John A. Sonquist, University of Michigan

## Abstract

Within the past few years computerized sequential algorithms have been developed and used to search out and display non-additive relationships implicit in survey data. The objective is the display of information relevant to the problem of how to specify a multiple classification analysis model purporting to explain the phenomenon in question. Experimentation with one such algorithm (AID) has led to a classification of simple additive and interactive multivariate models related to elementary Boolean operators. These models are described and illustrated; extensions to the AID algorithm to facilitate specification of complex models, for dealing with covariate models arising from "crucial" variables or over-time survey data and for imposing symmetry restrictions are discussed. Preliminary experimentation with the improved algorithm is reported.

Keywords: Survey, Multivariate analysis, model specification, statistical interaction, computer program, prediction, analysis of variance, multiple regression, multiple classification analysis, data analysis strategy, simulation.

\_\_\_\_\_

The task of interpreting a mass of non-experimental data, such as that generated by modern cross-sectional and longitudinal survey methods has remained a difficult one. Two reasons for this appear to stand out as particularly worthy of examination. One is that the statistical methods used by most data analysts are oriented toward deductive testing of isolated hypotheses rather than toward the more frequent task of extraction of information from data and the development of models on an inductive basis. In addition, the full power of "second generation" computing equipment has not yet been utilized (let alone "third generation"); it has the ability to extend the scope of the analyst's logic as well as perform computations rapidly for him.

The present paper is a modest attempt to surmount these obstacles by extending the idea of coupling the calculating power of modern computing equipment to a formal sequential algorithm, the objective of which is the generation of information about the data. In particular, our concern is with obtaining summary information organized for use by the analyst in the inductive stages of his research; that is, in the stages of model formulation and development. The objective of the procedure is the display of information relevant to the problem of how to formulate or specify a model purporting to explain the phenomenon in question.

One frequently appearing and concrete form taken by this general specification problem is the question of choosing and of specifying the form of the terms to be included in a least squares multiple classification equation (MCA). Such an equation is commonly used as a statistical model representing the simultaneous and direct explanatory effects of a correlated set of presumed causal factors on a single dependent variable. It is of the form

 $Y_{ijk...} = \bar{Y} + a_i + b_j + ... e_{ijk...(1)}$ 

where

- Y = the score (on the dependent variable) of individual k, who falls in category i of predictor A and into category j of predictor B, etc.;
- $\overline{Y}$  = the grand mean of the dependent variable;
- a = the "effect" of membership in the i-th category of predictor A expressed as a deviation from the grand mean and adjusted for the intercorrelations of predictor A with the other terms in the equation;
- b = the corresponding "effect" associated with predictor B;
- e = the error term for the k-th ijk...individual.

Any given term in the equation, say b, may represent the effect of a two-way or higher-order interaction term between explanatory factors not<sub>1</sub>otherwise represented in the equation.

For a more detailed exposition of this basic technique, see Andrews, Morgan and Sonquist (1967). A simple n-way analysis of variance model is not used because the number of observations in the cells of an n-way cross-classification of the predictors are not equal or proportional. This condition ordinarily occurs in survey or other non-experimental data in which there are non-zero correlations between explanatory characteristics.

It is this problem of choosing terms to be included in the equation to which we address ourselves. We shall review some of the previous work in this area, outline several improvements in an algorithm proposed several years ago, and present an illustration of the improved technique.

# Previous Work in This Area

In a previous paper before this group, two of the present authors reported preliminary results from a new method of data analysis. We had undertaken to design, program, test, and document a large scale computer program. The algorithm incorporated into the program was aimed at extracting information bearing on the need to introduce interaction terms into a multivariate analysis in which a set of correlated predictor variables were to be related simultaneously to a criterion. Starting from a consideration of some of the problems inherent in applying multivariate statistical techniques (particularly multiple regression) to cross-sectional survey data, we had concluded that the use of Multiple Classification Analysis (MCA) permitted adequate handling of most data characterized by intercorrelated predictors, nominal scales and non-linearities, but could not deal with interaction effects. The computerized procedure that was developed at that time attacked the problem of locating interaction terms by asking a different kind of statistical question of the data than is implied by the immediate choice of MCA.

