





College of Business Administration, University of Illinois at Chicago Circle

3

ALT 3422 AC

.

.



FILE COPY DIIC

ł 1 .

AFE Ł

18.2

4

The documant & is trasting mins ! for public rateries the same distribution is aclimited.

1.1 82

٠

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. 82-2 March 8, 1982

PREPARED FOR THE OFFICE OF NAVAL RESEARCH UNDER CONTRACT N00014-80-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

VV88

^{*}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 Fisher Iris Data
 - 6.1. A Univariate Example
 - 6.2. A Multivariate Example
- 7. Conclusions and Discussions

Acknowledgement

References

	Accession For	-
	NTIS URAAL DTIN 114 Uncurrend Justifientien	-
	By Distribution/	
	Avail 3 Avail 5 Dist Special	
L	A	

*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

ABSTRACT

Multi-sample cluster analysis, the problem of grouping samples, is studied from an information-theoretic 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 the univariate model with varying means and variances, and in the multivariate model with varying mean vectors and variance-covariance matrices. 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); Univariate model with varying means and variances, Multivariate model with varying mean vectors and variance-covariance matrices; maximum likelihood.

-11-

^{*}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

1. Introduction

In a previous paper, we introduced and developed Akaike's Information Criterion (AIC) for multi-sample cluster analysis, the problem of grouping samples. We derived the form of AIC in both univariate and multivariate analysis of variance models, where the assumptions of independence, univariate and multivariate normality, equal variances and variance-covariance matrices were fundamental to the analysis. We gave numerical examples and results which demonstrated the utility of AIC in identifying the best clustering alternatives. (See, Bozdogan and Sclove [5]).

In this paper, we shall continue to study the multi-sample cluster problem. However, here we shall develop Akaike's Information Criterion (AIC) for multi-sample cluster analysis with varying means and variances in the univariate model, and with varying mean vectors and variance-covariance matrices in the multivariate model, since often in practice the assumption of equal parameters within the model is a rather dubious requirement.

Many practical situations require the presentation of multivariate data from several structured samples for comparative purposes and the grouping of the heterogeneous samples into homogeneous sets of samples in which parameters might vary. For this reason 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

*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. should not. For examples of multi-sample clustering situations, we refer the reader to Bozdogan and Sclove [5].

In the statistical literature several conventional test procedures are available for testing whether or not several populations have equal variances, as required by the analysis of variance (ANOVA) model. If we have a reason to doubt this is the case, then we may first want to test the equality of variances. In the multivariate case the equality of covariance matrices is certainly more hazardous. Therefore, for this reason we may want first to test the equality of variances in the univariate case, and the equality of covariance matrices in the multivariate case. This is an important option to use in clustering groups or samples, and in general.

In the literature, however, there are several test procedures for testing the equality of variances, and covariance matrices. For example, in the multivariate case, the most commonly used test is <u>Box's M test</u> despite the fact that it is very expensive to compute it on a high speed computer, even on an IBM 370. Moreover, as in the case of MANOVA, these test procedures are not revealing or informative in multi-sample clustering problems. Therefore, in this paper we shall propose again 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

-2-

(1.1) AIC = $-2 \log_e$ (maximized likelihood)

+2 (number of independently adjusted parameters within the model). 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. The AIC avoids this spurious improvement 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 best model.

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

-3-

subsequent sections, that is, in Section 4 and in 5, we shall derive the AIC procedure for the univariate model with varying means and variances, and for the multivariate model with varying mean vectors and covariance matrices. In Section 6, we shall give numerical examples for both univariate and multi-variate multi-sample cluster analysis on the same real data set 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)
$$\underline{X} (n \times p) = \begin{cases} \underline{X}_1 \\ \underline{X}_2 \\ \vdots \\ \vdots \\ \underline{X}_K \\ \underline{X}_K \\ \underline{X}_K \end{cases}$$

be a single data matrix of K groups or samples, where \underline{X}_{g} (ngxp) represents the observations from the g-th group or sample, g=1,2,...,K, and n = $\sum_{\substack{g=1 \\ g=1}}^{K}$ ng. 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

-4-

restricted to being equal or proportional to other n_g 's. The total number of Kobservations is $n = \sum_{g=1}^{n_g}$. Let X_{gi} be the px1 vector of observations in group g=1 $g=1,2,\ldots,K$, and for individual $i=1,2,\ldots,n_g$.

