

MULTIPLE LINEAR REGRESSION VIEWPOINTS

A publication of the Special Interest Group on Multiple Linear Regression

MLRV Abstracts appear in CIJE, the ERIC System, and microform copies are available from University Microfilms International MLRV is listed in EBSCO Librarians Handbook.

# ISSN 0195-7171

## MULTIPLE LINEAR REGRESSION VIEWPOINTS

### Chaliperson

Editor Isadore Newman The University of Akron

Assistant/Editor Keith McNell New Mexico State University

Las Cruces, NM 88001 Cynthia Boone-Hawkins The University of Akron Akron, OH 44325-4205

Akron, OH 44325-4205

### Executive Secretary

Cover Artist

## EDITORIAL BOARD

### Walter Wengel Computer Systems Manager 2123 S. Arlington Heights Rd. Arlington Heights, IL 60005

Susan Tracz Department of Advanced Studies California State University Freeno, CA 93740

Samuel Houston Department of Mathematics and Applied Statistics University of Northern Colorado Greeley, CO 80639

Dennis Hinkle Virginia Polytechnic Institute Blacksburg, VA 24061 Andrew Bush Baptist Momonal Hospital Memonals TN

John William University of North Dakota Grandi Forkey ND 68201

Joe Ward San Antonio, TX 78228

Basil Hamilton North Texas State University Denton, TX 76201

Isadore Newman Research and Evaluation The University of Akron Akron, OH 44325-4205 Comparison of Conjoint Analysis, Multiple Regression Models with Person Vectors and Profile Analysis to Assess Important Factors Used to Select Colleges

> isadore Newman The University of Akron

John Frees Ashland University

The purpose of this study was to investigate the relative effectiveness of the traditional conjoint analysis approach to the multiple regression approach that includes person vectors profiles analysis. It was expected that the more sophisticated models would increase the effectiveness in terms of its shrinkage estimates and the accuracy of its predictability of two holdout groups. The data source consisted of a sample of 100 students who rated eight colleges on five attributes--quality of education, financial aid, quality of dorm life, student/faculty relations, and social aid.

PRESENTED AT THE AMERICAN EDUCATIONAL RESEARCH ASSOCIAITON 1990 ANNUAL MEETING, BOSTON, MASSACHUSETTS. Comparison of Conjoint Analysis, Multiple Regression Models with Person Vectors and Profile Analysis to Assess Important Factors Used to Select Colleges

#### Introduction

In recent years, many colleges and universities have faced increased competition for students. Thus, it has been increasingly important for an institution of higher education to be able to identify what factors are important to the students who chose to enroll in the institution.

Marketing research (Cattin & Wittink, 1982) has identified conjoint analysis as a very useful statistical technique in which one is interested in having the clients, students, or consumers prioritize a variety of items. Two other approaches also seem to be appropriate to use when attempting to assess the selection process of college-bound students: (1) multiple regression models with person vectors (Frass & Newman, 1989); and (2) profile analysig.

#### Objectives

.

This paper attempted to compare the ability of conjoint analysis, multiple regression models with person vectors, and profile analysis to produce information that could be used by college and university personnel to determine which factors were important to students when selecting a university or what type of students selected a given type of university.

#### Data Collection

The research instrument used to collect the data analyzed in this study focused on five institutional attributes reported to be of significance to students who martriculated to Ashland University. This list of attributes was developed through literature reviews (Tiernry 1980; Traynor, 1981; Kuh, Coomers, & Lindquist, 1984; Conant, Brow, & Mokwa, 1985), discussion with program advisors and students, and from the past experiences of admissions recruiters.

The five attributes included in this study were financial aid, social life, quality of dorm life, student-faculty relationships, and quality of education. Each of the five attributes had two levels. The two levels that were formed for each attribute were assigned a value of 0 or 1 in order to allow the researchers to quantitatively form hypothetical universities with various combinations of attribute levels. The attributes, levels, and values assigned to each level were as follows:

1.	Quality of education
	a) reputation is not well known = 0 b) reputation is well known = 1
2.	Student/Faculty relationships
	<ul> <li>a) faculty are accessible if sought = (</li> <li>b) faculty are extremely accessible = 1</li> </ul>
3.	Quality of dorm life
	a) below my expectations = 0 b) above my expectatons = 1

#### 4. Financial aid

Star Bart

**. 5.** 

a) little financial need is met = 0 b) most financial need is met = 1

 $\sum_{i=1}^{n} \frac{1}{i} \sum_{i=1}^{n} \frac{1}{i} \sum_{i$ 

Social Life

a) few social activities are available = 0
 b) many social activities are available = 1

Five attributes with two levels each would allow 32 different university profiles to be formed. With the assumption that interaction effects are negligible, the main effects could be estimated with only eight orthogonal arrays. The eight orthogonal arrays used in this study which were formed with the aid of the computer software entitled Conjoint Designer (Bretton-Clark, 1987), were listed in Table 1.

In addition to the eight orthogonal arrays, two arrays were designed to provide a means of assessing the degree of predictive validity. (See Table 1.) These two arrays were referred to as the "holdout universities" because they were not included in the estimation procedures.

The questionnaire was administered during the second week of the fall term of 1987 to freshman students enrolled in a freshman seminar course. The responses of 100 of the students were used in this study. See Fraas and Paugh (1989) for

#### Conjoint Analysis

The analysis conducted by the use of a software package (Bretton-Clark, 1987) produces a set of five regression coefficients plus a constant term for each student. That is, a separate regression analysis was performed on the data of each of the 100 students.

Each of the regression coefficients generated by the conjoint analysis for a given student indicated what would happen to the respondent's ratings of the universities when the attribute changed from the "zero" level to the "one" level. To illustrate the point, consider the regression coefficient value of 2.0 recorded for the financial attribute for respondent 1. If financial aid was to increase from the "little need being met" category to the "most need being met" category, the respondent's ratings of the universities would increase by 2.0 points on the 1 to 10 scale used on the questionnaire.

A relative importance figure was calculated for each attribute by dividing the sum of the five average regression coefficients into each of the average regression values. The five relative importance figures generated by this procedure were expressed as percentages.

### Table 1

**A** • Multiple Linear Regression Models

Universities	Quality of Education	Student/ Faculty Relation- ships	Quality of Dorm Life	Financial Aid	Social Life
	· · 0		0	0	0
e entre Aley <b>B</b>	1	0	0	1	1
	1	- <b>1</b>	1	1	0
1997 - 19	0		· <b>1</b>		1
Barra and Charles and B	0	1	0	1	0
14 - N. 19 - N. 19 - N. 19 - 19 - 19 - 19 - 19 - 19 - 19 - 19	1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	u servicu i ser	а., « О	1
an a	1	0	1	т. С. С	0
n. H	0	0	· · · <b>1</b>	1	1
		an an Araba Araba Araba	1) 第1 以前 1) · · · · · · · · · · · · · · · · · · ·		
Holdout			e de la companya de	5	
Universities I	1	<b>1</b>	1. <b>1</b>	0	1
	1	<b>1</b> .	0	1	0

Note. Each characteristic is composed of two levels. The zero value indicates the presence of the lower of the two levels.

6

and the second second

Results of the Conjoint Analysis

The relative importance figures indicated that financial aid was the most important attribute with a value of 26.24%. Financial aid was followed in importance by the quality of dorm life (21.29%), the quality of education (20.84%), the student/faculty relationships (16.63%), and the social life (15%). (See Table 2.)

: '

7

#### Predictive Validity

The observed and predicted ratings for the holdout universities were used to provide two estimates of the ability of the results of the conjoint analysis to predict student ratings. The first estimate was a correlation coefficient for the predicted and observed ratings. The second estimate was an average absolute difference value for the difference between the predicted and observed ratings. The correlation coefficient value and the average absolute difference for the observed and predicted ratings were .37 and 1.87, respectively.

#### Multiple Linear Regression Model

## With a Surrogate Person Variable

#### Model Structure

The second approach used to analyze the survey information required the construction of a multiple linear regression model that included a surrogate person variable. Before such a model is presented, however, a discussion of a model that includes the actual person variables may prove helpful. The variables included in the model that used person variables (Model 1) were as follows:

X = ratings of the eight hypothetical universities (values ranged from 1 to 10) X1 = quality of education

0 = "low" level; 1 = "high" level)
X2 = student/faculty relationship

(0 - "low" level; 1 - "high" level) X3 - quality of dorm life

(0 = "low" level; 1 = "high" level)

X4 = financial aid

(0 = "low" level; 1 = "high" level)

X5 =social life

(0 = "low" level; 1 = "high" level)

P1 = respondent 1

(1 if from respondent 1; 0 otherwise)

P2 = respondent 2

(1 if from respondent 2; 0 otherwise)

P99= respondent 99

(1 if from respondent 99; 0 otherwise)

The structure of the regression model with person variables was:

Y = aU + b1X1 = b2X2 = b3X3 = bb4X4 = b5X5 = b6P1 = b7P2 = . . .b104P99 = e (model 1)

The use of the person variable required by Model 1 is not practical due to their large number. Thus a multiple linear regression model designed to include a surrogate person variable was used. This surrogate person variable measured the impact of the 99 person variables required by Model 1.<sup>1</sup>

#### Table 2

### Conjoint Analysis Results

e.,		÷.	15		1	
11	λ.	1.1	1	5.7		

	Average	· • • • • •
Characteristic	Regression Coefficient	<pre>% of Relative Importance</pre>
Ripancial Aid	1 776	26.24
FINANCIAL MIG	<b>1.</b> //5	20.24
Quality of Dorm Life	1.440	21.29
Quality of Education	1.410	20.84
Student/Faculty Relationships	1.125	16.63
Social Life	1.015	15.00
n an	n n n n n n n n n n n n n n n n n n n	

Correlation coefficient between the predicted and observed ratings of the holdout universities = .37

Average absolute difference between the predicted and observed ratings of the holdout universities = 1.87

The value of the surrogate person variables was composed of an average rating for each person. The surrogate variables was represented in Model 2 by "X6." The values for this variable ranged from 2.625 to 8.5 for the 100 students.

The multiple regression model with the surrogate person variable (Model 2) used to analyze the survey information was as follows:

Y = aU = b1X2 = b2X2 = b3X3 = b4X4 = b5X5 = b6X6 = e (Model 2)

The regression coefficients for the university attributes that were generated by Model 2 were equal to the average regression coefficients for the conjoint analysis (See Table 3).

Before the regression coefficients could be statistically tested, the standard errors had to be corrected for the appropriate degrees of freedom. The number of degrees of freedom was 695, which was equal to the sample size of 100 (number of students) multiplied by 8 (number of colleges) minus 6 (number of attributes plus one). Each of the regression coefficients for the university attributes was statistically significant at the .01 level. The multiple correlation coefficient was .764; and 2 value was .58. the R

#### Predictive Validity

The regression coefficients generated by Model 2 were used to predict the ratings of the holdout universities. The

correlation coefficient for the predicted and observed ratings was .76. The average absolute difference between the predicted The same procedure applied to the second half of the data set resulted in a correlation coefficient value of .74 between the observed and predicted ratings. Again, this value shows little shrinkage (1.7%) from the multiple correlation coefficient of 753 for Model 2.

••

1 Refer to Pedhazur (1977), Williams (1977; 1980), Fraas

and the second sec

a da serie Esta da serie da serie

McDougall (1983), and Williams and Williams (1985a; 1985b) for discussions of a surrogate variable used to measure the amount of variation in the dependent variable associated with a set of person variables.

#### Comparison of the Results

The estimated impact of the university attributes on the student ratings by the conjoint analysis, and the multiple linear regression model with a surrogate person variable were identical. For both procedures, the order of importance was as follows: (1) financial aid, (2) quality of dorm life, (3) quality of education, (4) student/faculty relationships, and (5) quality of social life.

The multiple linear regression model with the surrogate person variable, however, produced a correlation coefficient value of .76 for the predicted and observed ratings of the holdout universities, as compared to the value of only .37 for the conjoint analysis.

The multiple linear regression model with the surrogate person variable also produced a lower average absolute difference between the predicted and observed ratings for the holdout universities than did the conjoint analysis. The average absolute difference values were 1.50 and 1.87.

The low  $\mathbb{R}^2$  values of the regression models that used the clusters as the independent variables indicated that the clusters were unable to explain the variation in the university ratings to any high degree. For this data set, the cluster information was of little assistance in identifying the importance of university characteristics as viewed by various groups of students.

## Quannal Analysis

The following description of quannal analysis is heavily based on Vantubergen (1966) and Newman and Carolyn Benz (1988). The third data analysis procedure applied to the data set was quannal analysis. The purpose of using this procedure was to determine whether certain types of people could be identified that favored different types of schools.

The factor analysis computer program used in this study was QUANNAL Vantubergen, 1966). This program places squared multiple correlation values in the principle diagonal as commonality estimates and conducts a Q-analysis. This approach is appropriate for the purpose of differentiating between people int erms of the shape of their profiles.

Five steps are used in a Q factor analysis.

Step 1 - An intercorrelation matrix is formed by correlating every person's ratings of the items with every other person's rating of items.

Thus, the eight ratings for respondent 1 were correlated with the ratings of the other 99 respondents. The same procedur $\epsilon$ was followed for each respondent.

Step 2 - The matrix of intercorrelations if submitted to factor analysis so that "persons" are variables and

items are observations. A principal axis solution is obtained. This result is submitted to a varimax rotation which produces orthogonal factors. On this basis, a factor represents a grouping of persons around a common pattern of sorting the items. Hence, a factor represents a type of "person" (Vantubergen, 1966).

Sub. No.	Two I	Factor II	Solution h	Sub.	Three	Factor	Solu	tion
1.	.22	.83	.95	No <sub>1</sub> .	I.30	II.	III	h <sup>2</sup>
2.	.92	.17	.88	2	.87	.16	.39	.93
3.	.98	13	.97	3.	.84	16	.50	.98
4.	.75	.49	.81	4.	.33	.37	.86	.99
5.	.82	.19	.71	5.	.37	.05	.90	.95
6.	06	.90	.82	6.	04	.91	.03	.83
7.	.86	.09	.76	7.	.97	.14	.14	.99
8.	.17	.92	.88	8.	02	.87	.39	.91
<b>%</b> Total				<b>%</b> Total				
Var.	48	34	82	Var.	34	32	27	93

:**\*** 

The factor analytic model constructs hypothetical types of "persons" based on the way the actual people interviewed rated the items. One can group people by assigning them to the type that they are most like, i.e., the factor on which they have the highest loading.

Step 3 - Each pattern of items associasted with each by weighting each item response of each item response of each of the persons most highly associated with a given factor by the degree to which they are loaded on that factor, the greater is the weight. These weighted responses are summed across each item separately. This procedure produces an item array of weighted responses for each factor in the rotated factor analysis solution selected. The arrays of weighted repsonses are then converted to z-scores (Vantubergen, 1966).

1

Hypothetical types constructed by the factor analytic model is based on a weighted pattern of the items (hypothetical types). The more a person's rating is like the hypothetical type, the more weight it receives in the average. The specific weight given is calculated as follows:

weight 
$$= \frac{r}{2}$$
 where:  $r = loading$ 

The weighted average is called an item factor array.

The persons used to estimate an array are highly associated with that type, but they are not associated to a high degree with any of the other types. For the persons selected, the square of the loading on that factor should approach the communality  $h_2$ . The arrays of weighted item ratings are converted to z scores. The array of z scores for each type is called the factor array.

Step 4 - The arrays of item z - scores for each factor (factor arrays) are ordered from most rejected for each factor. This provides a hierarchy of item acceptance for each factor or type of "persons" (Vantubergen, 1966).

The following are examples of hypothetical types of "persons" that the factor analytic model would construct:

••

	-	Types		
Items	<b>▲</b>	TT		
University 1	1.02	24	.72	
University 2	1.53	1.03	1.54	
University 3	.42	.31	-1.03	
University 4	06	.32	51	
University 5	-1.08	-1.35	-1.54	
University 6	.80	1.20	.5	
University 7	-1.20	.02	6	
University 8	.70	1.50	2.0	

When ordered in terms of the z-scores, the factor array becomes a hierarchy of items that are rated for each of the factors or types. The following is an example of the first typology (Type I):

Item	
University	2
University	1
University	6
University	8
University	3
University	4
University	5
University	7
	Item University University University University University University University University

Similar results were obtained for each type.