The technique (termed Automatic Interaction Detection [AID]) employs a sequential decision procedure to divide the sample into a mutually exclusive set of sub-groups through a series of successive binary partitions, each formed by combining observations which are alike on certain characteristics. At each stage every observation is a member of exactly one such subgroup. These groups are formed so that at each stage in the process their means account for more of the total sum of squares (i.e., reduce the error sum of squares more) than the means of any other pair of such sub-groups obtainable by the algorithm. Thus, at each stage

in the branching process, the set of groups developed at that point represents, according to the criteria of the model, the best currently available scheme for predicting the dependent variable from the information available. The branching process terminates when all existing sub-groups are either so homogeneous that the desired accuracy has been achieved, or no variable can be found which will enable an improvement in prediction sufficient to warrant its use.

An examination of the statistics computed during the tree-like partitioning process provides evidence in support of, or against, the introduction of additivity assumptions, as well as providing some indication of what kind of interaction term should be generated, if required.

Experimentation with the method on data which were constructed to violate additivity assumptions revealed that asymmetric tree structures were associated with the existence of predictors which interacted. On the other hand, additivity was associated with symmetric structures.

Interactive data were also found to be associated with the presence of differential changes in explanatory power displayed by corresponding predictors in different branches of the tree. Additivity was associated with similar changes in explanatory power.

Inspection of the profiles of mean values defining the effects of a given predictor in various branches of the tree also provided evidence of the appropriateness of making additivity assumptions. Similarity of effect profiles of a predictor in various parts of the tree structure implies additivity; but interactive data produce incongruent profiles.

AID was first proposed as a substitute for MCA. However, experience with the method in actual use and knowledge gained through experimentation with known data structures led to the development and publication of a strategy for using the AID and MCA techniques

3 This experimentation is reported in Songuist (1969a).

Preliminary results are reported in Morgan and Sonquist (1963b).

Details of the process, flow charts, the computer program, sample input and output, formulas, nine illustrative analyses and recommendations for interpreting the output of the program are reported in Sonquist and Morgan (1964).

jointly, to supplement each other.<sup>1</sup> In addition it appeared that substantial gaines in analytic power might be made from implementing a somewhat more 200phisticated sequential algorithm.<sup>2</sup>

Since its original implementation in a large scale program (Sonquist and Morgan, 1964) the algorithm has been translated by others to run on a number of other computers (Campbell, 1965; Aptakin, 1965; Land, 1965; Kay, 1966; Marks, 1966). Versions of the AID algorithm were also adapted to European equipment (Biervert, 1966, Arpi, 1967). The methods have been used by economists (e.g., Snowbarger, 1967; Gensemer, Lean and Neenan, 1965), by sociologists (Ross and Bang, 1965), by political scientists, (e.g., Sarlvik, 1968), by marketing researchers (e.g., Arpi, 1967), by engineers (e.g., Carlson, 1967), and by psychologists, (Caplan et al, 1966).

## Revisions in the Algorithm

Experience in using the original algorithm over a period of five years, discussions with others who have used it, and experimentation with contrived data have all led to a rather large collection of proposed improvements in the techniques. Most of these are minor from an analytic standpoint, such as obtaining output in improved form. They are obtained simply as a byproduct of the reprogramming necessary to incorporate the major changes proposed below.

The major revisions include extending the original algorithm to deal with covariance as well as a profile of means, providing for the imposition of an optional premium for symmetric partitioning, and extending the sophistication of the search algorithm to include a "look-ahead" to several successive partitions. In the latter case, the mechanism would not proceed blindly, always trying to maximize the explained variation in the current partition, but, like a chess player, would explore the possibility of sacrificing present payoff in favor of even greater gains from subsequent "moves." It would also attempt to provide the simplicity and parsimony of a symmetric model where the data appeared to warrant its use. Each of the revisions is discussed below in some detail.

## Covariance Search Routine

There are many situations in economic, sociological and psychological research where a multivariate analysis is required, but where there exists one dominant explanatory or control variable. The explanatory factor in question may be of particular importance to some theoretical edifice, it may be subject to control more easily than other factors, or it may simply have been measured much more readily and reliably than the others. Where the data come from an experiment and not a survey, the obvious procedure is covariance analysis.

However, with non-orthogonal survey data, one may want to search out subgroups in which there are different relationships between the dependent variable and this dominant explanatory variable. For instance, in much analysis of cross-section survey data, the economist is often interested in the effect of personal or family income on some behavioral variable, and on whether that "income effect" (as represented by the slope of the regression of the behavioral variable on income) varies with other circumstances. The answer to this question will help to determine whether it is necessary to disaggregate the data in models used for forecasting, and the optimal way to do it.

Sociologists and psychologists often face similar problems in which the purpose of the investigation requires isolating the effect of a particular variable under a wide variety of combinations of circumstances. For instance, intelligence, alienation and authoritarianism have each been the subject of repeated investigations in which the object has been to relate that particular factor to specific consequences in such a way as to specify the form of the relationship under various conditions and for particular types of people.