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

In this section, we shall just briefly mention how to obtain the total number of clustering alternatives for a given K, the number of <u>groups</u> or samples. For details we again refer the reader to Bozdogan and Sclove [5].

In general, the total number of ways of clustering K groups or samples into k clusters is given by

(3.1)
$$S(K,k) = \frac{1}{k!} \sum_{g=0}^{k} {\binom{k}{g}} {(-1)^{g}} {(k-g)^{K}}$$

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

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

(3.2)
$$\sum_{k=1}^{K} S(K,k)$$

For example, K=4 samples, if k is not specified a priori and varies, then there are in total 15 possible clustering alternatives to cluster K=4 groups or samples first into k=4 groups or samples, then k=3 groups or samples, k=2 groups or samples, and k=1 group or sample by using the equation (3.1)

-5-

respectively, and summing the results by using the expression (3.2) to obtain 15 as the total number of possible clustering alternatives.

Therefore, the total number of ways of clustering K groups or samples into k homogeneous groups or samples is given by equation (3.1), and the total number of possible clustering alternatives is given by the expression (3.2).

4. AIC For The Univariate Model

We now turn our attention to consider situations with several univariate normal samples. 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 μ_1 and variance σ_1^2 , the second with mean μ_2 and variance σ_2^2, \ldots , the Kth with mean μ_K and variance σ_K^2 . We may want to decide in this case if the variances of these K samples will be treated as equal or not, given no restriction on the population means. In terms of the parameters the univariate model is $\theta = (\mu_1, \mu_2, \ldots, \mu_K, \sigma_1^2, \sigma_2^2, \ldots, \sigma_K^2)$ with m=2k parameters, where k is the number of groups.

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 K n_{g} , $g=1,2,\ldots,K$, observations in the g-th group and $n = \sum_{\substack{g=1 \\ g=1}} n_{g}$. Denote the g=1unknown means of the groups by $\mu_{1},\mu_{2},\ldots,\mu_{K}$, and the unknown variances of the groups by $\sigma_{1}^{2},\sigma_{2}^{2},\ldots,\sigma_{K}^{2}$. Assume that the samples $(z_{11},z_{12},\ldots,z_{1n_{1}};\ldots;$

-6-

z ,...,z) are drawn randomly from K populations which are $N(\mu_g,\sigma_g)$. The K_1

basic null hypothesis of interest is given by

(4.1)
$$H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2$$
.

The alternative hypothesis is given by

H₁: Not all K variances are equal.

In the statistical literature, this is also known as the test of <u>homogeneity of</u> variances or <u>Bartlett's test</u>.

To derive the form of AIC subject to this constraint, we call the common unknown value of variances σ^2 . The likelihood function in this case is given by

(4.2)
$$L(\{\mu_{g},\sigma_{g}^{2}\};z) = (2\pi)^{-n/2} \frac{K}{\pi} (\sigma_{g}^{2})^{-n} g^{/2} \exp\{-\frac{K}{\sum_{g=1}^{K} (\frac{1}{2\sigma_{g}^{2}}) \sum_{i=1}^{ng} (z_{gi} - \mu_{g})^{2}.$$

The log likelihood is

$$(4.3) \quad l(\{\mu_{g}, \sigma^{2}\}) \equiv \log L(\{\mu_{g}, \sigma^{2}_{g}\}; z)$$

$$= -\frac{n}{Z} \log(2\pi) - 1/2 \sum_{g=1}^{K} n_{g} \log \sigma^{2}_{g} - \sum_{g=1}^{K} \frac{1}{2\sigma^{2}_{g}} \sum_{i=1}^{n_{g}} (z_{gi} - \mu_{g})^{2},$$

and the MLE's of μ and σ^2 are given by g

(4.4)
$$\hat{\mu}_{g} = \frac{1}{n_{g}} \sum_{i=1}^{n_{g}} z_{gi} = \bar{z}_{g.}$$
,

and

(4.5)
$$\hat{\sigma}_{g}^{2} = \frac{1}{n} \sum_{i=1}^{n_{g}} (z_{gi} - \bar{z}_{g.})^{2} = s_{g}^{2}, g=1,2,\ldots,K.$$

Substituting these back into (4.3) and simplifying, the maximized log likelihood becomes

(4.6)
$$l(\{\hat{\mu}_{g}, \hat{\sigma}_{g}^{2}\}; z) \equiv log L(\{\hat{\mu}_{g}, \hat{\sigma}_{g}^{2}\}; z)$$

= $-\frac{n}{2} log(2\pi) - 1/2 \sum_{g=1}^{K} n_{g} log s_{g}^{2} - \frac{n}{2}$.