Step 5 - The arrays of item z-scores (factor arrays) for each type are compared by subtraction for each pair of factors. This produces arrays of difference scores for each pair of factors. This provides the basis for differentiating one factor or type of person from another Vantugergen, 1966).

This is accomplished by comparing the types dealing with the following questions:

- What items differentiate one type from another type?
   What items differentiate one type from all other types?
- 3. What items or areas of agreement seem to cut across all of the types?

Question 1 is dealt with by comparing the array for all types taken two at a time. The 2-scores for each pair of universities are subtracted and ranked according to absolute differences. To illustrate, consider the following:

	TYPE II	Type I-Type II		
1.02	24	1.26	University	1
-1.20	.02	1.22	University	7
.70	1.80	.80	University	8
1.53	1.03	. 50	University	2
.80	1.20	.40	University	6
-1.08	-1.35	.27	University	5
06	.32	. 38	University	4
.43	.31	.12	University	3

Similar analyses are conducted for all other comparisons.

Question 2. Question 2 was addressed by examining those items that are higher (or lower) in the array for one type than they are in the arrays for all other types. This process is similar to the process followed in Question 1. That is, the Z scores of Type I are compared to the average Z scores for Types II and III.

Question 3. To the extent that the z-scores for all types are nearly equal, one assumes agreement. A consensus item would be one in which the difference between the largest z-score given that item by one of the types and the smallest z score is less than 1.00. In our example, the consensus items would be the following:

Rating of Universities	Maximum Difference	λverage 2-Scores λcross Types
University 5	.46	1.32
University 2	. 50	1.37
University 6	.70	.83
University 4	.83	.08

The average Z-scores of the consensus items and the Zscores of the differentiation items, which resulted from addressing Questions 1, 2, 3, are used to describe the types. That is, the universities corresponding to the aforementioned Zscores are used to identify types.

#### Results of Quannal Analysis

Three Q-factor analyses were computed. One analysis was based upon the ratings of the eight universities, the second on demographic variables, and the third on the university and

demographic variables together. On all three of the Q-factor analyses, only one typology emerged.

In the first analysis, all of the 100 subjects were identified in Type I. In the second analysis, 99 of the 100 were identified in Type I. In the analysis combining the universities and demographic variables, 98 of the subjects were identified in Type I. As one can see form these results, only one type consistently emerged; therefore, we were unable to use differences in types as predictor variables. A multiple regression analysis by Fraas on the impact of the demographic variable of the data further validates the homogeneity of this sample.

Since we were in a desperate search for more than one type, it was suggested that we try a cluster approach, which tends to produce more than one type. Ward's (1963) clustering program takes a set of N objects, which are measured on a number of different variables, and attempts to optimally group them from N to N-1, etc. The groupings are based upon maximizing the average intergroup distance, while minimizing the average intragroup distance.

The approach begins by defining each object as a group. These N groups are then reduced by one, until all persons have been classified into one of two groups. More detail of this approach can be found in SAS, as well as Veldman (1967).

Using the clustering program, three cluster analyses were completed. When using a cluster analysis, one has to decide on the number clusters one wants in the solution. The decision used for this study was that no cluster would contain less than five people.

The first cluster analysis, using the universities' ratings and the three demographics, produced four clusters with 27 people in cluster one, 56 in cluster two, 11 in cluster three, and 6 in cluster four. These four clusters accounted for 61% of the variance for all groupings. The second cluster analysis, based upon universities' ratings, produced three clusters with an  $\mathbb{R}^2$ equal to .55, with 58 individuals in cluster one, 36 in cluster two, and 7 in cluster three. The third cluster analysis, based

upon demographics alone, produced only two clusters with almost everyone loading on cluster one. Therefore, it was not considered.

The four clusters produced by the first cluster analysis were used as predictor variables to predict the ratings of each of the eight universities, the eight regression equations

Table	3
-------	---

: :

## Multiple linear Regression Results for Model 2

	Variable	R	egression efficients	T Value
	<b>X</b> 1	•	1.410	12.21*
	X2		1.125	9.74*
	X <sub>3</sub>	1.440	12.4	7*
	X4	1.775	15.3	7*
	X5	1.015	8.7	9*
n en	X6	1.000		
Constant	E		-3.36	

n = 800 $R^2 = .58$  $df_d = 695$ 

Contraction of the second second

. . .

一般を こ

## \* Statistically significant at the .01 level.

produced the following values: .12, .27, .17, .18, .18, .26, .34, and .28. When the clusters from the second cluster analysis containing three clusters, were used as predictor variables, they yielded the following  $R^2$  values: .03, .18, .14, .15, .16, .20, .30, and .18. Since the use of cross-validation procedures would produce even lower values, those procedures were not implemented.

٠.

and a second s

6.6

and the second second

#### Discussion

34

The conjoint analysis and the multiple regression model with a surrogate person vector produced identical estimates for the five university attributes. The multiple regression procedure that incorporated a surrogate person vector was better able to predict the holdout universities. Thus, these results seem to imply that if a university administration wants to obtain information on which university attributes are most important to their students, either conjoint analysis or a multiple regression model with a surrogate variable is an appropriate procedure.

With this data set the Q-factor analysis failed to provide useful information. The classifying of student by type did not allow for a high degree of explanation of the ratings of the various hypothetical universities. The use of Q-factor analysis, however, may provide insight into the university selection process by students if various groups are identifiable.

Three points should be noted with regard to future research. First, a multiple linear regression model with a surrogate person vector is a valuable procedure to use to determine which university attributes are important to students when selecting a university. The inclusion of the surrogate person variable did improve the researchers' ability to predict the ratings of the holdout universities. Further studies in this area with more detailed attributes would be informative.

Second, unless various groups of students rate the universities differently, Q-factor analysis obviously will not provide useful information. If such groups exist, however, the information may provide university administrators with some insight into what type of students prefer their particular university.

Third, the conjoint and regression analyses are really asking different questions that the Q-factor analysis. The conjoint and regression analyses are attempting to determine which of the university characteristics are most important. The Q-factor analysis attempts to determine if there are various typologies based on the students' university ratings. This third point leads to an often discussed conclusion. Determining the preferable research method is dependent upon the question of interest. In other words, the research question has to dictate the methodology.

and the second second

. . x w z

and the second s

.

#### References

- Bretton-Clark. (1987). Conjoint analyzer. New York. Conant, J., Brown, J., & Mokwa, M. (1985), Summer). Students are important consumers: Assessing satisfaction in a higher educaton context. Journal of Marketing Education, 13-20.
- Fraas, J. W. & Paugh, R. (1989). <u>Student perceptions of the</u> relative importance of selected attributes of an <u>Institution of higher education: a conjoint approach.</u> Ashland, Ohio: Ashland University, 1989. (ERIC Document Reproduction Service No. ED 312 960)
- Fraas, J. W., & McDougall, W. R. (1983). The use of one full full MLR movel to conduct multiple comparisons in a repeated measures design: An industrial application. Multiple Linear Regression Viewpoints, 12(1), 42-55.
- Fraas, J. W. & Newman, I. Conjoint Analysis: A study of the effects of using person var ables. Ashlan, Ohlo: Ashland University, 1989. (ERIC Document Reproduction Service No. ED 312 961)
- Green, P., Goldbert, S., & Montemajor, M. (1981). A hybrid utility estimation model for conjoint analysis. estimation model for conjoint analysis.
- Journal of Marketing, 45, 33-41. Green, P., Carroll, J. D., & Goldberg, M. (1981). A gener approach to product design optimization via conjoint **A** general
- analysis. Journal of Marketing, 45, 17-37. Green, P., & Strinivasan, V. (1978). Conjoinjt analysis in consumer research: Issues and outlook. Journal of Consumer Research, 5, 103-124.
- Kuh, G., Comes, M., & Lundquist, I. (1984, Winter). What prospective students really need to know about institutional quality. <u>College and University</u>, 167-175 Leigh, T., Mackay, D., & Summer, J. (1981). Realiability and
- validity of conjoint analysis and self-explicated weights:
- A comparison. Journal of <u>Marketing Research</u>, 21, 456-462. McNeil, K., Kelly, F., 4 MoNeil, J. (1975, <u>esting research</u> <u>hypotheses using multiple linear regressio</u>. Carbondale, IL: Southern IllInois University **QSS**.
- Newman, I. & Benz, C. R. Multivariate evaluation design: suggested technique for mental health systems. ApriI 25, 1987. Ohio Academy of Science, Canton, Ohio.
- Pedhazur, E. J. (1977). Coding subjects in repeated measures designs. <u>Psychological Bulletin</u>, <u>84</u>, 298-305.
- Tierney, M. L. (1964, November. Ten questions to ask when choosing a colle<del>ge. <u>Money</u>, 133-134.</del> Vantubergen, G. N. (1966). Quannal: A computer program for
- Q analysis. Mass Communication Research Bureau, School of Journalism, University of Iowa. Mimeo.
- Veldman, D. (1967). Fortran Programming for the behavioral <u>sciences.</u> New York: Holt, Rinehart, & Winston. Ward, J. d. (1963). Hierarchical grouping to potimize an objective function. American Statistical Association Journal, 58, 236-244.

Williams, J. D. (1989). Multiple comparisons in higher dimension designs. Monograph Series #5. <u>Multiple Linear</u> <u>Regression Viewpoints, 10</u>(3).

Regression Viewpoints, 10(3). Williams, J. D., & Williams, J. A. (1985a). Testing hypotheses in a repeated measures design on employee attitudes with large samples. <u>Multiple Linear Regression Viewpoints</u>, 13(2), 1-20.

<u>13(2), 1-20.</u> Wliams, J. D., & Williams, J. A. (1985b). Testing hypotheses in a repeated measures design: An example. <u>Multiple</u> <u>Regressions Viewpoints, 13(2), 35-46.</u>

:\*

•

rev.042291

 $\frac{1}{2} \left( \frac{1}{2} + \frac{1}{2} \right) = \frac{1}{2} \left( \frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) \left( \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) \left( \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) \left( \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) \left( \frac{1}{2} + \frac{1$ 

andra Stanton Stanton Stanton Stanton

 $\frac{1}{2} = \frac{1}{2} \frac{$ 

## Relationship Between Multiple Regression, Path, Factor, and LISREL Analyses

Randell E. Schumacher University of North Texas

#### Abstract

A basic knowledge of multiple regression concepts permits further understanding of path, factor, and lisrel analyses. Specifically, standardized partial regression coefficients (beta weights) as applied in path, factor, and lisrel analyses are presented. The multivariable methods have in common the general linear model and are the same in several respects. First, they identify, partition, and control variance. Second, they are based upon a linear combination of variables. And third, the linear weights can be computed based on standardized partial regression coefficients.

Multiple regression or the general linear model approach to the analysis of experimental data in educational research has become increasingly popular since 1967 (Bashaw and Findley, 1968). In fact today, it has become recognized as an approach that bridges the gap between correlational and analysis of variance thought in answering research hypotheses (McNeil, Kelly, & McNeil, 1975). Statistical textbooks in psychology and education often present the relationship between data analysis with multiple regression and analysis of variance (Draper & Smith, 1966; Williams, 1974a; Roscoe, 1975; Edwards, 1979). Graduate students taking an advanced statistics course are therefore provided with the multiple linear regression framework for data analysis. Given their knowledge of multiple linear regression techniques applied to univariate analysis (one dependent variable), their understanding can be extended to the relationship of multiple linear regression to various multivariate statistical techniques (Kelly, Beggs, McNeil, with Eichelberger 4 Lyon, 1969, pps 228-248; Newman, 1988). The article therefore expands upon this understanding and indicates the importance of the standardized partial regression coefficient (beta weight) in multiple linear regression as it is applied in path, factor, and lisrel analyses.

#### MULTIPLE REGRESSION

•

Multiple regression techniques require a basic understanding of sample statistics (n, mean, and variance), standardized variables, correlation (Pedhazur, 1982, pp 53-57), and partial correlation (Cohen & Cohen, 1975; Houston & Bolding, 1974). In standard score form the multiple regression equation is:

> z = bz y x

The relationship between the correlation coefficient, the unstandardized regression coefficient and the standardized regression coefficient is:



For two independent variables, the regression equation with standard scores is:

And the standardized partial regression coefficients are computed by:

<b>b</b>	r - r r yl y2 12	<b>b</b> -	r - y2	r r y1 12
	2	2		2
N.	1 - r		1 - :	r
	12			12

The correlation between the original and predicted scores is given the special name Multiple Correlation Coefficient. It is indicated as:

stanting y w R

.

And the Squared Multiple Correlation Coefficient is related as follows:

MULTIPLE REGRESSION EXAMPLE

A multiple linear regression example using a correlation matrix as input (SPSSX User's Guide, 3rd Edition, 1988, Chapter 13) is in the appendix. The results are:

1064

**2**.. R + br + Ь b r r 2 y2 y.123 3 y3 1 y1 = (.423) .507 + (.363) .481 + (.040) .276... 2 .40 R y.123

A systematic determination of the most important set of variables can be accomplished by setting the partial regression weight of each variable to zero. This approach and other alternative methods are presented by Kelly, Beggs, & McNeil et al (1969) and Darlington (1968).

In summary, regression techniques have been shown to be robust (Bohrnstedt & Carter, 1971); applicable to contrast coding (Lewis & Mouw, 1978); dichotomous coding (McNeil, Kelly, & McNeil, 1975); and ordinal coding (Lyons, 1971) research situations. Multiple regression can also be viewed as a special case of path analysis.

#### PATH ANALYSIS

Sewall Wright is credited with the development of path analysis as a method for studying the direct and indirect effects of variables (Wright, 1921, 1934, 1960). Path analysis is not a method for discovering causes, rather it tests theoretical relationships called "causal modeling". The specified model establishes causal relationships among the variables when:

- a. temporal ordering exists
- b. covariation (correlation) is present
- c. controlled for other causes

Model specification is necessary in examining multiple variable relationships. In the absence of a model, many different relationships among variables can be postulated with many different path coefficients being selected. For example, in a three variable model the following four relationships could be postulated:



The four different models have been considered without reversing the order of the variables. How can one decide which model is correct? Fath analysis doesn't provide a way to specify the model, but rather estimates the effects once the model has been specified "a priori". Path coefficients in path analysis take on the values of a product-moment correlation and/or standardized regression coefficients in a model (Wolfle, 1977). For example given model (d):

b = p 1 y1	Ъ 2	P y2		r " 12	P 12
THEN :					
x 1					
					•
X 2	v i	•	,	. :	

A path model is specified by the researcher based on theory or prior research. Variable relationships once specified, in standard score form, become standardized regression coefficients. In multiple regression, a dependent variable is regressed in a single analysis on all the independent variables. In path analysis one or more multiple regression analyses are performed. Path coefficients are computed based upon only the particular set of independent variables that lead to the dependent variable under consideration. As in regression analysis, path analysis can use dichotomous and ordinal data in the causal model (Boyle, 1970; Lyons, 1971).

#### MODEL SPECIFICATION

Path models permit diagramming how a particular set of independent variables lead to a dependent variable under consideration. How the paths are drawn determine whether the independent variables are correlated causes (unanalyzed), mediated causes (indirect), or independent causes (direct). The model can be tested for the significance of path coefficients (Pedhazur, 1982, pp 58-62) and a goodness-of-fit criteria (Marascuilo & Levin, 1983, pp 169-172; Tatsuoka & Lohnes, 1988, pp 98-100) which reflects the significance between the original and reproduced correlation matrix. This process is commonly called decomposing the correlation matrix (Asher, 1976, pp 32-34) according to certain rules (Wright, 1934).

#### PATH ANALYSIS EXAMPLE

A four variable path analysis program is in the appendix. In order to calculate the path coefficients for the model, two regression analyses were performed. The model with the path coefficients is:



:

The original and reproduced correlations are presented in matrix form. The upper half represents original correlations and the lower half the reproduced correlations which include the regression of paths linking independent variables to the dependent variable.