Another illustration is in the analysis of changes taking place over time.

<sup>&</sup>lt;sup>1</sup> The strategy developed was suggested in Andrews, Morgan and Sonquist (1967), and then elaborated further by Sonquist (1969a, 1969b).

<sup>&</sup>lt;sup>2</sup> See Sonquist (1967).

<sup>&</sup>lt;sup>3</sup> Some of the programming revisions might include the use of Fortran IV for compatibility with various manufacturers' equipment; a trichotomous partition option; improved input flexibility with regard to both data and control language; additional output including predictor summary tables; and improved variable transformation capabilities.

The initial value of a phenomenon under study clearly affects its value measured at a subsequent time. This is why the residuals from the regression of its current  $(t_2)$  value on its initial  $(t_1)$ value are often used as a measure of change, instead of the raw increments. However, this "initial value" effect might not be the same for all subgroups in the population. If, then, a single equation were to be fitted, a downward bias would be exerted on the correlations between change and those factors thought to be responsible for it. Thus, when residualizing a variable for study, a search should be made to determine if this initial value effect is homogeneous throughout the population.

For all these reasons, a variant of the earlier sequential data analysis algorithm has been developed in which the criterion for sequential subdivision of a sample is changed from the sum of squares explained by two subgroup means (instead of one pooled mean) to the sum of squares explained by two simple sub-group regressions (instead of one simple regression of the pooled data). The search algorithm otherwise has basically the same framework: each group potentially to be divided is examined using all feasible partitions based on each explanatory characteristic. The difference is that the quantity maximized when a group is divided into sub-groups is the variation explained by the regression of the dependent variable on the covariate within each sub-group. If the two best fitting regression lines differ as to intercept or slope, then the unexplained variation would be reduced, and the split with the largest difference between the two regressions would reduce it the most.

Given this conceptualization of the problem, a partition may be chosen so as to maximize any one of three quantities. One may evaluate reductions in unexplained variation due to differences between means, differences in regression lines, or both taken together. The original AID algorithm sought to maximize the sum of squares explained by means of the sub-groups resulting from the partition of the "parent" group, i.e., it maximized the expression

$$N_1 \overline{Y}_1^2 + N_2 \overline{Y}_2^2 - N \overline{Y}^2 \qquad (2)$$

The rationale for this is easily seen in Table 1.

It can be seen from Table 2 that a second quantity could also be maximized. This is the expression

$$(N_1-1)r_1^2s_{y1}^2 + (N_2-1)r_2^2s_{y2}^2 - (N-1)r_2^2s_y^2$$
  
(3)

This is the sum of squares of the two regression estimates around the two group means resultinf from the partition.<sup>2</sup> It reduces to:

$$\sum_{\substack{\Sigma \\ i=1}}^{2} \left[ \frac{\left[\Sigma\left(\overline{Y}-\overline{\overline{Y}}_{i}\right) (\overline{X}-\overline{\overline{X}}_{i}\right)\right]^{2}}{\Sigma\left(\overline{X}-\overline{\overline{X}}_{i}\right)^{2}} \right] - \frac{\left[\Sigma\left(\overline{Y}-\overline{\overline{Y}}\right) (\overline{X}-\overline{\overline{X}})\right]^{2}}{\Sigma\left(\overline{X}-\overline{\overline{X}}\right)^{2}}$$
(4)

A third quantity which may be maximized is obtained from the sum of these two and would represent the effect of both the variable used in the partition and the X covariate. The variation explained by the sub-group means and that explained by the regression are both maximized. This expression is

$$N_{1}\bar{\bar{Y}}_{1}^{2} + (N_{1}-1)r_{1}^{2}s_{y1}^{2} + N_{2}\bar{\bar{Y}}_{2}^{2} + (N_{2}-1)r_{2}^{2}s_{y2}^{2} - N\bar{\bar{Y}}^{2} - (N-1)r_{sy}^{2}$$
(5)

The existing algorithm has also been modified to permit an examination of the effects of one crucial categorical predictor in various parts of the sample without permitting it to be used in the partitioning process. Permitting the analyst to retain the ability to conceptualize the effects of the crucial variable in terms of slopes and intercepts as well as the ability to use profiles of sub-group means gives a desirable simplicity.

For a thorough discussion of this problem see Lord (1967).

<sup>&</sup>lt;sup>2</sup> If the order of the (k) classes of the explanatory characteristic is maintained, there are only (k-1) ways of forming two groups on the basis of that predictor. However, if one reorders them on the basis of the means of the dependent variable there are many more possible ways of forming the two sub-groups.