Since

(4.7) AIC =
$$-2 \log_e L(\frac{0}{2}) + 2m$$
,

where m is the number of parameters, and since

(4.8) -2 log L(
$$\{\mu, \sigma^2\}$$
) = n log(2π) + $\sum_{g=1}^{K}$ ng log s² + n,

then AIC becomes

(4.9) AIC (varying
$$\mu$$
 and σ^2) = n log(2π) + $\sum_{g=1}^{K}$ ng log s² + n + 2(2k).

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

(4.10) AIC* = (varying
$$\mu$$
 and σ^2) = $\sum_{g=1}^{K} n_g \log_e S^2 + 2(2k)$

where $S_g^2 = \frac{1}{n_g} \sum_{i=1}^{n_g} (z_{gi} - \bar{z}_{g.})^2$, g=1,2,...,K,

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 given by (4.9).

5. AIC For the Multivariate Model

As we mentioned in Section 1, that the assumption of equality of variances in one-way ANOVA, causes serious problems when we are testing the equality of several means. Parallel to this assumption, in the multivariate case the equality of covariance matrices even causes more serious problems. For this reason we may want first to test the equality of covariance matrices against the alternative that not all covariance matrices are equal. Therefore, throughout this section we shall suppose that we may have independent data matrices X_1, X_2, \ldots, X_K , where the rows of X_g (ngxp) are independent and identically distributed (i.i.d.) N_p($\mu_g, \underline{\Sigma}_g$), g=1,2,...,K. In terms of the parameters the multivariate model we shall consider is

 $\stackrel{\theta}{=} (\underset{a_1}{\underline{\mu}}, \underset{a_2}{\underline{\mu}}, \ldots, \underset{a_K}{\underline{\mu}}, \underset{\underline{\Sigma}_1}{\underline{\Sigma}}, \underset{\underline{\Sigma}_2}{\underline{\Sigma}}, \ldots, \underset{\underline{\Sigma}_K}{\underline{\Sigma}})$

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

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

(5.1) $H_0: \underline{\Sigma}_1 = \underline{\Sigma}_2 = \dots = \underline{\Sigma}_{\nu}.$

The alternative hypothesis is given by

H₁: Not all K covariance matrices are equal.

-9-

In multivariate analysis this is known as the <u>test of homogeneity of</u> <u>covariance matrices</u>.

To derive Akaike's Information Criterion (AIC) in this case the log likelihood function is given by

(5.2)
$$l(\{\underline{\mu}g,\underline{\Sigma}g\};\underline{Z}) \equiv log L(\{\underline{\mu}g,\underline{\Sigma}g\};\underline{Z})$$
$$= -\frac{np}{2} log(2\pi) - 1/2 \sum_{g=1}^{K} n_g log|\underline{\Sigma}g| - 1/2 \sum_{g=1}^{K} n_g tr \underline{\Sigma}g^{-1}\underline{A}g$$
$$- 1/2 \sum_{g=1}^{K} n_g(\overline{Z}g - \underline{\mu}g)^{*}(\overline{Z}g - \underline{\mu}g) \cdot .$$

The MLE's of $\underline{\mu}_g$ and $\underline{\Sigma}_g$ are

(5.3)
$$\hat{\mu}_{g} = \bar{z}_{g}, g=1,2,\ldots,K,$$

and

(5.4)
$$\hat{\Sigma}_{g} = \underline{A}_{g}/n_{g}.$$

Substituting these back into (5.2) and simplifying, the maximized log likelihood becomes