VARIABLE	Y	<b>X1</b>	X2	<b>X3</b> '	
Y	1.000	.507	.401	.276	
<b>X1</b>	.423	1.000	.224	.062	Original
X2	.362	.224	1.000	.577	Correlations
X3	.040	070	. 593	1.000	·

Reproduced Correlations

The offect	riginal ts: dire	correlation (DE),	ions can indirect	be complet (IE), spur	ely "rep cious (8	produced" if and correl	all ated
	i seteni Lectrolit	in (w. 111) Antier and		•	•		
r = 12	(******) (** (** <b>P</b> ) (*** (********************************			s* w *	ng sa sa sa Sa sa	· · 224	
r	P + 31 DE	P P 32 21 IE				062	
r - 23	P + 32 DE	P P 31 21 8	t.			577	
r - 1y -	P + ¥1 DE	PP Y221 IE	+ p p Y3 31 IE	+ p p p Y3 32 21 IE		507	•;
r - 2Y	P + Y2 DE	P P 13 32 IE	+ p p 1 21 8	+ p p p Y3 31 21 8		481	
r - 31	P + . 13 De	P P 1131	+ P P 12 32 8	+ p p p Y1 21 32 8	+ P P Y2 21 8	P = .276 31	

In summary, path analysis can be carried out within the context of ordinary regression analysis and does not require the learning of any new analysis techniques (Asher, 1976, p32; Williams, 1974b). The advantage of path analysis is that it enables one to specify direct and indirect effects among independent variables. In addition, path analysis enables us to decompose the correlation between any two variables into simple and complex paths of which some are meaningful. Path coefficients and the relationship between the original and reproduced correlation matrix can also be tested for significance.

#### FACTOR ANALYSIS

Path models and the associated test of significance between original and reproduced correlations are used in confirmatory factor analysis. Factor analysis assumes that the observed (measured) variables are linear combinations of some underlying source variable (factor). In practice, one estimates population parameters of the measured variables from a sample (with the uncertainties of model specification and measurement error). A linear combination of weighted variables relates to multiple regression in a single factor model and to a linear causal system (path analysis - "multiple" multiple regressions) in multiple factor models. Path diagrams therefore permit representation of the causal relationships among factors and observed (measured) variables in factor analysis.

In general, the first step in factor analysis involves the study of interrelationships among variables in the correlation matrix. Factor analysis will address the question of whether these subsets can be identified by one or more factors (hypothetical constructs). Confirmatory factor analysis is used to test specific hypotheses regarding which variables correlate with which constructs (Long, 1983).

Factor analysis assumes that some factors, which are smaller in number than the number of observed variables, are responsible for the covariation among the observed variables. For example, given a unidimensional trait in a single factor model with four variables the diagram would be (Kim & Mueller, 1978a, p 35):

b = .677	Y	d = .735 Y	U Y
I	X	d <b>9</b> 17	U
	1	1	1
b = .402 1 b = .800 2	X 2	d <b>6</b> 00 2	บ 2
b = .535	х	d = .643	U
	3	3	3

#### WHERE :

ī

3

b = standardised regression coefficient
The variance of each observed variable is therefore comprised of the proportion of variance determined by the common factor and the proportion determined by the unique factor, which together equal the total variance of each observed variable. Therefore:

2 2 2 8 = b + d = 1 i i i

The correlation between a common factor and a variable is:

r = b F,X i i

The correlation between a unique factor and a variable is:

đ

- 1

1. 2 1 2 1

inger en Art

a state state

The correlation between observed (measured) variables sharing a common factor is:

r	-	;	Ъ	1	
X	, X		5. <b>1</b>		1
<u> </u>	1			,	

And finally, the variance attributed to the factor as a result of the linear combination of variables is:

	2	
2 h =	8 b i =	2 R
	M	<b>F</b> .1234

Where: M = number of variables

2 b = squared factor loadings i

Note: 5 b = eigenvalue i b = communality

Ł

#### FACTOR ANALYSIS EXAMPLE

A single factor analysis program with four variables in a correlation matrix format is in the appendix. The path diagram is the same as above (Kim & Mueller, 1978a, p 35) with the weights as follows:

b = .677 b = .402 b = .800 b = .535 Y 1 2 3

And, factor scores computed as:

F = bY + bX + bX + bXy 11 22 33

Multiplying the coefficients between pairs of variables gives the following correlation matrix:

VARIABLE	Y	<b>X</b> 1	<b>X</b> 2	хэ
<b>Y</b>	2 b 1	.27	.54	.36
1980 (1892) (1986) <b>X1</b> (1986) 2003	.27	2 b 2	. 32	.22
**************************************	.54	. 32	2 b 3	.43
<b>X3</b>	.36	.22	.43	2 5 4

3 South States and Management States and States and

The common factor variance is:

The unique factor variance is:

2 = 8 (1 - b) = .54 + .84 + .36 + .71 = .61F.1234
M
4

In summary, factor loadings (variable weights) are standardized regression coefficients. As such, linear weighted combinations of variables loading on a factor are used to compute factor scores (Kim & Mueller, 1978b p 60). The weights are also the correlation observed (measured) variables and the factor the between (hypothetical construct). If the variable correlations (weights) are squared and summed, they describe the proportion of variance determined by that factor. This is traditionally known as an eigenvalue, but termed communality in factor analysis. When all variables are standardized, then the linear weights are called standardized regression coefficients (regression analysis), path coefficients (path analysis), or factor loadings (factor analysis). The factor analysis approach is distinguished from regression or path analysis in that observed variable correlation is explained by a common factor (hypothetical construct). In factor analysis therefore the correlation between observed variables is the result of sharing a common factor rather than a variable being the direct cause (path analysis) or predictor of another (regression analysis).

LISREL

Linear structural relationships (lisrel) are often diagrammed by using multiple factor path models where the factors (hypothetical contructs) are viewed as latent traits (Joreskog & Sorbom, 1986, pp I.5-I.7). The lisrel model consists of two parts: the measurement model and the structural equation model. measurement model specifies how the latent variables The or hypothetical constructs are measured in terms of the observed (measured) variables and describes their measurement properties (reliability and validity). The structural equation model specifies the causal relationship among the latent variables and is used to describe the causal effects and the amount of unexplained variance. The listel model includes or encompasses a wide range of models, for example: univariate or multivariate regression models," confirmatory factor analysis, and path analysis models (Joreskog & Sorbom, 1986, pp I.3, I.9-I.12). Cuttance (1983) presents an overview of several lisrel submodels with diagrams and diagrams and explanations. Wolfle (1982) presents an indepth presentation of a single model to introduce and clarify lisrel analysis. The lisrel program therefore permits regression, path, and factor analysis whereby model specification and measurement error can be assessed.

REASURPENT BRROR TO THE STATE AND A CONTRACT AND A

#### Fuller (1987) extensively covers lisrel and factor analysis models and especially extends regression analysis to the case where the variables are measured with error. Wolfe (1979, pp 48-51) presents the relationship between lisrel, regression and path analysis especially in regards to how measurement error effects the regression coefficient (path coefficient). Errors of measurement in statistics have been studied extensively (Wolfe, 1979). Cochran (1968) studied it from four different aspects: (1) types of mathematical models, (2) standard techniques of analysis which take into account measurement error, (3) effect of errors of measurement in producing bias and reduced precision and what remedial procedures are available, and (4) techniques for studying error of measurement. Cochran (1970) also studied the effects of error of measurement on the squared multiple correlation coefficient.

#### LISREL-FACTOR ANALYSIS EXAMPLE

A LISREL factor analysis program with a correlation matrix as input is in the appendix. The factor analytic model in matrix notation is:

 $\begin{array}{ccc} \mathbf{X} = \mathbf{L} \mathbf{x} + \mathbf{q} \\ \mathbf{x} & \mathbf{d} \end{array}$ 

Where: X = observed variables L = structural weights (factor loadings) x = latent trait (factor) q = error variance (unique variance) d

The LISREL results are:

L = LANBDA X (structural weights-factor loadings) A. ¥ = .677 X = .402 X = .800 X = .535 1 2 3 q = THETA DELTA (unique factor variance) **b**. d **Y** = .54 **X** = .84 X = .36 X = .71 3 1 2 2 2 b = LAMEDA X (common factor variance) α. Y = .46 X = .16 X = .64 X = .29 1 2 3

The concept of model specification and goodness of fit pertains to the original correlation matrix and the estimated correlation matrix. 3% The estimated correlation matrix is: 200 00 2 00 0000 The state of the second s

n 1986 - Jacob Contractor and State

3,

i general i general a ser l'interne

المراجع المنافع المراجع المراجع

tare in the state

State Charles

A MARKED ME ARACTED			3.3		1.28
in stand	.272	1 A	1.		<u>s</u>
σ	. 542	·*• <b>.321</b> *	1. 1. 1. 1	стан 1	$a^{(1)} \rightarrow b$
s na stander og de se de	.362	.215	.427	(1,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2	he de la
del de la compañía d			· · · ·		1000

The original correlation matrix is:

.276	.062	.577
.481	.224	7
.507	7~-	1 s
	<ul> <li>(a)</li> <li>(b)</li> </ul>	

The goodness of fit index (GFI) using the unweighted least squares approach (ULS) is then computed as:

2 1 - 1/2 trace (S -  $\sigma$ ) 🛾 GFI 💠 🚥 👘 2 the military and the the 1 - 1/2 (1.308 - 1.02) GFI 

GFI = 1 - .041

dina dia kaominina dia kao Ny faoritra dia kaominina di

S S 🖉 💻

LISREL-REGRESSION ANALYSIS EXAMPLE

A LISREL regression program with a correlation matrix as input is in the appendix. The regression model in matrix notation is:

 $\mathbf{Y} = \mathbf{G} \mathbf{X} + \mathbf{z}$ 

Y = dependent variable Where: G - gamma matrix (beta weights) X = independent variables

s = errors of prediction (error variance)

1. 粘。

The LISREL results are the same as in the previous regression program: Sec. March

2 R G r r y.123 1 y1 y3 **y2** 2 = (.423) .507 + (.363) .481 + (.040) .276R y.123 2 .40 R y.123

The appropriate statistical method to use is often an issue of debate. It sometimes requires more than one approach to analyzing data. The rationale for choosing between the alternative methods of analysis is usually guided by research hypotheses or questions.

The multivariable methods discussed have in common the general linear model and are the same in several respects. First, they identify, partition, and control variance. Second, they are based on linear combinations of variables. And third, the linear weights can be computed based on standardized partial regression coefficients.

The multivariable methods however have different applications. Multiple regression seeks to identify and estimate the amount of variance in the dependent variable attributed to one or more independent variables (prediction). Path analysis seeks to identify relationships among a set of variables (explanation). Factor analysis seeks to identify subsets of variables from a much larger set (common/shared variance). Lisrel determines the degree of model specification and measurement error. The different methods were derived because of the need for prediction, explanation, common variance, model and measurement error assessment type applications.

Multiple regression techniques are robust except for model specification and measurement errors (Borhnstedt & Carter, 1971). Multiple regression techniques are also useful in understanding path, factor, and LISREL applications. LISREL permits regression, path, and factor analyses whereby model specification and measurement error can be assessed. Lisrel also permits univariate or multivariate least squares analysis in either single sample or multiple sample (across populations) research settings. An understanding of multiple regression and general linear model techniques can therefore greatly facilitate ones understanding of the testing of research questions in multivariable situations.

## 

.

- 4

ų.

and good and good and a second a

a second second

MULTIPLE REGRESSION PROGRAM TITLE REGRESSION WITH CORRELATION MATRIX INPUT COMMENT VARIABLE MEANS=0; VARIANCES=1; CONSTANT=0 MATRIX DATA VARIABLES-Y X1 X2 X3/N-100 BEGIN DATA 1.000 .507 1.000 and the second se 1.000 .481 .224 and the statement .276 .062 .577 1.000 END DATA MATRIX-IN (\*) / MISSING-LISTWISE/ VARIABLES-Y X1 X2 X3/ REGRESSION MATRIX=IN(\*)/ DEPENDENT-Y/ ENTER X1 X2 X3/ 1997年1月1日日本主義主義 A Contract of the second and the second secon FINISH and the second ere e la c the terms to be a second second PATH ANALYSIS PROGRAM ONE 1.5 1. St. 188 (1875) 

A. VARIABLE 3 REGRESSED ON VARIABLES 1 AND 2 TITLE PATH ANALYSIS EXAMPLE WITH CORRELATION MATRIX INPUT

COMMENT VARIABLE MEANS=0; VARIANCES=1; CONSTANT=0 MATRIX DATA VARIABLES=Y X1 X2 X3/N=100 BEGIN DATA 1.000 .507 1.000 .000 .507 1.000 .481 .224 1.000 

.276 .062 .577 1.000 END DATA REGRESSION MATRIX=IN(\*)/ MISSING-LISTWISE/ VARIABLES-Y X1 X2 X3/ DEPENDENT=X3/ ENTER X1 X2/

FINISH

#### PATH ANALYSIS PROGRAM TWO

B. VARIABLE Y REGRESSED ON VARIABLES 1, 2, AND 3

TITLE PATH ANALYSIS EXAMPLE WITH CORRELATION MATRIX INPUT COMMENT VARIABLE MEANS=0; VARIANCES=1; CONSTANT=0 MATRIX DATA VARIABLES-Y X1 X2 X3/N=100 BEGIN DATA 1.000 .507 1.000 .224 1.000 .481 .276 .062 .577 1.000 END DATA REGRESSION MATRIX=IN(\*)/ MISSING-LISTWISE/ VARIABLES=Y X1 X2 X3/ DEPENDENT=Y/ ENTER X1 X2 X3/

FINISH

FACTOR ANALYSIS PROGRAM

```
TITLE FACTOR ANALYSIS EXAMPLE WITH CORRELATION MATRIX INPUT
COMMENT VARIABLE MEANS=0; VARIANCES=1; CONSTANT=0
MATRIX DATA VARIABLES-Y X1 X2 X3/N=100
BEGIN DATA
1.000
 .507
      1.000
       .224 1.000
 .481
        .062
              .577 1.000
 .276
END DATA
FACTOR VARIABLES-Y X1 X2 X3/
       MATRIX=IN (COR=*) /
       CRITERIA-FACTORS (1) /
       EXTRACTION-ULS/
       ROTATION-NOROTATE/
       PRINT CORRELATION DET INITIAL EXTRACTION ROTATION/
       FORMAT SORT/
       PLOT EIGEN/
       TINISH
```

#### LISREL FACTOR ANALYSIS PROGRAM

TITLE 'LISREL FACTOR ANALYSIS WITH CORRELATION MATRIX INPUT INPUT PROGRAM NUMERIC DUMMY END FILE The North Mark Mark Alexand Marka END INPUT PROGRAM USERPROC NAME=LISREL DATA FOR GROUP ONE DA NG-1 NI-4 NO-100 心惊惊,聋 LA 8 2 C 'Y' 'X1' 'X2' 'X3' . · · · Sec. 1 1. a. 2. KM SY 1.000  $S_{A}(t) \leq \frac{1}{2}$ .507 1.000 the house of , Ar site .481 .224 1.000 .062 .577 1.000 Constant Constant Constant .276 MO NX-4 NK-1 TD-DI, FR PH-ST The War Care and States LK 'FACTOR' PA LX 4 \* 1 OU ULS SE TV PC RS VA FS SS MI END USER an taan 11 te LISREL REGRESSION ANALYSIS PROGRAM And John Car TITLE 'LISREL REGRESSION ANALYSIS WITH CORRELATION MATIRX' INPUT PROGRAM and the second sec 2012 NUMERIC DUMMY END FILE A The Contract The the state of the second second

END INPUT PROGRAM USERPROC NAME-LISREL DATA FOR GROUP ONE DA NG-1 NI-4 NO-100 LA 'Y' 'X1' 'X2' 'X3' KM SY 1.000 .507 1.000 .224 1.000 .481 .276 .062 .577 1.000 MO NY-1 NX-3 PS-DI OU ULS SE TV PC RS VA SS MI TO END USER

医子宫 建立 计正式 网络马克马属德国德瓦

THE MARRING CLEARAN MY ...

r an tao an da anna an Airline an Airline an Airline. An 1975 - Airline Anna Airline an Airline an

- Asher, H.B. (1976). Causal modeling. Sage Publications: Beverly Hills: CA.
- Bashaw, W.L. & W.G. Findley (1968). Symposium on general linear model approach to the analysis of experimental data in educational research. Project No. 7-8096. U.S. Department of Health, Education, and Welfare, Washington, D.C.
- Bohrnstedt, G.W. & T.M. Carter (1971). Robustness in regression analysis. In H.L. Costner (Ed), Sociological Methodology, Jossey-Bass, pp 118-146.
- Boyle, R.P. (1970). Path analysis and ordinal data. American Journal of Sociology, vol. 75(4), 461-480.
- Cohen, J. & P. Cohen (1975). Applied multiple regression/ correlation analysis for the behavioral sciences. Lawrence Erlbaum: NJ.
- ochran, W.G. (1968). Errors of measurement in statistics. Technometrics, 10(4), 637-666.
- ochran, W.G. (1970). Some effects of errors of measurement on multiple correlation. Journal of the American Statistical Association, 65(329), 22-34.
- uttance, P.F. (1983). Covariance structure and structural equation modelling in research: a conceptual overview of LISREL modelling. Multiple Linear Regression Viewpoints, 12(2), 1-63.
- srlington, R.B. (1968). Multiple regression in psychological research and practice. Psychological Bulletin, 69, 161-182.
- caper, N.R. 4 H. Smith (1966). Applied regression analysis. John Wiley & Sons: New York, NY.
- iwards, A.L. (1979). Multiple regression and the analysis of variance and covariance. W.H. Freeman: San Franscisco, CA.
- iller, W.A. (1987) Measurement error models. John Wiley & Sons: New York, NY.
- vuston, S.R. & J.T. Bolding, Jr. (1974). Part, partial, and multiple correlation in commonality analysis of multiple regression models. Multiple Linear Regression Viewpoints, 5(3), 36-40.