Computational formulas are, of course, somewhat different and are not given here. In maximizing any of these expressions the last term is constant over all possible partitions and can be ignored. See Walker and Law (1953) pp 210-216.

<sup>&</sup>lt;sup>2</sup> See Walker and Law (1953), pp 242-244.

Source of Variation	d. f.	Sum of Squares	Mean Square
Observations around grand mean	N-1	$\sum_{j=1}^{k} \sum_{i=1}^{N_{j}} (Y_{ij} - \bar{Y})^{2} = SST$	s <sup>2</sup> y
Between group means	k-1	$\sum_{j=1}^{k} N_{j} (\bar{\bar{Y}}_{j} - \bar{\bar{Y}})^{2} = SSB$	MSB
Within Groups	N-k	$\sum_{j=1}^{k} \sum_{i=1}^{N_{j}} (Y_{ij} - \overline{Y}_{j})^{2} = SSW$	MSW

Analysis of Variance for Differences in Means

The covariate problem is illustrated below. Owning a home (as opposed to renting, etc.) is related not only to income, but to age, family size, and place of residence (urbanization). More important, a very large fraction of older people outside the large urban areas own a home regardless of their income; i.e., income differences have no effect at all on home ownership in this group. However, among young families with children, small increases in income appear to lead to substantial in creases in the probability that the family will soon own its home. Furthermore, among young people with no children at all, home ownership is rare at any income level. Clearly, economic projections of home ownership need to be based not just on aggregate statistics of income increases, but on who received them.

The effective use of survey data to shed light on such problems requires concentrating attention on the differential relation of income to economic activity as well as studying its effect in the aggregate.

These details are further illustrated by an examination of Figure 1. In the total sample the regression of Y on  $X_1$  is

 $Y = a_1 + b_1 X_1 + u_1$ (6)

However, when the sample is split on

variable X, into 2 groups, one with  $X_2 = 1,2$  and the other  $X_2 = 3,4$ , for the latter group this regression is

$$Y = a_3 + b_3 X_1 + u_3$$
(7)

and if, in addition, this group is split on variable  $X_4$ , then for  $X_4 = 4$  or 5, the regression of Y on  $X_1$  is

$$Y = a_7 + b_7 X_1 + u_7$$
 (8)

This is illustrated in Figure 1.

If we also have

and 
$$\Sigma u_1^2 > (\Sigma u_2^2 + \Sigma u_3^2)$$
 (9)  
 $\Sigma u_2^2 > (\Sigma u_4^2 + \Sigma u_5^2)$  (10)  
 $\Sigma u_3^2 > (\Sigma u_5^2 + \Sigma u_7^2)$  (11)

it is clear that the various effects of  $X_1$ on Y as revealed by the differences in the slopes (b<sub>1</sub>) and the intercepts (a<sub>1</sub>) are associated with the joint occurrence of the conditions denoted by the indicated values of variables  $X_2$ ,  $X_3$  and  $X_4$ . Thus one must devise a means of searching various sample sub-groups in order to learn whether these differential effects exist and what forms they take, under adequate constraints to reduce the probability of detecting differences which are spurious. It is this problem which lends itself to solution via a formal sequential decision process programmed to run on a

# TABLE 2.

Source of Variation	D. F.	Sum of Squares	Mean Square	
Regression estimates around $\overline{Y}$	1	$\sum_{i=1}^{N} (\hat{\bar{Y}}_{i} - \bar{\bar{Y}})^{2}$	(N-1)r <sup>2</sup> sy2	
Observations around regression estimates	N-2	$\sum_{i=1}^{N} (\mathbf{Y}_{i} - \mathbf{\hat{Y}}_{i})^{2}$	$\frac{N-1}{N-2} (1-r^2)s_y^2$	
Observations around $\overline{Y}$	N-1	$\sum_{i=1}^{N} (Y_i - \overline{Y})^2$	sy2	
Where $r^2 = \frac{\sum (\bar{x} - \bar{x}) (\bar{y})}{\sum (\bar{x} - \bar{x})^2 \Sigma}$	$\frac{\left[-\overline{Y}\right]^2}{\left(\overline{Y}-\overline{Y}\right)^2}$	$\hat{Y} = \bar{Y} + b_{yx}(X-\bar{X})$		
$s_y^2 = \frac{\Sigma (Y-\overline{Y})^2}{N-1}$		$\mathbf{b}_{\mathbf{y}\mathbf{x}} = \frac{\Sigma (\mathbf{x} - \overline{\mathbf{x}}) (\mathbf{y} - \overline{\mathbf{y}})}{\Sigma (\mathbf{x} - \overline{\mathbf{x}})^2}$		

# Analysis of Variance for Regression





Figure 1. Differential Effects of  $X_1$  on  $Y_1$ .