(5.5)
$$1(\{\underline{\mu}_{g}, \underline{\Sigma}_{g}\}; \underline{Z}) \equiv \log L(\{\underline{\mu}_{g}, \underline{\Sigma}_{g}\}; \underline{Z})$$

$$= -\frac{np}{2}\log(2\pi) - \frac{1}{2}\sum_{g=1}^{K} \log \left| \log \left| n_g^{-1} \underline{A}_g \right| - \frac{np}{2}.$$

1

-11-

Since

(5.6) AIC =
$$-2 \log_e L(\theta) + 2m$$

where m = kp + kp(p+1)/2 is the number of parameters, then AIC becomes

(5.7) AIC(varying
$$\mu$$
 and $\underline{\Sigma}$) = nplog(2π) + $\sum_{g=1}^{K} n_g \log |n_g^{-1}\underline{A}_g|$ + np
g=1 + 2[kp + kp(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

(5.8) AIC*(varying
$$\mu$$
 and $\underline{\Sigma}$) = $\sum_{g=1}^{K} n_g \log_e |\underline{A}_g| + 2[kp + kp(p+1)/2],$

where n_{cr} = sample size of group or sample g=1,2,...K,

 $|\underline{A}_g|$ = the determinant of sum of squares and cross-products (SSCP) matrix for group or sample g=1,2,...,K,

k = number of groups or samples compared, and

p = number of variables.

However, for purposes of comparison we retain the constants and use AIC given by (5.7).

6. Numerical Examples of Multi-Sample Cluster Analysis on Fisher Iris Data

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. A Univariate Examples

For the univariate numerical examples we shall illustrate our results on Fisher [6] iris data.

Example 6.1. Clustering of Irises by Groups: 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 6.1.

This data set has been quite extensively studied in classification and cluster analysis since it was published by Fisher [6], and still today, is being used as a "testing ground" for classification and clustering methods proposed by many investigators such as Friedman and Rubin [7], Kendall [8], Solomon [10], Mezzich and Solomon [9], 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.

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 using equation (3.1) and the

-12-

Sapal Sepal Petal Petal Petal Sepal Sepal Sepal Petal Petal Sepal Sepal Sepal Petal P	apal Sepal Petal Peta
5.1 3.5 1.4 0.2 7.0 3.2 4.7 1.4 6.3 3.3 6.0 2 4.9 3.0 1.4 0.2 6.4 3.2 4.5 1.5 5.8 2.7 5.1 1 4.7 3.2 1.3 0.2 6.9 3.1 4.9 1.5 7.1 3.0 5.9 2	מקרה שומרה ופחקרה שומו
4.63.11.30.25.32.34.01.30.36.32.95.05.03.61.40.26.52.84.61.56.53.06.624.63.41.40.36.33.34.71.64.92.54.515.03.41.50.24.92.43.31.07.32.96.314.42.91.40.26.62.94.61.36.72.55.814.33.71.50.25.02.03.51.06.53.25.124.83.41.60.25.93.04.21.56.42.75.314.33.01.10.16.12.93.61.35.82.55.025.84.01.20.25.62.93.61.35.85.125.74.41.50.46.63.04.51.56.53.05.515.13.50.35.82.74.11.07.73.66.9225.13.61.00.25.93.24.81.86.93.25.725.13.71.81.70.36.22.24.51.57.72.66.925.73.81.70.35.82.7 <td>5.1 3.5 1.4 0.3 4.9 3.0 1.4 0.3 4.6 3.1 1.5 0.3 4.6 3.4 1.5 0.3 5.4 3.9 1.7 0.4 5.0 3.6 1.4 0.3 4.6 3.4 1.5 0.3 4.6 3.4 1.5 0.3 4.4 2.9 1.4 0.3 4.9 3.1 1.5 0.3 4.3 3.0 1.4 0.3 4.3 3.0 1.4 0.3 4.3 3.0 1.4 0.3 5.7 4.4 1.5 0.3 5.7 4.4 1.5 0.3 5.7 3.8 1.7 0.5 5.1 3.6 1.0 0.3 5.1 3.6 1.0 0.3 5.1 3.6 1.0 0.3 5.2 3.4 1.6</td>	5.1 3.5 1.4 0.3 4.9 3.0 1.4 0.3 4.6 3.1 1.5 0.3 4.6 3.4 1.5 0.3 5.4 3.9 1.7 0.4 5.0 3.6 1.4 0.3 4.6 3.4 1.5 0.3 4.6 3.4 1.5 0.3 4.4 2.9 1.4 0.3 4.9 3.1 1.5 0.3 4.3 3.0 1.4 0.3 4.3 3.0 1.4 0.3 4.3 3.0 1.4 0.3 5.7 4.4 1.5 0.3 5.7 4.4 1.5 0.3 5.7 3.8 1.7 0.5 5.1 3.6 1.0 0.3 5.1 3.6 1.0 0.3 5.1 3.6 1.0 0.3 5.2 3.4 1.6