Joreskog, K.G. & D. Sorbom (1986). LISREL VI USER'S GUIDE: Analysis of linear structural relationships by maximum likelihood, instrumental variables, and least squares methods. Scientific Software: Mooresville, Indiana.

States and

An Sec Sec.

a second

:

1

NY CARACTER AND A CARACTER ANTE

and the second

- Kelly, F.J, D.L. Beggs, K.A. McNeil, T. Eichelberger, g L. Lyon (1969). Multiple regression approach. Southern Illinois University Press: Carbondale, IL. Texaster
- Kim, J. & C. Mueller (1978a). Introduction to Factor Analysis. Sage Publications, Beverly Hills: CA.
- Kim, J. & C. Mueller (1978b). Factor Analysis. Sage Publications, Beverly Hills: CA.
- Lewis, E.L. & J.T. Mouw (1978). The use of contrast coefficients. Southern Illinois University Press, Carbondale, IL. and the second sec
- Long, J.S. (1983). Confirmatory Factor Analysis. Sage Publications, Beverly Hills: CA.
- Lyons, M. (1971). Techniques for using ordinal measures in regression and path analysis. In H.L. Costner (Ed), Sociological Methodology, Jossey-Bass, pp 147-171.
- Marascuilo, L.A. 6 J.R. Levin (1983). Multivariate statistics in the social sciences: a researcher's addition guide. Brooks/Cole Publishing: Belmont: CA.
- McNeil, K.A., F.J. Kelly, J.T. McNeil (1975). Testing research hypotheses using multiple linear regression. Southern Illinois University Press: Carbondale, IL Content and the second
- Newman, I. (1988). There is no such thing as multivariate analysis: all analyses are univariate. Presidents analysis and the second secon Address at Mid-Western Educational Research Association, October 15, 1988, Chicago, IL.
- Pedhazur, E.J. (1982). Multiple regression in behavioral 🛲 🖘 👾 research: explanation and prediction (2nd ed.). Holt, Rinehart, & Winston: New York, NY.
- Roscoe, J.T. (1975). Fundamental research statistics for the behavioral sciences (2nd ed.). Holt, Rinehart, 6 as an and a Winston: New York, NY.
- SP88X Users' Guide, 3rd ed. (1988). McGraw-Hill: New Content of the second seco York, NY.
- Tatsuoka, M.M & P.R. Lohnes (1988). Multivariate analysis: techniques for educational and psychological research (2nd ed.). Macmillan Publishing Company: New York, NY.

- Williams, J.D. (1974a). Regression analysis in educational research. MSS Information Corporation: New York, NY.
- Williams, J.D. (1974b). Path analysis and causal models as regression techniques. Multiple Linear Regression Viewpoints, 5(3), 1-20.
- Wolfle, L.M. (1977). An introduction to path analysis. Multiple Linear Regression Viewpoints, 8(1), 36-61.
- Wolfle, L.M. (1979). Unmeasured variables in path analysis. Multiple Linear Regression Viewpoints, 9(5), 20-56.
- Wolfle, L.M. (1982). Causal models with unmeasured variables: an introduction to LISREL. Multiple Linear Regression Viewpoints, 11(2), 9-54.
- Wright, 8. (1921). Correlation and causation. Journal of Agricultural Research, 20, 557-585.
- Wright, S. (1934). The method of path coefficients. Annals of Mathematical Statistics, 5, 161-215.
- Wright, 8. (1960). Path coefficients and path regression: alternative or complementary concepts? Bicmetrics, 16, 189-202.

# The Case for Non-Zero Restrictions in Statistical Analysis

#### Keith McNeil New Mexico State University

One of the many advantages of MLR is its versatility and its ability to answer a vast array of questions. Unfortunately, most researchers fall into the habit of asking a small aubaet of very similar questions. The question being tested should be atated first, but can be identified from the full model and the restriction(s) placed on that full model. While the restrictions can take on any numerical value, almost all applications use the "default " value of zero:

- 1.  $a_1 = a_2$  or  $(a_1 a_2 = 0)$  (t-test)
- 2.  $a_1 = 0$  (Correlation)
- 3.  $\mathbf{s}_1 = \mathbf{a}_2 = \mathbf{a}_3 = \dots \mathbf{a}_1$  or  $(\mathbf{a}_1 \mathbf{a}_2 = \mathbf{a}_2 \mathbf{a}_3 = \mathbf{a}_3 \mathbf{a}_2 = \dots \mathbf{0})$
- (F-test)
- 4.  $(a_1 a_2) = (a_3 a_4)$  or  $((a_1 a_2) (a_3 a_4) = 0)$ (interaction)

The focus of this paper will be on the utility of making a non-zero restriction. Why the zero restriction occurs so frequently will be questioned and hopefully researchers and statisticians will see how the zero restriction limits the conclusions of the research. The argument will be made for making non-zero restrictions, resulting of course, from "non-zero" research hypotheses. The argument will be made for each of these statistical procedures: two group t teat, Pearaon correlation, aingle population mean, one-way analysie of variance, and interaction.

#### Two Group t Test

Perhaps the most widely used design compares the performance of two groups. The research hypothesie takes the following form: Research Hypothesis 1: For a given population, the New treatment ie better than the Traditional treatment on Y. (See Note 1 for discussion of directional hypothesie testing.)

Full Model:  $Y = a_1N + a_2T + E_1$ 

Where Y = criterion of interest,

N = 1 if subject in New treatment; 0 otherwise, and

T = 1 if subject in Traditional treatment; 0 otherwise. The research hypothesis implies that the sample mean for N should be greater than the sample mean for T, or  $a_1 > a_2$ , or  $a_1 - a_2 > 0$ . Restriction:  $a_1 = a_2$ , or  $(a_1 - a_2 = 0)$ Forcing the restriction into the full model results in: Restricted Model:  $Y = a_1N + a_1T + E_2$ But since the two vectors (N and T) are multiplied by the same

weights, the vectors can be added first. But N + T equals the Unit vector (or everyone). Therefore:

Restricted Model:  $Y = a_1U + E_2$ 

There are two linearly independent pieces of information in the full model. Forcing the one restriction on the full model results in one linearly independent piece of information in the restricted model. (See Note 2 for test of significance.)

A significant drop in the  $R^2$  from the full model to the restricted model results in a significant F. If the sample means are in accord with the anticipated result, then Research Hypothesis 1 can be held as tenable and the conclusion would be: For the given population, the New treatment is better than the Traditional treatment on Y. But all that has been said is that the New treatment is better than the Traditional treatment. We do not know how much better; all we know is that the two treatments are not equally effective.

But what if the cost of the two treatments is not the same? The Traditional treatment has surely been somewhat effective in the past. The New treatment will surely require some additional cost in the form of special inservice, purchase of new materials, acceptance by teachers, students, and community, etc. Before the Traditional treatment is replaced by the New treatment, perhaps the researcher should demonstrate that there is, say, more than a five-point superiority of the New treatment over the Traditional treatment.

When a non-zero research hypothesis is proposed, other researchers and statisticians often ask for the justification for the actual non-zero value chosen, as they should. But why should more justification be required for a non-zero value than for a zero Or looking at the issue form the other side, why are value? researchers allowed to test a zero value with little or no When one realizes that zero is only one of justification. an infinite number of values, then one realizes that the same amount of justification should be required of a zero value as of a nonzero value. Furthermore, when one attempts to justify the zero value restriction, one may realize that zero is not the value of interest. Those researchers who have been defaulting with zero should know how to choose a value, but may not. It is not the intent of this paper to illustrate how one determines the magnitude of the value tested in the research hypothesis, although a few suggestions will be provided.

In the case where there was an expectation of a five-point superiority, the research hypothesis would be: Research Hypothesis 2: For a given population, the New treatment is more than five points better than the Traditional treatment on Y.

Full Model:  $Y = a_1N + a_2T + E_3$ The research hypothesis implies that the sample mean for the New treatment is more than five units greater than the sample mean for the treatment or,  $a_1$  greater than  $(a_2 + 5)$  or  $(a_1 - a_2 > 5)$ Restriction:  $a_1 = a_2 + 5$ , or  $(a_1 - a_2 = 5)$  or  $(a_2 = a_1 - 5)$ Restricted Model:  $Y = a_1N + (a_1 - 5)T + E_4$  $Y = a_1N + a_1T - 5T + E_4$  $(Y + 5T) = a_1(N + T) + E_4$  $(Y + 5T) = a_1U + E_4$ 

There are two linearly independent pieces of information in the full model. Forcing the one restriction on the full model results in one linearly independent piece of information in the restricted model. (See Note 1 for test of significance.)

Notice that the full model in Research Hypothesis 1 is exactly the same as the full model in Research Hypothesis 2. The number of restrictions is also the same, resulting in the same number of What is different, though, is the nature of degrees of freedom. the restriction and hence the restricted models are different. The two research hypotheses are both "correct" and equally "valid" they just test two different hypotheses. Research Hypothesis 2 provides a more definitive conclusion.

"cost" of any treatment may be difficult to The actual But one must remember that Research Hypothesis 1 determine. reduces to the default assumption that the "costs" are equal. The choice of a research hypothesis leading to a restriction of  $(a_1 - a_2)$  $a_1 = 0$ ) should be defended as much as a research hypothesis leading to a restriction of  $(a_1 - a_2 = a_2 = a_1 - a_2 = a_2 = a_1 - a_2 = a_2 - a_1 - a_2 = a_2 - a_2 -$ The restriction  $(a_1 - a_2 = 0)$  has become a widely used default value, but we must realize that it is only one of an infinite number of values.

#### Pearson Correlation

The usual application of the Pearson correlation hypothesis 18:

Research Hypothesis 3: For a given population, the linear correlation between X and Y is greater than zero.

Full Model:  $Y = a_0U + a_1X + E_a$ 

The research hypothesis implies that the slope of the line of beat fit in the sample is positive, or a, > 0.

 $a_1 = 0$ **Reatriction:** 

Restricted Model:  $Y = a_0U + 0X + E_a$ 

$$Y = a_0 U + E_e$$

There are two linearly independent pieces of information in the full model. Forcing the one restriction on the full model results in one linearly independent piece of information in the restricted model.

If the F test is significant, then one concludes that the research hypothesis is tenable, that the linear correlation between X and Y is greater than 0, or that the change in Y per unit change in X is greater than 0; but we do not know how much greater than There may be reasons for wanting to know if the correlation is 0. greater than a particular value. For instance, if the correlation under consideration is either a validity coefficient or 8 definitely want reliability coefficient, then we would 8 correlation coefficient above some specified value, such as: Research Hypothesis 4: For a given population, the linear correlation between Y and the Retest of Y is greater than .80. Full Model:  $Y = a_n U + a_n R + E_n$ 

Restriction: Restricted Nodel  $R_2 = .64$ The research hypothesis implies that the restricted model  $R^2$  will be  $(.80)^2$  or .64. The formula in Note 2 can be used when testing this hypothesis for significance.

Consider Research Hypothesis 5: For a given population, there is more than a .6 unit change in Y for every unit change in X. In this case the models would be:

Full Model:  $Y = a_0U + a_1X + E_8$ Restriction:  $a_1 = .6$ Restricted Model:  $Y = a_0U + .6X + E_9$ 

 $(Y - .6X) = a_0U + E_9$ 

Notice that the full model in Research Hypotheses 3 and 4 is exactly the same as Research Hypothesis 5. The number of restrictions is also the same; resulting in the same number of degrees of freedom. What is different, though, is the nature of the restriction and hence the restricted models are different. The three research hypotheses are all "correct" and equally "valid" - they just test three different hypotheses. Research Hypotheses 4 and 5 provide more definitive conclusions.

The desired correlation (reliability, validity, etc.) may be difficult to determine, but should be no more difficult to justify than justifying the default value of 0. Just because  $a_1 = 0$  has been used in the past does not justify its use, particularly with hypotheses about reliability and validity.

#### Single Population Mean

The usual application of the single population mean hypothesis is:

Research Hypothesis 6: For a given population, the population mean is greater than a particular value, 8.

Here 8 is some meaningful value, depending on the given circumstances. Maybe the researcher wants to establish that the population mean height is greater than 72 inches. Or possibly the researcher is concerned that a four-choice, 100 item multiple choice test score is greater than a chance score of 25. Note that in these two examples (and in most hypotheses regarding a single population mean), the value of zero makes no sense. Suppose that a researcher wanted to establish that the population of freshman at a particular University had a mean College Board Score above the national average of 450:

Research Hypothesis 7: The population of freshmen at University X has a mean College Board Score greater than the national mean of 450.

Full Model: College Board Scores =  $a_0U + E_{10}$ 

The research hypothesis implies that the sample mean is greater than 450, or  $a_0 > 450$ 

Restriction:  $a_0 = 450$ 

Restricted Model: (College Board Scores) = 450 U +  $E_{11}$ , or (College Board Scores - 450) =  $E_{11}$ 

(See bottom of Note 2 for test of significance and McNeil, 1973 and McNeil, et al., 1975, p 315 for further details.)

The desired mean may be difficult to determine (i.e., it may require some thought or knowledge of the phenomenon under consideration), but no more difficult than justifying the default mean of 0. Indeed, using a mean of 0 in this example makes absolutely no sense at all, and that is why it doesn't appear in the literature.

#### One-way Analysis of Variance

The usual application of the multiple group F test (one-way ANOVA) is:

Research Hypothesis 8: There is at least one difference in the means on Y between the 1 populations.

Full Model:  $Y = a_1G_1 + a_2G_2 + ...a_1G_1 + E_{12}$ The research hypothesis implies that not all the sample means are equal, or that  $a_1$  not equal  $a_2$  not equal  $\ldots a_1$  for at least one pair of means, or  $(a_1-a_1 \text{ not equal } 0; a_2-a_1 \text{ not equal } 0; \dots a_{i-1} - a_i \text{ not}$ equal 0 for at least one pair of means.) Restriction:  $a_1 = a_2 = \dots a_1$ ; or

 $a_1 - a_1 = 0; a_2 - a_1 = 0; \dots a_{j-1} - a_j = 0$ By replacing all the coefficients with a common coefficient,  $a_0$ , we arrive at the following restricted model:

 $Y = a_0G_1 + a_0G_2 + \dots + a_0G_1 + E_{13}$   $Y = a_{00}G_1 + G_2 + \dots + G_{11} + E_{13}$   $Y = a_0U + E_{13}$ Restricted Model: Restricted Model: Restricted Model:

When the F test is significant then the restriction is rejected and the research hypothesis is accepted as tenable. But the research hypothesis just indicates that the i means are not Since most researchers are not satisfied with that all equal. information (confirming that the research hypothesis wasn't very interesting in the first place), most researchers turn to posthoc comparisons to find out where the differences lie. These posthoc comparisons are basically t-test comparisons and are thus like Research Hypothesis 1. (See Williams 1974). The suggestion here is to avoid asking a research hypothesis that you aren't interested in, and to go directly to non-zero research hypotheses that will yield satisfying information.