# TABLE 3.

					AB	
		_	_			
Configuration	AB	AB	AB	Y=0	Y=1	Y=2
Column(j)	1	2	3	4	5	6
Row(1)						
1		0	0	مa	CS	<u></u>
2	1	0 0	ň	CS .	BICOND	(M) BTCOND
2	2	0	ő	CS	(M) BICOND	BICOND
5	ő	1	0	CS CS		(M)CS
7	1	1	ő	~	SC A	(H)C5
5	2	1	0 0		ЪС В	(M) BT COND
7	2	2	0	(1)05		
8	1	2	0	(M) CS	(1)05	
9	2	2	0	(1)05		(H) 50
10	2		1	C S		
10	1	Ň		~	A SC	(H)C3
12	2	0	1		ыс Б	
13	2	1	1	FYOR	R C	(M) BICOND
15	1	1	1	SC	30	A CS
14	1 2			30	A CS	PICOND
15	<sup>2</sup>	2	1	A (M) EVOD	C5 P	(M)SC
10	1	2		D		(11)30
19	2	2		M)SC		A SC
10	2	2	2	(M) SC		
20	1	0	2		(11)05	M)SC
20	1 2	0	2	(1)05		(M)SC
21		1	2	M EVOR	(H)30 p	(M) SC
22	1		2	D		(1)50
23	2		2	K (M)SC		A SC
24	2	1 2	2	(FI) SC	A (M) EVOD	80
25	1	2	2	EAUK (M) EVOD	EXOP	80
20	2	2	2	C	SC	30 A
21	4	2	2	30	30	A

# Logical Models for Two Dichotomous Predictors and a Trichotomous Dependent Variable

<sup>a</sup> A means additive; CS means cumulative upwardly and substitutive downwardly; SC is the inverse of CS; (M) means "modified;" BICOND means "biconditional;" R means reversal and EXOR refers to "exclusive or". Table adapted from Sonquist (1969a).

#### large scale computer.

# <u>Extending the Search Algorithm</u> - <u>A</u> <u>Multi-step Look-ahead</u>

Experimentation with the original algorithm to determine its behavior under known conditions was carried out using contrived data. This process of working out tests under a variety of different conditions led to the development of a typology of multivariate models. From this viewpoint, multivariate interactive and additive AID and MCA models could be viewed as eighty-one variants of the same basic structure, which, in its simplest form (two dichotomous predictors and a trichotomous dependent variable) could be defined almost entirely in terms of the fundamental operators of Bool-ean algebra, (see Table 3), Eliminating permutations and inverses reduces the number of models to seven non-trivial ones, (see Table 5)<sup>2</sup>.

This typology also delineated exactly certain limitations of the original AID algorithm, some of which could not have been noticed earlier. For instance, while the published procedure was found to be capable of dealing adequately with many two-way interactions, others were identified as being difficult for it to deal with (Sonquist and Morgan, 1964). These were seen to consist of interactive relations characterized by consistency, i.e., by balance or symmetry. One such example is the biconditional model, illustrated in Table 5.

#### Table 4.

Occurrence of Effect "C" Associated with the Joint Occurrence or Joint Absence of Two Causal Factors, A and B.

ndition "A"	В	Condition	"B" Not-B
A	С		Not-C
t-A N	lot-C		С
t-A N	lot-C		

A more extensive discussion of these models is given in Sonquist (1969a)

<sup>2</sup> The eighth is a constant in all cells.

Although the frequency with which variants of this model actually occur in real data is not known, it is notable that at least one realization has received considerable attention in the recent sociological literature, the concept of status inconsistency. The problem of developing analysis techniques to deal with this class of models is of importance for economic, educational and psychological research as well as for sociology.

It can be seen that the earlier sequential partitioning algorithm which examines only the "zero-order" effects of A and B separately could not discover the consistency effect in these data. There are really two A "effects" and they cancel each other out in the total group. Moreover, the additive assumptions implied by the choice of Multiple Classification Analysis would also tend to conceal the real state of the world.

However, the extended AID algorithm partitions the sample tentatively, first on one causal variable and then on the other (as well as making tentative partitions on other variables). The actual partition is made so as to maximize the effects of several successive partitions. This makes it possible first to reveal a consistency effect to the analyst by means of appropriate output, then to make an appropriate partition, and finally, to continue with the rest of the sequential search procedure.