TABLE 6.1

. '

MEASUREMENTS ON THREE TYPES OF IRIS

 \mathcal{I}

-13-

expression (3.2) in Section 3, we obtain in total 5 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. In terms of the parameters, using the univariate model $\underline{\theta} = (\mu_1, \mu_2, \dots, \mu_K, \sigma_1^2, \sigma_2^2, \dots, \sigma_K^2)$ as our underlying model with varying means and variances for clustering the iris groups, from a simple run on the computer by using MINITAB package, we obtained the AIC's for each of the 5 clustering alternatives of each of the four variables separately. We report our results on each of the four variables, respectively, as follows.

TABLE 6.2. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE SEPAL LENGTH

Alternative	Clustering	nlog _e (2π)	K 2 ∑nglog _e Sg g=1	n	k	2(2k)	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	-218.710 -136.019 - 79.394 -188.536 - 57.603	150 150 150 150 150	3 2 2 2 1	12 8 8 8 4	218.971 ^a 297.662 354.287 245.145 ^b 372.078

TABLE 6.3. THE AIC'S FOR IRISES BY GROUPS ON VARIABLE SEPAL WIDTH

Alternative	Clustering	nlog _e (2π)	K 2 ∑nglogeSg g=1	n	k	2(2k)	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	-329.102 -262.503 -292.416 -319.093 -250.132	150 150 150 150 150	3 2 2 2 1	12 8 8 8 4	108.579 ^a 171.178 141.265 114.588 ^b 179.549

-14-

Alternative	Clustering	nlog _e (2 n)	K 2 ∑nglog _e Sg g=1	n	k	2(2k)	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	-313.055 12.795 70.394 -215.414 169.888	150 150 150 150 150	3 2 2 2 1	12 8 8 8 4	124.626 ^a 446.476 504.075 218.267 ^b 599.569

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

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

Alternative	Clustering	nlog _e (2π)	K 2 ∑nglog _e Sg g=1	n	k	2(2k)	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	-519.344 -245.374 -181.176 -398.271 - 82.454	150 150 150 150 150	3 2 2 2 1	12 8 8 8 4	-81.663 ^a 188.307 252.505 35.410 ^b 347.227

AIC(varying
$$\mu$$
 and σ) = nlog_e(2 π) + $\sum_{g=1}^{k}$ n_glog_e S_g² + n + 2(2k)
g=1

^aFirst Minimum AIC

bSecond 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 6.2, 6.3, 6.4, and 6.5, we see that the smallest minimum AIC occurs at the submodel 1 in Table 6.5 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 [6]).

Hence, we note here that we are clustering the irises by groups or species under a more general model rather than using the ANOVA model as our underlying model which we cosidered in a previous paper on multi-sample cluster analysis.

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 multivariate model given in terms of the parameters $\theta = (\mu_1, \mu_2, \dots, \underline{\Sigma}_1, \underline{\Sigma}_2, \dots, \underline{\Sigma}_K)$ as the underlying model for clustering these three iris groups or species. On the iris data, running SPSS MANOVA program, we obtain the following sum of squares and products (SSCP) matrices for each of the clustering alternatives. These are:

(1) (S) (VE) (VI)

$$\underline{A}_{(S)} = \begin{bmatrix} -1 & -1 \\ 6.0882 & 4.8616 & .8014 & .5062 \\ 4.8616 & 7.0408 & .5732 & .4556 \\ .8014 & .5732 & 1.4778 & .2974 \\ .5062 & .4556 & .2974 & .5442 \end{bmatrix} \begin{bmatrix} -1 \\ |50 \\ \underline{A}_{(S)}| = 1.949E-6 \\ 10g_e & (1.949E-6) = -13.148 \end{bmatrix}$$