#### Interaction

Interaction is usually viewed only as a potentially contaminating factor when trying to explain main effects. That is, most researchers hope that there is no interaction so that they can proceed with interpreting main effects. But the interaction research hypothesis may be important in and of itself. Indeed, whenever an F has been computed for the interaction, the interaction research hypothesis has been tested. The usual interaction research hypothesis in a 2x2 design is as follows: Research Hypothesis 9: For a given population, the difference on Y between Treatment 1 and Treatment 2 is not the same on Level 1 as on Level 2.

Full Model:  $Y = a_1(T_1 * L_1) + a_2(T_1 * L_2) + a_3(T_2 * L_1) +$  $a_4(T_2*L_2) + E_{14}$  $m_{2} = (R^{2}_{F} - R^{2}_{R}) / (11_{F} - 11_{R})$ Where  $T_{1} = 1$  if in Treatment 1; 0 otherwise,  $T_2 = 1$  if in Treatment 2; 0 otherwise,  $L_1 = 1$  if in Level 1; 0 otherwise,  $L_2 = 1$  if in Level 2; 0 otherwise,

 $(T_1 * L_1) = 1$  if in Treatment 1 and Level 1, etc.

The research hypothesis implies that the two differences are not the same, and that in the sample  $(a_1 - a_2)$  not equal  $(a_2 - a_4)$ , or  $[(a_1 - a_2) - (a_2 - a_4) \text{ not equal 0}].$ 

Restriction:  $(a_1 - a_3) = (a_2 - a_4)$ , or  $[(a_1 - a_3) - (a_2 - a_4) = 0]$ . By placing the one restriction on the full model, one arrives at the following restricted model (See Note 3 and McNeil, et al., 1975):

Restricted Model:  $Y = b_1T_1 + b_2T_2 + b_3L_1 + b_4L_2 + E_{15}$ 

Acceptance of the non-directional research hypothesis leads. to a non-directional statement. All that can be concluded is that the differences are not the same. Hence we don't even know if the differences are greater at Level 1 or Level 2, let alone the magnitude of the difference of the differences. We have just conducted a non-directional test of interaction; a directional test of interaction is reflected in the following:

Research Hypothesis 10: For a given population, the difference on Y between Treatment 1 and Treatment 2 is greater at Level 1 than at Level 2.

Full Model: 
$$Y = a_1(T_1 * L_1) + a_2(T_1 * L_2) + a_3(T_2 * L_1) + a_4(T_4 * L_4) + E_{44}$$

The research hypothesis implies that the difference between  $T_1$  and T<sub>2</sub> is greater at Level 1 than at Level 2, or in the sample  $(a_1 - a_2)$ higher than  $a_2 - a_4$ ) or  $[(a_1 - a_3) - (a_2 - a_4) > 0]$ .

Restriction:  $(a_1 - a_3) = (a_2 - a_4)$  or  $[(a_1 - a_3) - (a_2 - a_4) = 0]$ Restricted Model:  $Y = b_1T_1 + b_2T_2 + b_3L_1 + b_4L_2 + E_{17}$ 

A significant F for Research Hypothesis 10 provides more insight than would one for Research Hypothesis 9. We know that the differences are greater at Level 1, but again we do not know how much greater. If cost or theory dictate, aay, a difference greater than six before a decision is made, the following Research Hypothesis would be appropriate:

Research Hypothesis 11: For a given population, the difference on Y between Treatment 1 and Treatment 2 is more than 6 units at Level 1 than at Level 2.

 $Y = a_1(T_1 + L_1) + a_2(T_1 + L_2) + a_3(T_2 + L_1) +$ Full Model:  $a_4(T_2*L_2) + E_{10}$ 

The research hypothesis implies that the difference between  $T_1$  and T<sub>2</sub> is greater at Level 1 than at Level 2 by more than 6 units, or in the sample  $(a_1 - a_2)$  higher than  $(a_2 - a_4 + 0)$  or  $[(a_1 - a_2) - (a_2)]$  $- \mathbf{a}_{4} > \mathbf{0}_{1}.$ 

Restriction: 
$$(a_1 - a_2) = (a_2 - a_4) + 6;$$
 or  $(a_1 - a_2) = (a_2 - a_3) > 6$ 

Restricted Model:  $(Y - 6) = b_1T_1 + b_2T_2 + b_3L_1 + b_4L_2 + E_{10}$ Research Hypotheses 9, 10, and 11 all test an interaction question, but in slightly different ways. In all three hypotheses, there are four linearly independent pieces of information in the full model. Forcing the one restriction on the full model results in three linearly independent pieces of information in the restricted model. Notice that the full models in Research Hypotheses 9, 10, and 11 are exactly the same. The number of restrictions is also the same; resulting in the same number of

degrees of freedom. What is different, though, is the nature of the restriction and hence the restricted models are different. The three research hypotheses are all "correct" and equally "valid" -they just test three different hypotheses. Research Hypothesis 11, though, provides a more definitive conclusion, because as in the previous examples, a non-zero restriction was made.

Note 1. All the Research Hypotheses in this paper (except the oneway ANOVA) are directional Research Hypotheses. This follows the author's contention that a directional Research Hypothesis provides conclusive information whereas a non-directional Research Hypothesis provides no conclusive information. The Full and Restricted models are the same for the directional and nondirectional hypotheses. The non-directional Research Hypothesis allows the researcher to conclude that  $a_1$  does not equal 0, while the directional Research Hypothesis allows the researcher to conclude that  $a_1 > 0$  (McNeil & Beggs, 1971). With reference to the non-zero restriction, of, say 6, the non-directional Research Hypothesis allows the conclusion that  $a_1$  not equal to 6, while the directional Research Hypothesis allows the conclusion that  $a_1 > 6$ . The directional Research Hypothesis allows a more definitive conclusion using the same data and the same degrees of freedom.

The general F test for testing two regression models is Note 2.  $F(m_1, m_2) = (R^2_p - R^2_p) / (11_p - 11_p)$ 

 $(1 - R^2_{\mu}) / (N - 1i_{\mu})$ 

Where:  $R^2_{R} = R^2$  of the full model,  $R^2_{R} = R^2$  of the restricted model,

 $1i_{\rm F}$  = pieces of linearly independent information in the full model,

 $11_R = pieces of linearly independent information in the$ restricted model,

- $m_1 = (11_p 11_p), and$
- $m_{z} = (N 11_{p})$

This test cannot be used when either the restricted model has no predictors, when the criterion variable is different in the two models, or when the Unit vector is not in the restricted Model. In these cases, the F test must rely upon the sum of the squared scores in the error vector, E in both the full model (ESS,) and the restricted model (E88,):

 $F = (ESS_{g} - ESS_{p}) / (11_{p} - 11_{R})$ 

#### $(ESS_{p}) / (N - 11_{E})$

Note 3. The interaction examples all assumed equal N. The concepts still apply to the unequal N situation, although the Note 3. restricted models will be different. (See Williams, 1972.)

n an	SUMMARY	
SIGNIFICANCE TEST	USUAL RESTRICTION	SUGGESTION
Pearson Correlation	zero	non-zero based on theory or cost
difference between two means	zero	non-zero based on theory. or cost
difference between means (one-way f)	only zero	ignore omnibus F, go an with planned comparisons based on theory or cost
interaction	almost always zero	non-zero based on theory or cost
single population Mean	always non-zero	use more often

#### Bibliography

- McNeil, K. A. Testing an hypothesis about a single population mean with multiple linear regression. <u>Multiple Linear</u> <u>Regression Viewpoints</u>. 1973, 4(1), 7-14.
- McNeil, K. A., & Beggs, D. L. Directional hypotheses with the multiple linear regression approach. Paper presented at the meeting of the American Educational Research Association, New York, February 1971.
- McNeil, K. A., Kelly, F. J., & McNeil, J. T. <u>Testing Research</u> <u>Hypotheses Using Multiple Linear Regression.</u> Cardondale: Southern Illinois University Press, 1975.
- Williams, J. D. Two way fixed effects analysis of variance with disproportionate cell frequencies. <u>Multivariate Behavioral</u> <u>Research</u>, 1972, 7, 57-83.

14 14

3 P

Williams, J. D. <u>Regression Analysis In Educational Research</u>. New York: M88 Information Corporation, 1974.

C:CASENON

MULTIPLE LINEAR REGRESSION VIEWPOINTS VOLUME 18, NUMBER 1, FALL 1991

:\*

.

# Considerations, Issues and Comparisons In Variable Selection and Interpretation In Multiple Regression

;·

Susan Tracz, Ric Brown, and Rebecca Koprive California State University, Freeno

The selection of independent variables when utilizing multiple linear regression in a study is an involved and complex process. The availability of a variety of computer programs usually referred to as "stepwise" procedures affords users numerous options about which they often have little understanding. The purpose of this paper, then, is twofold: first, to present the major uses of regression analyses, the advantages and disadvantages of selection procedures and some caveats for researchers and those who teach statistics, and secondly, to present, compare and contrast several variable selection techniques using two data set.

Huberty (1989) suggests that the concept of variable selection may have some worth in terms of parsimony, explaining relationships, lowering the cost of data collection, and, sometimes, parameter estimation. Variable selection procedures called stepwise procedures are available on all the major statistical computing packages including SAS, SPSS, and BMDP. Even novice researchers can easily run numerous stepwise 1000 Huberty (1989), however, continues by saying that procedures. RX. stepwise analyses have been basically used for three purposes: selection and deletion of variables, 2) assessing relative variable importance, and 3) a combination of selection and variable ordering.

1 1. 30 %

Given this information, it is not surprising to find numerous articles in the literature and theses and dissertations in university libraries that have used and misused stepwise procedures despite the many published caveats concerning its appropriateness. Perhaps one reason for the frequent misuse of stepwise procedures is the mistaken perception that the results of a stepwise procedure will yield the "best" equation. According to Hocking (1983), "there is not likely to be a best equation in multiple regression" (p. 226). This is because the use of differing criteria may result in the selection of different sets of variables (Draper & Smith, 1981). Pedhazur (1982) more specifically stated that such methods as all possible regressions, forward selection, backward elimination, stepwise S. 63 7 selection and blockwise selection can be utilized with differing 乙四南 criteria which will result in differing solutions depending on a the criteria. Morris (1989) sums up these ideas by saying that . Ant "there is little theoretical justification for expecting any stepwise procedure to be best" (p. 2).

The goal of stepwise regression is to choose a subset of variables from a larger set for the purpose of parsimony, prediction, explanation, and/or theory-building. However, since 1. 36.00 the criteria used in selecting variables are statistical, measurement error or randomness may lead to the selection of one variable instead of an equally viable alternative variable. Cohen and Cohen (1975) expounded on this issue saying that "problems include capitalization on chance because of simultaneous tests, sample specificity and trivial differences in partial relationships leading to choosing one variable over another" (p. 103).

When predictor variables are intercorrelated, "there is no satisfactory way to determine relative contributions of the variables on R squared" (Edwards, 1984, p. 107) and "the idea of independent contribution to variance has no meaning" (Darlington, 1968, p. 169). Huberty (1989) reiterates these points by noting that various subsets of a given size can yield nearly the same  $R^2$ value. Pedhazur (1982) states that the R' in variance partitioning is sample specific and that nearly identical regression equations can have radically different  $R^2$  values. Furthermore, an incremental  $R^2$  may be statistically significant but substantially meaningless. Pedhazur (1982) argues that the incremental partitioning of variance may be used to control one variable while studying another variable only in causal modeling, and even then the results are of limited value in determining policy.

Another problem to be dealt with is the interpretation of the regression coefficients. Huberty (1989) cautions that the order in which a variable is entered into a model should not be used to assess its relative importance. "The interpretation of regression coefficients as indices of effects of independent variables on the dependent variable appeals to researchers because it holds the promise for unraveling complex phenomena. Examination, however, is important because the apparent simplicity is deceptive" (Pedhazur, 1982, p. 221). Pedhazur (1982) warns that the absence of a theoretical model makes the meaningful interpretation of the estimated regression coefficients impossible. The types of specification errors that can occur are numerous including omission of relevant variables, inclusion of irrelevant variables, interactions among variables, and the hierarchy of polynomial terms (Cohen & Cohen, 1975; Pedhazur, 1982; Peixoto, 1990).

When so many caveats against it have been published, the continued wide usage of stepwise procedures is difficult to understand. Variable selection techniques in regression analysis can be discussed in terms of parsimony, prediction, explanation and theory-building, and selection techniques are problematic in all of these areas.

Parsimony involves finding "a smaller set of predictor variables that do an accurate job of predicting, nearly as well as the total set of variables" (Morris, 1984, p. 1). Obviously, parsimony is helpful to researchers who reap benefits in terms of economy of data collection costs and time. However, the Criteria for the selection of the best variables must be weighed on a continuum between internal (parsimony) and external (cross validation) accuracy (Morris, 1984). A prior decision made in the name of parsimony can have a tremendous impact on the results of regression analyses used for prediction, explanation and theory-building.

Pedhazur (1982) states that "for prediction, the goal of regression is to optimize prediction of criteria" (p. 136). The selection of variables for this purpose should account for as much of the variance as possible while balancing practical considerations such as cost and ease of administration. While Morris (1989) finds "particularly 'pernicious' ... a situation with a naive researcher ascribing the best prediction equation from the results of a stepwise program" (p. 1), Pedhazur (1982) argues that "prediction may be accomplished in the absence of theory, but explanation is inconceivable without theory" (p. 174).

The goals of many researchers in terms of explanation have been to identify major variables and determine their relative 1082 importance (Pedhazur, 1982). This suggests that stepwise . Start techniques may be plausible initially. The stepwise programs 312 basically perform a hypothesis formulation function (McNeil, Kelly, & McNeil, 1975). However, "problems arise with the 1.87 stepwise approach, since a great many hypotheses are being tested the resulting best model will most likely be drastically overfit ...... with replication relatively unlikely" (p. 364).

Cohen and Cohen (1975) state that "a research strategy of treating all independent variables simultaneously is most appropriate when no logical or theoretical basis for considering any variable to be prior to any other either causal or relevant in terms of research goals" (pp. 97-98). However, despite this seeming endorsement, they continue by saying "a dim view is taken of stepwise in exploratory research because orderly advance is more likely in the social sciences when researchers use theory to provide hierarchical ordering formed by causal hypotheses rather than computers ordering independent variables" (p. 103).

Given all the problems of sample specificity, interpretation of regression weights, and varying R values, the question arises when is it actually appropriate to use stepwise procedures. Huberty (1989) says that in cases where a large ratio of sample size to variables exists, generalizability of stepwise regression is enhanced, but an external analysis or a cross validation should also be conducted. Thorndike (1978) agrees arguing that "when a fairly large number of predictor variables are available it is advisable to use a stepwise approach, but cross validate" 3 200 (p. 167). Finally, Cohen and Cohen (1975) state that the 1 L distrust of stepwise procedures decreases if: "1) the research goal is predictive not explanatory; 2) N is very large for a ÷. given number of independent variables (40 to 1); and, 3) cross validate" (p. 104). Perhaps, Huberty (1989) offers the best 12 2 advise when he says that "thorough study and sound judgement are suggested for choosing variables at the outset" (p. 62), and that 1 **%**6 "the data analyst should allow the findings at each stage to influence the direction through subsequent stages" (Allen & Cody, 112 cited in Huberty, 1989, p. 65). 1 St.C

The numerous stepwise procedures available in the major statistical computing packages are so easy to execute, however, that users quickly learn to rely on them, and there is a great temptation for researchers, especially novice researchers, to assume that a stepwise procedure will yield the best model which will stand up to the test of cross validation. Again, this is simply not true. Stepwise procedures actually yield many best models depending on the procedure used and the criteria employed, and it is up to the researcher to decide which one to use and why. In short, stepwise procedures are no substitution for thinking and theorizing. This paper, will now present, compare and contrast several variable selection techniques using two data sets. In the first example, the results of various stepwise techniques from the SAS package will be compared. In the second example, the results of several stepwise regressions used to answer various research questions will be compared.

The first example consists of a dummy data set of 30 subjects used for classroom teaching purposes. The dependent variable is graduate grade point average (GPA), and the four independent variables are the Graduate Record Exam Quantitative subscale (GREQ), the Graduate Record Exam Verbal subscale (GREV), the Miller's Analogy Test (MAT), and a faculty rating of graduate student performance [RAT]. :(This data set is available from the authors upon request).