In general, such a two-split scanning algorithm appears capable of providing information adequate for the analyst to identify the two-way interactions existent in the data. It also appears able to provide leads or clues to the existence of three-way interactions. This is seen to be a simple extension of the way in which the present algorithm provides clues to the existence of twoway interactions. Thus, an algorithm which examines the cross-classifica-

For an example, see Blalock (1966).

# TABLE 5

Seven Logical Models for Two Dichotomous Predictors and a Trichotomous Dependent Variable



4. Additive

Н	М
м	Ĺ

5. Biconditional

Н	L
L	н

6. Modified

DICON	ditional
Н	М
L	Н

7. Reversal

H	М
L	м

tion of p predictors simultaneously appears able to reveal terms composed of p raw variables regardless of the symmetry of the term. However, such an algorithm also appears capable of revealing a term involving p + 1 raw variables if the term is asymmetric.

For instance, if we have the three variable asymmetric model,"if A and B and C, then Y = 0, otherwise Y = 4," the algorithm using a twosplit strategy would produce the sequence of partitions illustrated in Figure 2.

Of course the amount of computing required to search out combinations of three or more variables increases as an exponential function of the number of variables considered simultaneously. Hence constraints have to be put on the process to permit the eliminations of unpromising leads and thus the examination of the subsequent partitions. However, this does not appear necessary in the three variable case.

# Extending the Search Algorithm -Premium for Symmetry

The original algorithm represented a step in the direction of specifying a statistical model so it fits the data rather closely. The introduction of the look-ahead principle moves further in this direction. However, it can be anticipated that increasing the "wiggling" ability of a model being fitted will also increase the probability of the analyst basing his theoretical model on data largely reflecting idiosyncratic characteristics of the sample under investigation. Hence additional constraints are needed that would tend to guard the analyst against over-fitting his model. One such constraint is a premium for symmetry. Thus, when a look-ahead is employed a capability can be provided for increasing the probability that if a given predictor is used to make a partition in a given way on one branch of a partition sequence it will also tend to be used similarly in the parallel branch. This principle of constraint toward symmetry is illustrated in Figure 3.

Once groups four and five have been created using variable B the symmetry question arises. The proposed partition of group two, the "Not-A's" into new groups could be accomplished using, say, variables B, C, or D, but not variables E, F, or G, since the latter show insufficient explanatory power. In the previous algorithm, the choice of a variable on which to base a partition would have been to compare B, C, and D and then choose the one capable of producing the greatest reduction in the unexplained sum of squares. The present proposal would alter the algorithm to require that if C or D were chosen over B, it would have to achieve a certain ratio of explanatory power when compared with that resulting from a partition based on B identical to the one already performed in the paral-lel group. The comparison ratio would be supplied by the analyst at the time of execution of the program. Setting the ratio to 1.0 would simply cause the regular algorithm to take effect. Setting it larger than 1.0 would bias the algorithm toward symmetry; setting it at less than 1.0 would tend to prevent symmetry. For instance, a value of 1.25 would require that a non-symmetric partition explain 25 percent more variation than a symmetric partition in 9rder to be actually used in a split.

# Preliminary Examination of the New Algorithm

As a preliminary investigation into the power of the new algorithm with respect to the look-ahead and covariance options, tests were made using all combinations of the seven logical models of Table 5 applied to both means and slopes (see Tables 6a and 6b).

- In the case of dichotomous predictors the identical partition is the only possible one. This is not the case where the variable has three or more classes. We focus on total symmetry, i.e., forming identical sub-groups based on the same predictor.
- <sup>2</sup> Appropriate values of the symmetry premium for actual use are yet to be worked out by experimentation.
- <sup>3</sup> For each of the means model in Table 5a, a dependent variable was generated, Y<sub>1</sub> = a + e<sub>1</sub>, i = 1,...,100, where (a) represents the mean of a given cell (25 observations were generated for each call) and e<sub>1</sub> is a random error e<sub>1</sub> ~ (N(0,.5). The e<sub>1</sub> remains constant for all 7 models. A covariate x<sub>1</sub> ~ N(0,1), i = 1, ...,100, was also generated, and for all 49 combinations of means and slopes the dependent variable y<sub>1</sub> = a + bx<sub>1</sub> + e<sub>1</sub>, i = 1, ...,100, was formed, where (a) and (b) are the mean and slope values for the appropriate cell. The meansalone model, y = a + e<sub>1</sub> was also evaluated, making the total number of experiments fifty-six.



