall, M.G. (1966), "Discrimination and Classification," in P.R. Krishnaiah (Ed.), Multivariate Analysis, New York: Academic Press.

- [14] Mezzich, J.E., and Solomon, H. (1980), <u>Taxonomy and Behavioral Science</u>, New York: Academic Press.
- [15] Solomon, H. (1971), <u>Numerical Taxonomy</u>, <u>Machimatics</u> in the Archaeological and Historical Sciences, Edinburgh: Edinburgh University Press, 62-81.
- [16] Williams, W.H., and Goodman, M.L. (1971), "A Statistical Grouping of Corporations by Their Financial Characteristics," <u>Journal of Financial</u> and Quantitative Analysis, 1095-1104.
- [17] Wilks, S.S. (1932), "Certain Generalizations in the Analysis of Variance," <u>Biometrika</u>, <u>24</u>, 471-494.

REFUR I MULLIDER I A I RUR FAME	READ DISTRUCTIONS
ALFONT NUMBER	IL ACCIPIENT'S CATALOS NUMBER
Technical Report 82-1 $/$	49
TITLE (and Subliday	S. TYPE OF REPORT & PERIOD COVI
Multi-Sample Cluster Analysis Using Akaike's	Technical Report
Information Criterion	6. PERFORMING ORS. REPORT NUMB
Authodya	B. CONTRACT ON GRANT HUMBERY OF
namparsum buzdugan and scanley L. Schove	N00014-80-C-0408
Hamparsum Bozdogan and Stanley L. Sclove       N00014-80-C-04         • PERFORMING ORGANIZATION NAME AND ADDRESS       III. PROGRAM ELEMENT.         University of Illinois at Chicago Circle       AMEA & WORK UNIT N         University of Illinois at Chicago Circle       III. REPORT DATE         Box 4348, Chicago, IL 60680       III. REPORT DATE         January 30, 15       III. REPORT PAGES         Unclassified       III. RECURTY CLASS. (or Unclassified         Statistics and Probability Branch       IIII. RECURTY CLASS. (or Unclassified         Arlington, VA 22217       III. REPORT         APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNLIMITED.         7. DISTRIBUTION STATEMENT (of the debuted entered in Block 26, If different free Report)         Unlimited distribution	18. PROGRAM EL ENENT, PROJECT, T
University of Illinois at Chicago (ingle	
Box 4348 Chicago, IL 60680	
- CONTROLLING OFFICE NAME AND ADDRESS	12. REPORT DATE
· · · · · · · · · · · · · · · · · · ·	January 30, 1982
	32+ii
MONITORING AGENCY HANE & ADDRESSII different frem Controlling Office	) 18. SECURITY CLASS. (of shis report)
Office of Naval Research	UNCIASSITIED
Statistics and Probability Branch	IL DECLASSIFICATION DOWNERAD
Ariington, VA 22217	
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above ontored to Stock 20, 11 alternate Unlimited distribution	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the abovest entered in Stock 20, 11 allerent Unlimited distribution	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the abovest entered to Stock 20, 11 atternet Unlimited distribution	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above onlored to Stock 20, 11 allorer Unlimited distribution	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above of an electric to Electr 20, 11 allocant Unlimited distribution L SUPPLEMENTARY NOTES	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the choicest entered in Sheek 20, 11 efforement Unlimited distribution SUPPLEMENTARY HOTES EXEV SORDS (Continue on revuese olds 11 necessary and identify by block made Multistample cluster analysis: Akaike's Infor	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above of the Block 20, 11 different Unlimited distribution SUPPLEMENTARY NOTES REV SORDS (Continue on revues of a 11 eccessory and identify by block mark Multi-Sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the abstract entered to Effect 20, 17 different Unlimited distribution SUPPLEMENTARY HOTES KEY FORDS (Continue on revuew olds 17 necessary and identify by block made Multi-sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above of the Stock 26, 17 Allorent Unlimited distribution Unlimited distribution SUPPLEMENTARY NOTES KEY CORDS (Continue on revuese of it accessory and identity by Most man Multi-Sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL OISTRIBUTION STATEMENT (of the above at an above to stored to store 20, 11 allorent Unlimited distribution Unlimited distribution Unlimited distribution E SUPPLEMENTARY NOTES EXEV SORDS (Continue on reverse olds II according and identify by block main Multi-Sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above and in Stock 20, 11 alternate Unlimited distribution Unlimited distribution Unlimited distribution E SUPPLEMENTARY HOTES KEY SORDS (Continue on reverse olds 11 according on Identify by Mont media Multi-Sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m Multi-sample cluster analysis, the problem of Multi-sample cluster analysis, the problem of	IMITED.
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above entired in Steek 20, 11 different Unlimited distribution Unlimited distribution Unlimit	IMITED. IMI
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the deduced entered in Block 20, If different Unlimited distribution Unlimited distribution Unlim	IMITED. Mainten mation aximum likelihood grouping samples, is stud kaike's Information Criter alue of the likelihood wit at value. The multi-sampl
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the above entered in Block 20, If different Unlimited distribution Unlimited distribution SUPPLEMENTARY HOTES EXEV CORCE (Continue on reviewe and If anothy by Most main Multi-sample cluster analysis; Akaike's Infor Criterion (AIC); ANOVA Model; MANOVA Model; m Multi-sample cluster analysis, the problem of from an information-theoretic viewpoint via A (AIC). This criterion combines the maximum v the number of parameters used in achieving th cluster problem is defined, and AIC is develo	IMITED. IMI
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL DISTRIBUTION STATEMENT (of the aboves entered in Steek 20, if different Unlimited distribution Unlimited distribution Multi-sample cluster analysis, the problem of from an information-theoretic viewpoint via A (AIC). This criterion combines the maximum view of parameters used in achieving the cluster problem is defined, and AIC is develo form of AIC is derived in both univariate and vaniance models Wreaterion	IMITED. IMITED. mation aximum likelihood grouping samples, is stud kaike's Information Criter alue of the likelihood wit at value. The multi-sampl ped for this problem. The multivariate analysis of ented and resulte are shown
APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNL CONTINUENTION STATEMENT (of the electron entered in Stock 20, 11 eliferent Unlimited distribution CONTINUENTARY HOTES CONTINUENTARY HOTES CONTINUENT CONTINUENTARY HOTES CONTINUENT AND A MODELS AND A MODELS AND A MODELS IN A A MULTI-SAMPLE CLUSTER ANALYSIS, THE PROBLEM OF FORM AN INFORMATION-THEORETIC VIEWPOINT VIA A (AIC). This criterion combines the maximum v the number of parameters used in achieving th Cluster problem is defined, and AIC is develo form of AIC is derived in both univariate and variance models. Numerical examples are pres	IMITED. IMI

.

i

-

1

.

7		
to demonstrate the utility alternatives.	of AIC in identifying the best clustering	
$\setminus$ .		
	-	
	•	
		•
· · · · ·		
	• •	
	•	

٠

٠ ٠