The intercorrelations among these variables and the associated probabilities are presented in Table 1.

Table	2_1	Correlations	and prob	abilities	<u>(N = 30</u>
Vari	bles	GREO	GREV	MAT	RAT
GPA	(r) (p)	.61 .0003	.58 .0008	.60 .0004	.62 .0003
GREQ	(r) (p)		.47 .009	.27 .15	.51
GREV	(r) (p)			.43 .02	.41 .03
MAT	(r) (p)				.52 .003

As can be seen the dependent variable GPA is highly correlated with all of the independent variables. All the independent variables are also highly correlated with each other except for the combination of GREQ and MAT (r = .27) and possibly GREV and RAT (r = .41). Therefore, pairs of unique information have been set up between GREQ and MAT and between GREV and RAT.

Five different analyses were run using this data set. The first was a full model with all four dependent variables using the forced solution, PROG REG. This model was significant (F = 11.13, p <.0001, R<sup>2</sup> = .64, adjusted R<sup>2</sup> = .58). The parameter<sup>25</sup> estimates, t values and probabilities appear in Table 2. In this model the t values for GREQ and MAT are significant, while those for GREV and RAT are not.

Table 2	Results of full mo	<u>odel u</u>	sing the	and the Planner of a tak	
forced sol	ution in PROC REG to	pred	ict GPA ST		11.2585 15 T M
from all i	ndependent variables	<b>1</b>	a to the second s	anda: Brook ol v	
Variable	Parameter Estimate	£	P	strate. bra en	11.7
Intercept	-1.738	-1.83	- 08 - C		1.21 <b>4</b> - 199
GREQ	.004	2.18	.04		
GREV	.002	1.45	.16		に合い。 に合い。
MAT	.021	2.19	.04	an an an Anna a Anna an Anna an	
RAT	.144	1.28	.21		
			- *		10.0

1.

1.25

1.

1.18

1.1

77.00

1.

The next analysis which was performed was a forward selection. This program identifies a subset of variables which will be as efficient as the entire set of variables for the set predicting GPA. In this case, the significance level for entering a variable into the model has been set on the lenient side to .15. The variables were entered into the model in the following order: RAT, GREV, MAT, and GREQ. The R<sup>2</sup> values for each new model and the change in  $R^2$  are presented in Table 3. The  $R^2$  for the full stepwise model is .64, as in the full model, since all the variables were entered into the model.

Table 3 Resulting F from the forward select GPA from allindependent	<u>'s and chang</u> ion method t variables	<u>es in R's</u> to predict
Variable Entered <u>into the Model</u>	R <sup>2</sup>	Change <u>in R<sup>2</sup></u>
RAT	.39	-
GREV	.52	.13
MAT	.57	.05
GREQ	. 64	.07

The third analysis was a backward elimination. The procedure starts with all the variables entered into the model - ¥ 6**.** and then eliminates variables. The significance level for retaining a variable in the model has been set to .05. Again the full model had an R<sup>4</sup> of .64. The variable, RAT, was removed first (R<sup>4</sup> = .62) and then GREV (R<sup>4</sup> = .58), so the best model with Again the GREQ and MAT only included has an  $R^2$  of .58. The results appear 18 in Table 4.

Table 4       Resulting R's and changes in R's         from the backward elimination method to predict         GPA from all independent variables					
Variables Included			Variables Removed	R <sup>2</sup>	Change <u>in R<sup>2</sup></u>
GREQ, GREV	/, MAT,	RAT	-	.64	-
GREQ, GREV	, MAT		RAT	. 62	.02
GREQ, MAT			RAT, GREV	.58	.04

The fourth analysis used the stepwise method. This procedure differs from the forward selection method in that variables entered on earlier steps do not necessarily remain in he model on subsequent steps. After a variable is added, other ariables in the model are inspected to determine if they still roduce a significant F statistic. If the F is not significant, he variable is deleted from the model on that step. For this ase, the significant level for entry into the model was set to 15, and the significance level for remaining in the model was et to .05. The results for this analysis appear in Table 5. he variable, RAT, was entered into the model first ( $R^2 = .39$ ), hen GREV ( $R^2 = .52$ ) and then MAT ( $R^2 = .57$ ). Finally, MAT ( $R^2 =$ 52) was removed from the model because the F value for that ariable was not significant, so the resulting best model ncluded RAT and GREV ( $R^2 = .52$ ).

Table	5 Result	ing R <sup>2</sup> s and c	hanges	in R <sup>z</sup> s			
from_a	<u>from all independent variables</u>						
Step	Variable Entered	Variable <u>Removed</u>	<u>R<sup>2</sup></u>	Change <u>in R<sup>2</sup></u>			
1	RAT	dim	.39				
2	GREV	-	.52	.13			
3	MAT	-	.57	.05			
4	-	MAT	.52	.05			

Finally, the last stepwise procedure used was the maximum  $R^2$  thod. This procedure adds variables that maximizes  $R^4$ . The sults of this procedure are presented in Table 6. This ocedure went through five steps and arrived at a model which cluded all four independent variables ( $R^2 = .64$ ). However, it uld be argued that the best model is determined on the basis of e C(P) statistic. The optimal model is the one for which the P) statistic approaches the number of predictors. In this se, the researcher should stop at step 4 since the C(P) atistic is then equal to 4.63 which is closest to the number of edictor variables or four.

Table method	<u>6 Resulting R<sup>2</sup> and C(P)</u> to predict GPA from all	from t) indepe	ndent_varia	R <sup>2</sup> bles
Step	Variables in the model	R <sup>2</sup>	C(P)	
1	RAT	.39	16.74	
2	GREV, RAT	.52	9.69	
3	GREV, Mat, Rat	.57	7.77	
4	GREQ, GREV, MAT	. 62	4.63	
5	GREQ, GREV, MAT, RAT	. 64	5.00	

65 22 Table 7 presents a summary of the results of all the rable / presents a summary of the results of all the procedures. The full model, forward selection, and maximum R<sup>2</sup> method all include all four predictor variables and give an  $R^2$  of the .64. What is curious is that for the procedures which select only two variables the solutions are quite different. The 94 2 stepwise procedure ends up with RAT and GREV  $(R^2 = .52)$ , while the backward elimination ends up with GREQ and MAT  $(R^2 = .58)$ . The forward, stepwise and maximum R' methods all enter RAT into a maximum R' maximum the model first because this variable has the highest correlation say with GPA (4 = .62). The next variable entered is GREV. The and set correlation between RAT and GREV is .41. In the other "best" two variable solutions the correlation between the two predictors, and we way GREQ and MAT is .27. It is important to note that these are the these lowest two correlations among all the predictor variables. When, variables are highly intercorrelated and one variable is entered into a model first, the next variable entered will add the most unique information, i.e., has the lowest correlation with the first variable. In other words, variables are really entered as pairs (GREQ 4 MAT,  $R^2 = .58$ ; GREV 4 RAT,  $R^2 = .52$ ). Also, in which some situations the procedures, namely forward selection, stepwise, maximum  $R^2$ , did not produce the maximum  $R^2$  for the two variable models even though most users think they do. This is because the algorithms in these procedures don't really check all the possibilities.

<u>Table 7 Comparis</u>	on among the best models of the	<u>[1]</u>
model and stepwis	e results	
		_2

Procedure	Variables in the model	<b>R</b>	
Full model	GREQ, GREV, MAT, RAT	.64	A.S. AND
Forward selection	RAT, GREV, MAT, GREQ	. 64	
Backward elimination	GREQ, MAT	.58	5 <b>63</b> 5
Stepwise procedure	RAT, GREV	.52	
Maximum R <sup>2</sup>	GREQ, GREV, MAT, RAT	.64	1985 - 1985 -

1.20

In light of this information, what advise can be given to researchers using stepwise procedures? First of all, users of computer packages should know the limitations of the procedures they use. Secondly, researchers should always study the correlation matrix before looking at other results. A thorough knowledge of the intercorrelations may lead researchers to force certain variables into their models first.

In the next example, the results of stepwise regressions are used to answer different research questions. In this example, data from 65 first time, post-myocardial infarction and first time, post-coronary bypass patients were used to study attributions, self-efficacy, and outcome expectations as predictors of depression. The dependent variable was a 20 item scale called the Center for Epidemiological Studies - Depression [CES-D]. Attribution was measured by two instruments: a 9 item behavioral attribution scale [BEHATT] measuring the causes of heart disease that an individual can change, such as smoking, drinking, etc., and an 8 item nonbehavioral attribution scale (NONBATT) measuring the causes of heart disease that are less controllable, such as heredity, luck, etc. The self-efficacy scale [SELFEFF] has 19 items and measures behaviors that individuals have some degree of confidence that they can change. Outcome expectancy 1 [OUTEXP1] was a 19 item scale rating how important patients believe changing particular behaviors are in preventing future heart attacks. Outcome expectancy 2 [OUTEXP2] was a 19 item scale rating the extent of a patient's belief that if behaviors are changed future heart disease will be prevented. A series of four research questions was asked by individual members of a group of researchers and medical practitioners who each advocated a different modelling approach. The data was then analyzed using combinations of forced and stepwise procedures.

In the first analysis, the question was asked whether the set of attribution or the set of self-efficacy and outcome expectation yielded the largest  $\mathbb{R}^2$ . The results of this analysis consisting of two regression models which entered all variables simultaneously appears in Table 8. These two regression models produce very similar  $\mathbb{R}^2$  values (.26 for the attribution variables and .32 for the self-efficacy and outcome expectation variables), and the weights for four of the five variables were significant. In general, it was found that individuals were less depressed about their heart condition if they believed they had some control in the matter.

Variable Regression Model 1 OUTEXP2		E 15.6*	alara a <b>sai</b> th Taite a saithean Taite a saithean
Selfeff Outexp1	26 12	<b>4.6</b> * 1.0	
R <sup>*</sup> = .32 <u>Regression Model 2</u> BEHATT	- 37	9.6*	
NONBATT $R^2 = .28$	.31	6.5*	

\*\*\*\* £

1 N 1 **1 1** 

0.0**6** 

1. 19 6. 19 6 - 19 6 6

1. A. A. A.

" . OG

3 . . . .

· · · · ·

10.23

\*p < .05

In the second analysis, the question was asked which set of variables explains the most variance after one set was already forced into the model. When the self-efficacy and outcome expectation variables were entered into the model first, the  $R^2$  was .32. After the attribution measures were added the  $R^2$  increased by .08 to .40. When the attribution measures were forced into the model first, the  $R^2$  was .28. After the self-efficacy and outcome expectation variables were added, the  $R^2$  increased by .12 to .40. The results of both analyses were fairly similar.

The third analysis was a forward stepwise regression using all five independent variables. These results appear in Table 9. In this case, the two behavioral attributions added significantly to outcome expectancy 2 in predicting depression.

Lysis 1 - Rest stapying regr	<u>ilting R<sup>-</sup>s and changes</u>
	<u>Change in R<sup>2</sup></u>
.19	.19*
.31	.12*
.37	.06*
.40	.03
.40	.0,0
	tepvis 3 - Rest stepvise tegr ith all inder R <sup>2</sup> .19 .31 .37 .40 .40

\* p< .05

The fourth analysis took a more theoretical approach. Some theory suggests that attributions precede behaviors. Following this reasoning two analyses were performed. For the first model, the behavioral attribution variable was forced into the model followed by the stepwise addition of the self-efficacy and outcome expectation variables. For the second model, the nonbehavioral attribution scale was forced into the model followed by the stepwise addition of the self-efficacy and outcome expectation variables. The results appear in Table 10. Only the significant additions of the stepwise procedures are reported. In both cases, outcome expectancy 2 was the only significant contribution to the attribution variable in oredicting depression. Again the resulting R<sup>2</sup> values (.31 and .27) from these two models are quite similar.

Table 10 Analysis 4 Two combinations of forced attribution and stepwise outcome expectancy/self- efficacy regression analyses to predict depression				
Variables Regression_Model_1	<u>R<sup>2</sup></u>	<u>Change in R<sup>2</sup></u>		
BEHATT	.18	.18		
OUTEXP2	.31	<b>.13</b>		
Regression Model 2				
NONBATT	.14	.14		
OUTEXP2	.27	.13		

In summary, although one could argue in favor of each of hese four analyses, the last analysis seems most reasonable ince it was based on theory. This example does show, once gain, that the research question must dictate the research ethodology.

It is hoped that researchers will realize that although ultiple linear regression is a powerful and flexible statistical echnique and although stepwise computer procedures are otentially useful and facilitative, using these techniques and rocedures to meaningfully explain data is a complex process.

For non-experimental research, it is difficult if not impossible o untangle the effects of various variables. Sound thinking, neoretical framework and understanding of the analytical methods re necessary to avoid illogical or unwarranted conclusions" Pedhazur, 1982, p. 175). "Any meaningful analysis applied to omplex problems is never routine. The clarifying of ontroversies in social science research will not be enhanced by oplying all sorts of techniques" (Pedhazur, 1982, p. 171).

### References

State and the second

(Amore)

- 12 S

Cohen, J., & Cohen, P. (1975). <u>Applied multiple regression/</u> <u>correlation analysis for the behavioral sciences</u> . Hillsdale, New Jersey: John Wiley & Sons.	nonbraat followst outscos foly the
Darlington, R. B. (1968). Multiple regression in psychologi research and practice. <u>Psychological Bulletin</u> , <u>69</u> (3), 182.	.cal 161-
Draper, N., & Smith, H. (1981). Applied regression analysis Ed.). New York: Wiley.	(2nd
Edwards, A. L. (1984). <u>An introduction to linear regression</u> <u>correlation</u> (2nd Ed.). New York: W. H. Freeman and Company.	<b>inand</b> in the second <b>inand</b> in the second single states in the second single states in the second second second second second second second second second second second second second second second second second
Hocking, R. R. (1976). The analysis and selection of variab in linear regression. <u>Biometrics, 32</u> , 1-49.	) <b>les</b>
Huberty, C. J. (1989). Problems with stepwise methods be alternatives. <u>Advances in Social Science Methodology</u> , 43-70.	tter 1,
McNeil, K. A., Kelly, F.J., & McNeil, J.T. (1975). <u>Testing</u> <u>research hypotheses using multiple linear regression.</u> Carbondale, IL: Southern Illinois University Press.	n už ti stati ti stati
Morris, J. D. (1989). <u>Alternative variable selection strate</u> <u>in classification problems</u> . Paper presented at the mee of the American Educational Research Association, San Francisco.	gies ting
Nie, N. H., Huss, C. H., Jenkins, J. G., Steinbrenner, K., G Bent, D. H. (1975). <u>Statistical Dackage for social sci</u> (2nd Ed.). New York: McGraw Hill.	ences 🤐
Pedhazur, E. J. (1982). Multiple regression in behavioral research. (2nd Ed.) New York: Holt, Rinehart and Wins	iton.
Peixoto, J. (1990). A property of well-formulated polynomia regression models. <u>American Statistician</u> , <u>44</u> (1), 26-30	
Thorndike, R. M. (1978). <u>Correlational procedures for resea</u> New York: Gardner Press, Inc.	<u>rch.</u>

• •

## Case Influence Statistics Available in SAS Version 5

John T. Pohiman Southern Illinois University, Carbondale

#### Abstract

Case influence statistics are a useful diagnostic tool for identifying high leverage cases in a sample. A case's influence on a solved regression model depends on that case's residual and its location in the distribution of the predictor variables. Cases with large residuals and located in extreme ranges of the predictor variables' distributions will be most influential. Case influence is illustated with an SAS analysis of a simple data set.

The REG program in version 5 of the Statistical Analysis System (SAS) provides a collection of case influence statistics described by Belsley, Kuh and Welsch (1980), and Freund and Littell (1986). Influence statistics are designed to aid in the detection of cases which are highly influential in the estimation of the regression coefficients. A case's influence on the regression solution is estimated by deleting that case from the sample and recomputing the coefficients. If the coefficients change considerably upon deleting a case, that case is deemed influential. Generally, cases which have large residuals and are in extreme ranges of the predictor variables' distributions will be most influential.