-16-

$$\underline{A}_{(VE)} = \begin{cases} 13.561 & 4.362 & 9.066 & 2.7436 \\ 4.362 & 4.825 & 4.05 & 2.019 \\ 9.066 & 4.05 & 10.82 & 3.582 \\ 2.7436 & 2.019 & 3.582 & 1.9162 \end{cases} \xrightarrow{-1}_{[50]} \underbrace{|50]_{A_{(VE)}}| = 1.8053E-5}_{10g_e} (1.8053E-5) = -10.922 \end{cases}$$

$$\underline{A}_{(VI)} = \begin{bmatrix} 19.813 & 4.5944 & 14.861 & 2.4056 \\ 4.5944 & 5.0962 & 3.4976 & 2.3338 \\ 14.861 & 3.4976 & 14.925 & 2.3924 \\ 2.4056 & 2.3338 & 2.3924 & 3.6962 \end{bmatrix} \begin{bmatrix} -1 \\ |50 \\ A \\ (VI)| = 1.2244E-4 \\ 10g_e (1.2244E-4) = -9.0079 \end{bmatrix}$$

 $\frac{A}{(S, VE)} = \begin{bmatrix} 40.901 & -5.9433 & 74.361 & 28.144 \\ -5.9433 & 22.69 & -41.404 & -15.291 \\ 74.361 & -41.404 & 208.02 & 79.425 \\ 28.144 & -15.291 & 79.425 & 31.62 \end{bmatrix} \begin{bmatrix} -1 \\ |100 & \underline{A}_{(S, VE)}| &= 3.3118E-4 \\ 10g_e & (3.3118E-4) &= -8.0128 \end{bmatrix}$

(3) (S, VI) (VE)

. '

$$\frac{A}{(S, VI)} = \begin{pmatrix} 88.469 & -8.4997 & 177.42 & 73.311 \\ -8.4997 & 17.29 & -42.351 & -17.414 \\ 177.42 & -42.351 & 434.61 & 184.69 \\ 73.311 & -17.414 & 184.69 & 83.45 \end{pmatrix} \begin{pmatrix} -1 \\ |100 & \underline{A}_{(S, VI)}| = .0025193 \\ 10g_{e} (.0025193) = -5.9838 \end{pmatrix}$$

-17-

 Γ

(4) (VE, VI) (S) $\frac{A}{(VE, VI)} = \begin{bmatrix}
44.264 & 12.322 & 45.245 & 16.699 \\
12.322 & 10.962 & 14.137 & 7.9228 \\
45.245 & 14.137 & 67.476 & 28.584 \\
16.699 & 7.9228 & 28.584 & 17.862
\end{bmatrix} \begin{vmatrix}
100 & \underline{A}_{(VE, VI)} & = 3.1476E-4 \\
10g_e & (3.1476E-4) & = -8.0637
\end{cases}$ (5) (S, VE, VI) $\frac{A}{(S, VE, VI)} = \begin{bmatrix}
102.6 & -6.0197 & 189.78 & 76.884 \\
-6.0197 & 28.307 & -49.119 & -18.124 \\
189.78 & -49.119 & 464.33 & 193.05 \\
76.884 & -18.124 & 193.05 & 86.57
\end{bmatrix} \begin{vmatrix}
-1 \\
150 & \underline{A}_{(S, VE, VI)} & = .0018787 \\
10g_e & (.0018787) & = -6.2772
\end{cases}$

After carrying out all our computations for each of the clustering alternatives (using the Matrix Algebra Routines in the SPEAKEASY interactive computer package), we obtain the AIC's from (5.7). The results are shown in Table 6.6.