Figure 2. ABC Implies Y = 0, Else Y = 4



Figure 3. Symmetric Partition.





7 Logical Models Applied to Means



7 Logical Models Applied to Slopes of Y on X



Since this was a preliminary study, "masking" factors such as intercorrelations and noise in the predictors were not considered. However, to provide a test of the algorithm's ability to pick out explanatory factors from a noisy background two uncorrelated dichotomous dummy variables were generated along with the two variables defining the model. Each experiment was made using the two real factors and the two noise factors as predictors, utilizing the "look-ahead" with two splits (creation of three groups).

For the means-alone cases, the sum of squares between the three terminal groups was maximized, i.e., from Table 1,

$$\sum_{i=1}^{3} N_{i} \overline{Y}_{i}^{2} - N\overline{Y}^{2}$$
(12)

For the means and slopes cases, the quantity maximized between the three goups was the sum of squares resulting from both the reduction in means and regression, i.e., from Table 2,

$$\sum_{i=1}^{3} \left[ N_{i} \overline{Y}_{i}^{2} + (N_{i} - 1) r_{i}^{2} s_{yi}^{2} \right]$$

$$- \left[ N \overline{Y}^{2} + (N - 1) r_{y}^{2} s_{y}^{2} \right]$$
(13)

As was expected, with the exception of the biconditional model, the look-ahead with two splits yielded results similar to those from no lookahead; the final groups in both cases were identical, but occasionally the look-ahead would lead to splits on the variables in reverse order from a parallel analysis using no look-ahead.

The biconditional-means model still presented some problems, however. For the no-slope version of the model as well as the univariate, cumulative, modified cumulative, additive and reverse slope models, the algorithm with no look-ahead could make incorrect splits; that is, it could use the dummy variables by mistake. it is significant, however, that the look-ahead identified the models correctly. However, it occasionally made subsequent partitions on the dummy variables. This proved to be capable of remedy by an adjustment of the reducibility criteria controlling the termination of the partitioning process. Proper choice of this criterion still permits legitimate partitions to take place, but prevents subsequent spurious ones. For the no-lookahead case, this fraction p apparently should be in the range .005  $\leq$  P  $\leq$  .008. For the onestep look-ahead (two splits) this fraction P apparently has a lower bound of P1  $\leq$  .016. The look-ahead apparently does the job it was designed for.

Our findings from this initial experimentation with covariance models suggest that the differences in means may be much more powerful in determining what split is to be made than are differences in slopes. In fact, in most of the cases where (a) the two groups had no differences in means and (b) one of these groups had slope zero, the two regressions resulting from tentative partitions were not sufficient to meet the reducibility criterion with, or without, a look-ahead. For instance, in the example given below in Table 7, the algorithm would split the sample on variable A, but would not split either of the resulting groups.

Table 7.

An	Exampl	.е	of	а	Remaining	Problem
----	--------	----	----	---	-----------	---------



1 The criterion is the fraction of the total sum of squares from the total input sample that a split (or sequence of splits in the look-ahead case) must explain in order for the split actually to be made. The maximized function [equation (12) or (13)] must be greater than p times TSS. The behavior of the algorithm is in keeping with previous results from the original algorithm. For p < .005, there is a tendency to split on the dummy variables after the correct splits have been made. For p<sub>i</sub> < .016 the look-ahead would occasionally split on a dummy variable even before making the correct splits.

<sup>&</sup>lt;sup>1</sup> Further work is to include tests with predictors of varying levels of intercorrelation and skewness as well as the assessment of the algorithm activity to deal with correlated "noise."

While this result is not entirely unexpected, it may imply that maximizing differences in means alone or in regression slopes alone may be a more powerful tool than the combination of the two. Further experimentation with these models using only the slopes is obviously indicated.

# Significance of These Findings

This extended algorithm represents a continuation of our attempts to develop better methods for adequate handling of the problem inherent in applying multivariate statistical techniques in the analysis of cross-sectional survey data with large numbers of cases (1000 or more). The covariance capability is relevant for the analysis of panel as well as one-time cross-sectional data. The increasing volume of survey and other non-