Figure 1 presents a scatter diagram which illustrates case influence for a simple linear regression model in which a dependent variable (Y) is regressed on one predictor (X). The ten data points denoted with the symbol ( $\bullet$ ) yield the regression equation

' = 1 + 1X.

S. SHEW

- 1 fr

G. Ma

1.512

at all the second

A. AMBLES MORE

The ten data points denoted with the letters A to J are then used, one at a time, to augment the original sample of ten observations. Ten augmented samples of size 11 are thus created. The first augmented sample is composed of the 10 original data points plus point A. The second augmented sample consists of the 10 original observation plus point B, and so on to the tenth augmented sample using case J along with the original observations. The influence of the ten lettered data points is determined by comparing the regression coefficients obtained when the lettered data point is included in the analysis with the coefficients obtained after deleting that data point. Table 1 shows the results of this analysis.

Insert Figure 1 About Here

The second and third columns in Table 1 contain the regression coefficients obtained when cases A to J augment the original sample of 10 cases. The last two columns of the table show the change in the regression coefficients due to the presence of each lettered case. Note that the largest change in the slope coefficient occurs for cases F and J. Cases F and J have the largest deleted residuals and are the most disparate cases in the distribution of X. Cases F and J are the most influential cases. Case J has a strong positive influence on the slope coefficient to be .231 units higher than it would be if case J were not in the sample. Case F, to the contrary, has an identically strong negative influence on the slope coefficient.

Insert Table 1 About Here

### INFLUENCE STATISTICS AVAILABLE IN PROC REG

- .

The influence statistics described here are available in the SAS REG procedure as options. SAS provides the statistics HAT DIAG H, DFBETA and DFFITS. For this illustration assume that the general linear model is fit to a data set, namely

where Y is a vector of values on the response variable, X is an nx(p+1)

matrix of values on the independent variables with a leading unit vector, B is the vector of regression coefficients and E is a residual vector. Letting XT denote the transpose of X, the ordinary least squares regression coefficients are given by

Β = (**ΧΤΧ)<sup>-</sup> <sup>1</sup>ΧΤΥ,** 

(4) A 1 (4.1)

and the predicted values of Y are produced by

- X(XTX)- İXT Letting H = X(XTX)- İXT, then

The matrix H is the projection matrix for the predictor space in that it operates on Y to yield Y', and is termed the <u>hat matrix</u>. H is of order n×n and of the same rank as X. The main diagonal values of H, h<sub>11</sub>, are measures of the dispersion of case 1 from the centroid of the predictor variable space. Two cases with the same value of h<sub>11</sub> are on the same probability contour of the multivariate distribution of the predictor variables. In fact, h<sub>11</sub> is a linear transformation of the Mahalanobis distance of case 1 from the centroid of X (Weisberg, 1980, p. 105). The h<sub>11</sub> values measure
the potential for a case to be influential. The actual influence exerted by a case will also depend on that case's residual.

The DFBETA statistics are measures of the influence each case has on each of the regression coefficients. For each case the will be a separate DFBETA value for each regression coefficient in the model, including the intercept. The DFBETA for case I on coefficient j is

 $DFBETAj(1) = \frac{b_j - b_j(1)}{[S^2(1) (X^T X)^{11}]^{1/2}}$ 

:

where by is the regression coefficient for predictor j estimated from the total sample, bj(i) is the regression coefficient for variable j estimated in the sample with case i deleted,  $S^2(i)$  is the error variance estimate from the sample with case i deleted and  $(x^Tx)^{11}$  is the i-th diagonal element of  $(x^Tx)^{-1}$ .

The DFFITS statistic is a scaled measure of the influence of case i on the predicted value of Y. Since all of the regression coefficients are used to produce a predicted Y value, DFFITS becomes an aggregate measure of the influence of case i on the entire regression equation. The DFFITS statistic for case i is given by

$$Y'_{i} - Y'_{i}(i)$$
  
= \_\_\_\_\_\_\_\_\_  
[S2(i) bii]1/2

DFFITS(1)

••

where Y'i is the predicted Y for case I based on the total sample, Y'i(i) is the predicted Y based on the regression equation estimated without case I in the sample, and hill is the 1-th diagonal value of H. The DFFITS statistic is very similar to Cook's D (Cook, 1979), another measure of influence available in the REG program and also in the SPSSx regression program. Cases with DFFITS values greater than  $2[(p+1)/n]^{1/2}$  are considered to be high leverage cases (Belsley et al., 1980, p. 28).

## ILLUSTRATION WITH A DATA SET

Appendix A provides a SASLOG and LISTING for a sample regression model based on 24 cases. Page 1 in Appendix A contains the model statement (SASLOG line 30) which requests the regression of attitudes toward school (ATTSCH) on INCOME and IQ. The INFLUENCE option is requested for the model.

1111

125

Page 2 in the Appendix contains the parameter estimates for the model, followed by the influence statistics. The studentized residuals (RSTUDENT) and the HAT DIAG H present the two important sources of case influence. Case 6 has the largest studentized residual (2.9823) and case 14 also has a large studentized residual (-1.5497). The DFFITS value for case 14 is (-1.5747), and this is the largest value, in absolute terms,

in the sample. The negative value of DFFITS for case 14 means that the predicted Y for case 14 is increased when case 14 is deleted from the sample. Conversely, the presence of case 14 in the sample causes that case's predicted value to be reduced.

The DFBETA statistics are then presented for each regression coefficient, for each case. Case 14 is also the most influential case for estimating each of the regression parameters individually: INTERCEP DFBETA = -.5455, INCOME DFBETA = -1.4997 and IQ DFBETA = .9250. As with the DFFITS statistic, the sign of the DFBETAs indicate the direction of influence on the regression coefficients for case 14. Case 14's presence in the sample causes the y-intercept to decrease, the regression coefficient for INCOME to decrease and the coefficient for IQ to increase. On page 5 of the Appendix the regression equation is estimated with case 14 deleted from the sample, and indeed the changes in the coefficients are as suggested by the DFBETA diagnostics for case 14.

# HANDLING INFLUENTIAL CASES

Once the influential cases have been identified the analyst must decide what to do with them. The first step should be to determine if the influential cases are correctly coded. Typographical errors made while entering the data can produce highly influential cases. If data errors are detected, clearly the proper course of action is to correct the data values. If the correct data values are not available then deletion of such

cases is reasonable.

建建 化水合合

However, if the analyst determines that a case is correctly coded and still highly influential, three alternatives are available: 1. delete the case from the sample, 2. retain the case in the sample but note that the case is influential, or 3. revise the model to accommodate the influential case.

in the state

28.48

3000

1423

It is a questionable practice to delete cases from a sample simply because they are unusual. In fact, unusual cases often point to weaknesses in our models and may suggest improvements in our theories. For example, if a researcher fit a linear model to a nonlinear relationship many of the data points would be found to have large residuals and therefore might be highly influential. Deletion of unusual cases in this example would lead to the interpretation of an incorrect model. When a case is deleted from a sample it is presumed that the model is correct and the offending case is invalid. Our models should be burdened to fit our data; our data should not be obliged to fit our models. Data should not be deleted to better fit our models unless we have compelling evidence that the data is wrong.

The least squares criterion can itself be the cause of an influence problem. A case's influence is proportional to the square of its residual when OLS estimation is used. A researcher might try fitting a model using a criterion other than OLS. The SAS version 5 package has a

procedure that fits models using the least absolute value error (PROC LAV). Unfortunately, this procedure is not available in version 6 of SAS. This program minimizes the sum of the absolute deviations from the model, thereby tempering the influence of high residual cases. If the coefficients estimated with OLS and LAV criteria are comparable, the model may be considered sufficiently robust for interpretation. Page 4 in the Appendix shows the LAV solution for the same model estimated earlier using OLS. The only coefficient that is changed markedly is the y-intercept. The coefficients for INCOME and IQ are approximately the same as their OLS counterparts. One might, therefore, conclude that the OLS estimates are fairly robust in this sample.

REFERENCES Beisley, D. A., Kuh, E., and Weisch, R. E. (1980). <u>Regression diagnostics</u>. New York: John Wiley and Sons, Inc.

Cook, R. D. (1979). Influential observations in linear regression. Journal of the American Statistical Association, 74.169-174.

get in the general second

All All All

- e y - 20

医白白白 囊囊囊膜

1000 2

in the second

Freund, R. J. and Littell, R. C. (1986). <u>SAS system for regression 1986</u> edition. Cary, N. C.: SAS Institute Inc.

Weisberg, S. (1980). <u>Applied linear regression</u>, New York: John Wiley and Sons, Inc.



Figure 1. Scatter Diagram Illustating influence

· ·

Case	Regression C	oeficients	Influence of	Case on
		Slope	Intercept	Slope
	•			
٨	1.625	.846	.625	154
B	1.435	.913	.435	087
С	1.182	1.000	.182	.000.
D	.913	1.087	087	.087
Ε	.692	1.154	308	.154
F	1.920	.769	.920	231
6	1.652	.870	.652	130
Η	1.273	1.000	.273	.000
I	.870	1.130	130	.130
J	.538	1.231	462	.231

Table 1. Influence of Cases A-J on Model Coefficients

1

•

Note: The regression equation for the original 10 cases is Y' = 1 + 1X.

### Appendix Page 1

#### **·SASLOG FOR THE INFLUENCE ILLUSTRATION**

1 DATA ONE; 2 OPTIONS LS = 70 NUMBER; 3 INPUT SUBID GENOER IQ NERLTH GRADE INCOME ATTECH; NOTE: DATA SET NORK. ONE HAS 24 OBSERVATIONS AND 7 UNAIABLES. NOTE: THE DATA STATEDENT USED 0.07 SECONDS AND ONK. 29 PROC RED; 30 NODEL ATTECH = INCOME IQ /INFLIENCE; NUTE: THE PROCEELINE NEG USED 0. 15 SECONDS AND 416K AND PRINTED PROES 1 TO 2. 31 PROC LINU; 22 MODEL MITTECH = INCOME IQ; NOTE: LAN IS NOT SUPPORTED BY THE AUTHOR OR BY SAS INSTITUTE INC. NOTE: THE PROCEEDINE LINU USED 0.16 SECONDS AND SO20K AND PRINTED PAGE 3. 33 DATA THO: 34 SET CHE; 35 IF SUBIÓ HE 14;

.• \*

NOTE: DATA SET HOAK.THD HAS 23 OBSERVATIONS AND 7 UNRIABLES. NOTE: THE DATA STATEMENT USED 0.04 SECONDS AND 424K.

36 PROC RED; 37 MODEL ATTISCH = INCOME IQ; NOTE: THE PROCEDURE RED USED 0. 10 SECONDS AND 440K AND PRINTED PROE 4.

Na se a catego

Appendix Page 2

\* S.V.

innly with			7 W 7		•	1. 1. 1.			
	5		aun of			5	E E		
HODE. ENCON	85°			× B	62.01459 4430664	<b>4</b> 3.	0	1000.	
	Ĕ	222	044600				<b>2</b> 8	Υ. Υ	
			-						
	8	٤w	ST LINE						
INTERCEP INCOME IQ		9-0	1722.001 h		8.23978049 0.16067102 0.0823277	• • • • • • • • • • • • • • • • • • •		0.9001	
٥				, –		8			
8	ē	ł		<b>E</b> :	Ξ	MTIO			
- 4 6 4	• <b>? ?</b> ?		588 -97	8831		282.0			
* 10 0	179		N Q N	6 8 A					
	<b>PP</b>	<b>8</b> 8		88		BIR.	282		
22	i u d			287					
<b>88</b>	4 1 1	JB		¥5		286	0.000	0.1128	
121	<b>†?</b> 7		<b>7</b> 93	582	0.2000		1267.0		
223	- 7	<b>B</b> Ę		588					
28	1		99 9	\$\$2	6.0	533	0.135		

٠

8

8.0

8

.

: 1		٦.				
			8			

.

٠

•

	RESIDURL	ASTUDIT	HAT DIAG H	COU PATIO	DFFITS	I ATTERCEP OFTEETAB
21	7.3604	1, 1207	0.2520	1.2909	0 6522	-0 2000
22	-1.5951	-0.2199	0.1402	1.2267	-0.0000	
23	-3.6772	-0.5110	0.0495	1.1714	-0, 1166	-0.0628
24	-1.1594	-0.1529	0.0605	1.2279	-0.0300	0.0140

.

086		INCOME OPTIETNS	IQ DFBETNS
	1	-0.2443	0.2573
	2	0.0360	-0.0075
	3	0.2091	-0.0658
	•	-0.0075	0.1724
	5	0.0345	-0.0061
	Đ	0.2725	0.0099
	7	0.2375	-0.0409
		-0.0245	-0.1670
	y 10		-0.1629
		-0.0059	-0.0401
			-0.0121
		0.0354	-0.1049
			-0.300
	19		0,9250
	17 14		-0.1099
			0.077
			0.0001
	20		
	<i>8</i> 0 21		-0.000
	22		0.0027
	2		
	X		
•		J.UIII	

•

# Appendix Page 4

#### LIN REFERSION FREEDURE FOR EPERCENT UNRINGLE ATTSCH

在这时代

As y L

ia, inge Sugare

•

.....

#### 

INTER		-0.23255814
INCOME	۶	1. 13953489
		0.09302325

#### (NOTE: THE COEFFICIENT ESTIMATES ARE UNIQUE.)

RESIDUAL SUM OF RESILLITE VALLES = 120.744 10605 ROLLETED TOTAL SUM OF RESILLITE VALLES = 272.00000000 NUMBER OF DESEMANTIONS IN DATA SET = 24 DEP UNRIABLE: ATTSCH ANALYSIS OF UNRIANCE

RUNE	DF	SUM OF SQUARES	NEAN Scurre	F VALUE	<b>PROB</b> >F
HOOEL EPROR	220	4425.92197 1095.73020	2212.95099 <u>54.7655097</u> 5	40.392	0.0001
C TOTAL	22	5521.65217			
ROOT	MBE	7.401791	A-ED.THE	0.6016	
	<b>IERN</b>	33.56522	ndj n-sq	0.7017	
C U		22.05197			

## PINNETER ESTIMITES

VARIABLE	DF <sup>1</sup>	PROVETER ESTIMATE	STREAD DROA	t for ho: Piareter-o	PROB > ITI
INTERCEP	1	4.03922071	8.45705999	0.479	0.6301
Incore		1.37254205	0.21574806	6.222	0.0001
Iq		0.03666642	0.09724585	0.379	0.7085

· `

# Some Parallels Between Predictive Discriminant Analysis and Multiple Regression

٠.

John D. Morris Florida Atlantic University

Carl J. Huberty University of Georgia

Parallels 1 Marsh

The purpose of this paper is to outline some important similarities in, and differences between, predictive discriminant analysis (DA) and multiple regression (MR). The areas covered, chosen for their importance and need for clarification, are estimates of model accuracy, hypothesis testing, and non-least equares models. Some of the parallels are well known, some are less well known, and some appear to have not yet been considered at all.

It is well known that when 1) only two groups are involved, 2) the two population predictor covariance matrices are assumed equal, and 3) the two prior probabilities of group membership are taken to be equal, the popular "minimum chi-square rule" (Tatsuoka, 1971, p. 218) associated with discriminant analysis (DA) is equivalent to predicting a dichotomous criterion variable via multiple regression (MR) methods and classifying a subject into the group for which the predicted criterion is nearer the actual. An especially enlightening examination of this and some other multivariate techniques from the general perspective of MR is provided by Flury and Riedwyl (1985).

However, a precaution about the equivalence of two-group classification and multiple regression with a dichotomous criterion is appropriate. In a two-group situation, there is one linear discriminant function (LDF) and there are two linear classification functions (LCFs); an LDF and an LCF are simply linear composites of

the predictors. It is true in a two-group context that the regression weights are proportional to the single set of LDP weights. When a linear regression function (LRP) or an LDP is used for classification purposes a cut-off criterion needs to be determined with an LRP it is midway between the two values by which the dichotomous criterion is coded, with an LDP it is midway <u>between</u> the LDP means for the two groups. With the use of LCPs, there is no cutoff per se; rather a unit is classified into the group with which is associated the larger LCP score. It turns out that the respective LCP weight differences are proportional to the corresponding LDP and (therefore) the LRP weights.

Input scores for an LRF, an LDF, and an LCF are typically predictor variable measures. [As stated above, any of the three linear composite types may be used for a two-group classification problem.] It turns out that another, still equivalent, approach to two-group classification may be employed. Here, one uses LDF scores for each unit as input for an LCF; we thus have, in essence, a single predictor score for each unit.