IADLE 0.0. INE AIL 3 FUR IRISES BY GROUPS UN ALL VARI.
--

Alternative	Clustering	nplog _e (2π)	K _1 ∑nglog _e ng <u>Ag</u> g=1	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	-1,653.89 5 -1,251.675 -1,144.480 -1,463.770 - 941.580	600 600 600 600 600	3 2 2 2 1	84 56 56 56 28	132.829 ^a 507.049 614.244 294.954 ^b 789.144

n = 150 plants, p = 4 variables

m = kp + kp(p+1)/2 parameters

AIC(varying μ and $\underline{\Sigma}$) = nplog_e (2 π) + $\sum_{g=1}^{K}$ nglog_e|n_g A_g| + np + 2m g=1 aFirst Minimum AIC

^bSecond Minimum AIC

-18-

Hence, looking at the Table 6.6, 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.

7. Conclusions and Discussion

From our numerical results in Section 6, we see that AIC and consequently minimum AIC procedures can successfully indeed identify the best clustering alternatives when we cluster samples into homogeneous sets of samples both in the univariate and the multivariate models with varying parameters.

In our previous paper on multi-sample cluster analysis (Bozdogan and Sclove [5]), we considered ANOVA and MANOVA as our two underlying models where the assumption of equal variances and covariances were used to cluster the groups or samples for multi-sample data. There, we used AIC also in identifying the best clustering alternatives in clustering the iris groups or species. We obtained the same results in determining the three types of iris species and that if we were to cluster any one of the two iris groups, we should cluster <u>I. versicolor</u> and <u>I. virginica</u> together as one homogeneous

-19-

group, and we should cluster <u>I. setosa</u> completely separate as one heterogeneous group.

In summarizing the results of AIC-values for the multivariate case only from the previous and this paper, we obtain the following table.

Alternative	Clustering	AIC(varying μ and Σ)	AIC(common <u>Σ</u>)
1	(S) (Ve) (Vi)	132.829 ^a	242.524ª
2	(S, Ve) (Vi)	507.049	652.824
3	(S, Vi) (Ve)	614.244	750.334
4	(Ve, Vi) (S)	294.954 ^b	439.124 ^b
5	(S, Ve, Vi)	789.144	788.994

TABLE 6.7. THE AIC'S FOR IRISES BY GROUPS ON ALL VARIABLES UNDER TWO MULTIVARIATE MODELS

aFirst Minimum AIC

^bSecond Minimum AIC

Comparing the AIC's in Table 6.7 above, we see that AIC(varying μ and $\underline{\Sigma}$) values are much less than the AIC(common $\underline{\Sigma}$) values for each of the clustering alternatives except for the last clustering alternative (i.e., alternative 5) in clustering the iris groups or species. Since according to the definition of AIC, the model with the minimum AIC is chosen to be the <u>best model</u>, then the above results suggest that when we are clustering iris data, and in general, we should use different covariance matrices rather than using equal covariance matrices in data analysis.

As we mentioned in the introduction of this paper, in practice the assumption of equal covariance matrices within the model is a rather dubious requirement.

Thus, in concluding, we see that the use of AIC shows how to combine the information in the likelihood with an appropriate function of the number of parameters to obtain estimates of the information provided by competing alternative models. Therefore, the definition of MAICE gives a clear formulation of the principle of parsimony in statistical model building or comparison as we demonstrated by numerical examples. 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

and an Sub line in a second

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.

-22-

REFERENCES

- [1] Ambramowitz, M., and Stegun, I.A. (1968), <u>Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables (Nat. Bur. of Stand. Appl. Math. Ser., No. 55), 7th printing. U.S. Govt. Printing Office, Washington, D.C.</u>
- [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), An Introduction to the Statistical Analysis of Data, Boston: Houghton Mifflin Company.
- [5] Bozdogan, H., and Sclove, S.L. (1982), "Multi-Sample Cluster Analysis Using Akaike's Information Criterion," Technical Report No. 82-1, ONR Contract N00014-80-C-0408 (NR042-443), Quantitative Methods Department, University of Illinois at Chicago Circle, Chicago, Illinois 60680.
- [6] Fisher, R.A. (1936), "The Use of Multiple Measurements in Taxonomic Problems," <u>Annals of Eugenics</u>, <u>7</u>, 179-188.
- [7] 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.
- [8] Kendall, M.G. (1966), "Discrimination and Classification," in P.R. Krishnaiah (Ed.), <u>Multivariate Analysis</u>, New York: Academic Press.
- [9] Mezzich, J.E., and Solomon, H. (1980), <u>Taxonomy and Behavioral Science</u>, New York: Academic Press.
- [10] Solomon, H. (1971), "Numerical Taxonomy," <u>Mathematics in the Archaeological and Historical Sciences</u>, Edinburgh: Edinburgh University Press, 62-81.