## References

- Andrews, F.M., Morgan, J.N., and Sonquist, J.A. <u>Multiple Classification Analysis</u>. Ann Arbor: University of Michigan, Institute for Social Research. 1967.
- Aptakin, Peter. <u>Automatic Interaction Detector</u>. New York: Service Bureau Corporation, Computing Sciences Division. 1965.
- Arpi, Bo. En Modell for Optimal Marknadssegmentering. Lund Business Studies No. 4. Studentlitteratur Lund. 1967.
- Biervert B., Dierkes, M., and Walzel, A. <u>Automatischer Split zur Optimalen Kombination</u> <u>Erklarender Faktoren</u>. Forschunsstelle fur Empirische Sozialökonomik und Rechenzentium, University of Cologne, March, 1968.
- Blalock, H.M. The identification problem and theory building: The case of status inconsistency. <u>Amer. Soc. Rev.</u>, 31, 1, 52-61.
- Campbell, Robert H. <u>Modifications to the Auto-</u> <u>matic Interaction Detector COMCOM/Simulation</u> Memo No. 26. Program operational on the MIT-FMS system. Cambridge: Massachusetts Institute of Technology. 1965.
- Caplan, N., Lippitt, R., Gold, M., Mattick, H., Suttles, G., and Deshaies, D. <u>The Outcome</u> of the Chicago Youth Development Project. New York: Boys Club of America Convention. 1966.

experimental data and the growing needs of researchers in many disciplines to utilize it in solving practical problems makes the development of adequate data-handling methods even more imperative. The procedures are frankly oriented toward the inductive phases of research in which model-building is the objective rather than toward the deductive model-testing phase. It is felt that a focus here will be of benefit both to theory builders and to those who must find immediate solutions to pressing problems.

- It is possible that this behavior results largely from the sizes of the mean and slope effects relative to the noise level. This remains a problem for further exploration.
- Carlson, W.L. <u>Identification of the Problem</u> <u>Driver from Driver Records: A Preliminary</u> <u>Analysis</u>. Ann Arbor: Highway Safety Research Institute, University of Michigan. 1968.
- Gensemer, Bruce L., Lean, Jane A., and Neena, William B. Awareness of marginal income tax rates among high-income taxpayers. <u>National Tax Journal</u>, XVIII, September, 1965, 258-267.
- Kay, Kevin, <u>Automatic Interaction Detector:</u> <u>"AID" Translater into CDC 3600 Fortran and</u> <u>Compass,</u> Technical Report 46. East Lansing: Michigan State University, Computer Institute for Social Science Research. 1966.
- Land, K.C. <u>The Sociology Program Library, A</u> <u>User's Manual</u>. Austin: Population Research Center, Department of Sociology, University of Texas. 1968.
- Lord, F.M. Elementary models for measuring change. In <u>Problems in Measuring Change</u>, C.W. Harris (ed.) Madison: University of Wisconsin Press, pp. 21-39. 1962.
- Marks, E. NCHS AID Program. Unpublished manuscript. Wharton School, University of Pennsylvania. 1965.

- Morgan, J.N., and Sonquist, J.A. Problems in the analysis of survey data: And a Proposal. <u>Journ. Amer. Stat. Assoc</u>., 1963(a), 58, 415-434.
- Morgan, J.N., and Sonquist, J.A. Some results from a non-symmetrical branching process that looks for interaction effects. <u>Proceedings</u> of the Social Statistics Section, American Statistical Association, 1963(b).
- Ross, J., and Bang, S. Predicting the adoption of family planning. <u>Studies in Family</u> <u>Planning</u>, No. 9. 1966.
- Sarlvik, Bo. Socio-Economic Predictors of Voting Behavior: Research Notes from a Study of Political Behavior in Sweden. Unpublished Manuscript. 1968.
- Snowbarger, Marvin. <u>An Interaction Analysis of</u> <u>Consumer Durable Expenditures.</u> Unpublished <u>Ph. D. thesis, University of Michigan.</u> 1966.

- Sonquist, J.A. <u>Automatic Interaction Detection</u> and <u>Multiple Classificcation Analysis: The</u> <u>Validation of a Search Strategy.</u> Unpublished Ph. D. dissertation, University of Chicago. 1969(a).
- Sonquist, J.A. Finding variables that work. <u>Public Opinion Quarterly</u>, 1969(b), Vol. XXXIII, No. 1, pp. 83-96.
- Sonquist, J.A. Simulating the research analyst. Social Science Information, 1967, Vol. 6, No. 4, p. 207-215.
- Sonquist, J.A., and Morgan, J.N. <u>The Detection</u> of Interaction Effects. Ann Arbor: Survey Research Center Monograph No. 35, Institute for Social Research, University of Michigan. 1964.
- Walker, H.M., and Lev, J. <u>Statistical Infer-</u> <u>ence</u>. New York: Holt, Rinehart and Winston. 1953.