-23-

REPORT DOCUMENTA	TION PAGE	READ INSTRUCTIONS
AEPORT NUMBER	1. GOVT ACCESSION NO.	& RECIPIENT'S CATALOG NUMBER
Techical Report 82-2	AD-A113 422	
TITLE (and Substite)		1. TYPE OF REPORT & PERIOD COVERED
Multi-Sample Cluster Analy	sis With Varying	Technical Report
Parameters Using Akaike's	Information	- PERFORMING ORG. REPORT NUMBER
Criterion		
AUTHORICA		8. CONTRACT OR GRANT NUMBER(*)
Hamparsum Bozdogan and Sta	nley L. Sclove	N00014-80-C-0408
PERFORMING ORGANIZATION NAME AND A	ODRESS	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
University of Illinois at	Chicago Circle	
Box 4348, Chicago, IL 6068	0	•
. CONTROLLING OFFICE NAME AND ADDRE	¥	12. REPORT DATE
		March 8, 1982
-		23+jj
SUBNITORING AGENCY NAME & ADDRESS	dillorent from Controlling Office)	15. SECURITY CLASS. (of this report)
Office of Naval Research		Unclassified
Statistics and Probability	Branch	IS. DECLASSIFICATION/DOWNGRADING
Arlington, VA 22217	· · · · · · · · · · · · · · · · · · ·	
APPROVED FOR PUBLIC RELEASE	E: DISTRIBUTION UNLIN	IITED.
APPROVED FOR PUBLIC RELEASE	E: DISTRIBUTION UNLIN	NITED.
APPROVED FOR PUBLIC RELEASE CONTINUENTION STATEMENT (of the element Unlimited distribution	E: DISTRIBUTION UNLIN	NITED.
APPROVED FOR PUBLIC RELEASE CONTRIBUTION STATEMENT (of the domain Unlimited distribution	E: DISTRIBUTION UNLIN	IITED.
APPROVED FOR PUBLIC RELEASE CONSTRUCTION STATEMENT (of the chorese Unlimited distribution - SUPPLEMENTARY NOTES	E: DISTRIBUTION UNLIN autorod in Block 20, 18 different fre	IITED. - Report)
APPROVED FOR PUBLIC RELEASE CONTRIBUTION STATEMENT (of the domain Unlimited distribution BUPPLEMENTARY NOTES KEY CORDS (Contract on course of a ll note Multi-sample cluster analy Univariate model with vary with varying mean vectors likelihood.	E: DISTRIBUTION UNLIN autorod in Block 20, 11 dittores in satured in Block 20, 11 dittores in sis; Akaike's Information ing means and variance and variance-covarian	Ation Criterion (AIC); ces; Multivariate model nce matrices; maximum
APPROVED FOR PUBLIC RELEASE Destrimention statement (of the doment Unlimited distribution B SUPPLEMENTARY NOTES Multi-sample cluster analy Univariate model with vary with varying mean vectors likelihood. AfstRACT (Common on control doe H more Multi-sample cluster analy from an information-theore (AIC). This criterion comb the number of parameters u cluster problem is defined form of AIC is derived in	E: DISTRIBUTION UNLIN metered in Block 20, 11 different from visis; Akaike's Informi- ring means and variance and variance-covarian- dent variance-covarian- visis, the problem of e etic viewpoint via Aka- bines the maximum val- used in achieving tha- i, and AIC is developed the univariate model	Accession Accession ation Criterion (AIC); ces; Multivariate model nce matrices; maximum grouping samples, is studio aike's Information Criterio ue of the likelihood with t value. The multi-sample ed for this problem. The with varving means and

.

· - -

T

•

•

•

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (Then Date Entrol)

variances, and in the multivariate model with varying mean vectors and variance-covariance matrices. Numerical examples are presented and results are shown to demonstrate the utility of AIC in identifying the best clustering alternatives. Γ

5 N 2102- "R. 314- 9601

. '

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE/When Jan Interest