When generalizing from a two-group problem to a k-group problem, it is advisable to forget the LRF and LDF approaches and focus on the LCF approach, with predictor measures as input scores.

#### Estimates of Model Accuracy

Estimation of the cross-validated accuracy of a prediction model offers similarities and differences between MR and DA methods. In

\$ #Free 1.224

Parallels 3

both DA and MR the researcher must decide what type of crossvalidated accuracy is of curpern. For instance, is interest in simply estimating an accuracy index parameter from the associated statistic, that is, estimating the index of accuracy ( $\mathbb{R}^2$  or percent of "hits", respectively) that would obtain in the population from that same index in the sample, or is interest in the accuracy that would obtain on application of sample optimized weights to alternate samples from the same population? The concern in this paper will be with the latter type of accuracy.

As in an estimate of cross-validated  $\mathbb{R}^2$  in MR, a judgment of DA "hit-rate" based on the calibration sample is optimistically biased in reference to application to alternate samples. To estimate a cross-validated result in MR, another decision that must be made is whether interest is in relative accuracy, as manifested in the correlation of Y and  $\hat{Y}$ , or in absolute accuracy, as manifested in the MSE. In either case, several formula estimates are available (see Huberty & Mourad, 1980; Romborn, 1978). It is probable that in most predictive uses of MR in the behavioral sciences, such as in personnal selection, concern is with relative accuracy.

Unlike in MR, the concern in predictive DA is in classification accuracy; this is implicitly a concern of absolute accuracy. A formula estimate for cross-validated hit-rate in the general k-group case has largely eluded methodologists. However, a useful, although complicated, formula estimate for cross-validated hit-rate in the two-group case was derived by McLachlan (1975). According to that

#### and a start of the

#### Parallels 4

estimator, the hit rate, 
$$P_g$$
 for group g, where  $g = 1$  or 2 is:  
 $\hat{P}_g = 1 - P(-D/2) - f(-D/2) \{(p - 1)/(Dn_g) + D[4(4p - 1) - D^2]/(32m) + (p - 1)(p - 2)/(4Dn_g^2) + (p - 1)[-D^3 + 8D(2p + 1) + 16/D]/(64mn_g) + D[3D^6 - 4D^4(24p + 7) + 16D^2(48p^2 - 48p - 53) + 192(-8p + 15)]/(12288m^2)\},$ 

where P is the standard normal distribution function, i.e., P(-D/2)is the area to the "left" of -D/2, f is the standard normal density function, D is the Mahalanobis distance, p is the number of predictor variables,  $n_g$  is the number of subjects in group g, and  $n = n_1+n_2-2$ . While the formula looks formidable, with patience, it is <u>calculable</u> with a hand-held calculator. Moreover, as the last term in the multiplier for f(-D/2) is usually very small, one may choose to ignore it, making the formula even more tractable. If the researcher with an orientation toward MR notas that  $D^2 = R^2 N(N-2)/(1-R)^2 n_1 n_2$ , then the McLachlan estimator of cross-validated hit-rate can be obtained from the  $R^2$  resulting from regressing the dichetomous criterion on the predictors.

One slightly "unnerving" aspect of the McLachlan estimator is that it can yield estimated hit-rates that are <u>larger</u> than those that are estimated from the known positively biased process of reclassifying the calibration sample (Morris & Buberty, 1986; 1987). This is unlike the case in MR where the "shrunken" multiple correlation is necessarily less than the value of the multiple correlation derived from the calibration sample. The explanation for

#### Parallels 5

this apparent paradox between methods is that estimators of the cross-validated multiple correlation are functions of the corresponding calibration sample multiple correlation, and are therefore <u>quaranteed</u> to yield smaller values than the sample value. In this sense, the McLachlan hit-rate estimator is not parallel to the MR formula estimators. While it is an estimator of crossvalidation hit-rate, it is not a function of the calibration sample generated hit-rate. Rather, it is a function of the Nahalambis distance between groups, as well as other variables. That is, it does not simply estimate a parameter from a function of the corresponding statistic as do NR formula estimators.

Contraction of the second

An alternate nonparametric approach to estimating crossvalidated hit-rate, which has a wide following in the DA literature, is the "lasve-on-out" procedure (Buberty, 1984; Buberty & Mourad, 1980; Lachenbruch & Mickey, 1968; Mantaller & Tukey, 1968). In this method, a subject is classified by applying the rule derived from all Se except the one being classified. This process is repeated "roundrobin" for each subject with a count of the overall classification accuracy used to estimate the proce-validated accuracy.

Clearly the same "round-rubin" procedure can be used to estimate either relative or absolute accuracy in the use of MR, and has appeared in that context, with perhaps the earliest reference due to Gollob (1967). In a system intended to select optimal MR predictor variable subsets, Allen (1971) coined the procedure "PRESS," and he appears to be the source most often cited in the MR literature.

The apparent computational difficulties due to the inversion of N matrices can be avoided in both MR and DA by using a matrix identity due to Bartlett (1951). This identity is cited and used explicitly in introducing the technique in the DA context by Lachembruch and Mickey (1968), but was not mentioned by Allen in the first introduction of PRESS (1971) nor in its presentation in a later text (Allen & Cady, 1982, p. 254), although the same identity was implicitly used. Moreover, Allen doesn't cite the DA literature and the parallel application of the PRESS procedure. It appears that this resampling process was "invented" independently in the MR and DA literatures.

#### Full vs Restricted Model Rypothesis Testing

A technique that is well known and widely used by MR researchers is that of hypothesis testing through contrasting full and restricted prediction models. The power of this method, its generality, and its applicability to a <u>very</u> wide arena of theoretical questions in science is no doubt part of the reason for the establishment of the MERSIG within AERA.

The same types of model contrast "explanatory increment" questions can be asked and seem to be at least as much potential interest when the criterion is classification accuracy. However, we know of <u>no</u> examples of this technique being used in the literature. There seems to be no reason not to test the difference in proportion of correct classifications (hit-rate) between full and restricted

models to examine maningful hypotheses, just as is done using the  $\mathbb{R}^2$ in MR. The appropriate test statistic is McNemar's (1947) contrast between correlated proportions. Moreover, as the index, "I", of increase in classification accuracy over chance (see Buberty, 1984, p. 168) is distributed similarly, it becomes apparent that such a test would also be applicable to that statistic.

An example of such a test from a study in which the absequent high school dropout of a sample of 76 children was predicted from data available in fifth grade will now be presented. The six predictor variables were gender, race (two levels), number of elementary schools in which the child had been a student, the number of grades the child had repeated, the family structure (living with at least one natural parent and no other soult, or not), and the child's total number of fifth grade absences. As we have evidence of the relationship between both gender and race and the criterion of high school dropout, the hypothesis to be tested concerned the significance of the increment to classification accuracy afforded by adding the four "non-organismic" variables (number of elementary schools, number of grades repeated, family structure, and the total number of fifth grade absences) to the prediction model containing only gender and race.

Classifying the calibration sample, the proportion of correct classifications for the total model was 75% and for the model including only gender and race it was 63%. A 2x2 table illustrating the number of hits and misses for both models is:

Parallels 8

		All Pres	lictors
:		MISS	HIT
Gender and Race	HIT	9	39
	MISS	10	18

The test statistic, z = 1.73, would typically be considered nonsignificant (P = .08) and therefore offers no evidence that the addition 

accuracies for these two three predictor veriable models (member of elementary schools, runber of repeats, and family structure, 798; meter of elementary schools, meter of repeats, and meter of absences, 790) were each greater than for the total six variable predictor model. Thus, unlike the multiple correlation coefficient in MR, even with remainder "internal" estimates of classification hit-rate, accuracy does not necessarily monotonically increase as one adds predictor variables. A different persective concerning contrasting reduced and full model predictor vertable subsets my therefore be necessary for DA applications.

One may argue, however, that the group Mater entrante of

accuracy should be used in any case . An dilitizarision of the impact cross-validated estimates. wenter and a state water a state of whether the state of the state of the state of the state of the state of the

regression, have received a great deal of internition and the Literature (e.g., Dan Marguen, 1978, Normal District Papel & Lunneborg, 1985/ Rosepter, 1979) manel accession with the (Capabell, 1980; Dipilio, 1976; 1977; 1979; Marthur & Billion (1987). As the Denefit to predict the estimate of main methods a struction of whether the convert is printine or atmolute antipapy the results for DA bend to be a subservict tinger dot and they the

intallels 9

that using a cross-validated estimator main neve is that the leaveone-out estimator for the hit rates involved in the involtes is tested above were 644 for the full six-wardenie and all for the three variable model, with a resulting man statistic of a. 2.45, which is of course significant at the .02 lines! . When any the researcher would most likely done to a different dentitieting the significance of the impression out to place solutional yariolas using

(**196**)

Star I Cont Starts A CONTRACTOR The second second second second second second second second second second second second second second second se and the second second Non Least-aguades presidention for the state land intering sides

parallel to the ones of absolute approxy in the fit and (forsis & Buberty, 1997) anharten presitet was elevened to the set links brder cettain linked dimensioners bounders and anti- anti- anti- antijust as likely to coally Without an intermedication tabant when to use the technique. Radio methods are der dram the minutes that they have been purported to be for gitter are the or the A suggested

nethod for choosing between alternate predictor weighting algorithms, including ridge and least equares, has been advanced for the DA case ... by Morris and Ruberty (1987), and for the ME case by Morris (1995). Computer programs for both analysis types are evaluable.

La Stand Land Call

and the second in the second second second second second second second second second second second second second

# Of Miniland

The second second and the second second second second second second second second second second second second s Allen, D. A. (1971) . The product on all of the line of the or iteration for · Belecising one distance of the state of th Rentischy, Department of Stational Lines ..... Allen, D. A., & Cady, F. B. (1982) And Later Manual date by A Collegatory Delingers, CAL Medianorshi Bartlett. H. C. ("Dille An investe antiput all introme theining in disorininant analysia, Apple and detimate and all state ice, 22,

constant out lanter

- canonical variate analysis, foundation the statistical Bankeny, 20, Jolla ..... in environmental particular by Darlington, R. B. (1978). Rechter ver uner metaligen Parchelogical
- the millerin, 25, 1200-1285. ( Pay the Marshall Later 20 DiPillo, P. J. (1976). The appl station and blas to diamrininant andly the Comparison to the south states of the states of
- apallysis, Constantiate topic ir, and the second strates and the second
- Statistics, <u>A8</u>, 1447-1487.
- disorialnest Supplies and the supplies with the second

- **Thus**lets 11

- Campbell, N. A. (1980). Shrunken consumized in allerationst and
- DiPillo, P. J. (1977). Runther applications of applies the discriminant DiPillo, P. J. (1979), Blased distribution enably in Systuation of
- the optimum produce filling of standing the standing of the standing the
- Flury, B., & Risderil H. (1985). " units, the Durant Amongroup
  - Derression, When And Black Street Street, DR. 20-20-

- Gollob, H. P. (1957, September). Cross-valuate constant and a state size\_coo. Paper presented at the manting of the American Psychological Association, Washington D.C.
- Aberty, C. J. (1984). Issues in the use and interpretation of discriminant analysis. Perchological guilterin, of 155-171.
- Buberty, C. J. & Mourad, S. A. (1980). Estimation in sullimple correlation/prediction. Educational and Print Stories Mansulament, 40, 101-112.
- Lachenbruch, P. A., & Mickey, M. R. (1958). Betimetion of error range in discriminant analysis. Technometrics, 10, 1-44
- McLachlan, G. J. (1975). Confidence intervals for the domainational probabilities of misallocation in disgriminant analysis. Bicmetrics, 31, 161-167.
- McNamar, Q. (1947). Note on the sampling error of the differences between correlated propettions or percentages. Environmentally 12, 153-157.
- Morris, J. D. (1982). Ridge regression and some alternate weighting Conjus: A comment on Derlington. Betterbeiten auf 91, 203-210.
- Morris, J. D. (1983). Stepvine sidge regression: A computational clarification. <u>Probological Buildetin, 91. 960-966.</u>
- Morris, J. D. (1986). Microcoputer selection of a predictor weighting algorithm. Multiple Diseas Manuel Manuel 15 53-68.

# 11 alatiana

Norris, J. D., & Haberty, C. J. (1987). Selecting a two-group

Classification weighting algorithm. Multitation Methevioral Research, In press. Mosteller, P., & Bukey, J. W. (1968). Data analysis, including

statistics. In G. Lindsey, & R. Archaen (Bis.), <u>Bandsock of</u> <u>Bocial permission</u>, <u>Vol 2.</u> <u>Banding Mans.</u>, <u>MARLED, 1998</u> Pagel, M. D., <u>5</u> Lumsboorg, C. R. (1905). <u>Bepletical restlation</u> for ridge pegression. <u>Betabolocical Buldenick</u> 1970 state-bas. Roseboom, W. W. (1975). <u>Betisstion of offerenced Bildenics</u> sailuible oorrelation: <u>A classification</u>. <u>Betabolocical Activity</u> (1976). <u>55</u>, 1349-

Rombon V. M. (1979), Widge sepressions about the bigul - ment? Envirolenter D. D. (Jeffe), GG, 280-289, 200 Sepression (Contract) Tataucka, M. M. (1971), <u>Bulledover Loverand Marcha</u> (Contract), Son Wiley.

and a second and a second and a second 
#### MULTIPLE LINEAR REGRESSION VIEWPOINTS VOLUME 18, NUMBER 1, FALL 1991

•

# The Use of Regression Diagnostics to Improve Model Flt: A Case of Role Strain and Job Stress

Susan Tuli Beyerlein and Michael Martin Beyerlein University of North Texas

#### Abstract

This paper illustrates the importance of using regression diagnostics to improve model fit when using standard multiple regression statistical packages such as SASPC. This study examined the relationship between employee perceptions of their work environments and perceived job stress. The analysis was theory driven rather than exploratory in nature, and was performed using SASPC multiple regression procedures. Variables were coded to reduce possible collinearity. Various regression diagnostics were examined to detect the presence of outliers, influential observations, residual correlation, and collinearity (e.g., VIFs, DFFTIS, the C<sub>p</sub> criterion, HAT (leverage) values, and the Durbin-Watson test). These values, coupled with the various regression procedures yielded a final, best nine-variable model of  $R^2 = .48$ , significantly larger than the initial value of  $R^2 = .27$ . Future research in this area could be strengthened through 1) an examination of the path analytic and LISREL models in the literature that attempt to model indirect effects, 2) possible incorporation of select, higher-order terms from these studies, and 3) utilization of the regression diagnostic procedures outlined in this paper.

#### Inconction

Role conflict and role ambiguity are two stressors that have been linked to various health and physical outcomes. Role conflict involves conflicting task assignments initiated by superiors of equal rank and authority. Role ambiguity concerns the lack of clarity regarding job assignments, work objectives, and others' expectations. Kahn and others (1964) found that men who experience role conflict and role ambiguity on the job exhibit more tension and less job satisfaction than men whose roles are congruent or unambiguous. Research shows that role conflict correlates with a number of other outcomes including poor performance (Liddel & Slocum, 1976), poor peer relationships (French & Caplan, 1972), and turnover (Brief & Aldag, 1976; Hammer & Tosi, 1974). Role ambiguity has been linked to ineffective coping, as well as turnover.

Underwillization and job future ambiguity are two additional job stressors that have been shown to impact perceived job stress (Caplan, Cobb, French, Harrison, & Pinneau, 1980). Underwillization of abilities involves the lack of opportunity on the job to use skills and knowledge acquired in school or from previous experience and training. Job future ambiguity concerns levels of certainty regarding future career plans, opportunities for promotion, future value of current job skills, and future job responsibilities. These four variables, that is, role conflict, role ambiguity, underwillization, and job future ambiguity, plus years on the job and gender were chosen from a larger set of variables because of strong theoretical connections to stress, and after correlation analysis suggested they were the best set for predicting perceived job stress.

The present study involved a survey of staff members at a large, southwestern university. Respondents were white collar workers in various clerical, secretarial and administrative positions. A total of 660, 14 page surveys were sent through the campus mail system, and 134 were returned for a response rate of 20.3 percent. Twenty-three cases were omitted because of missing data. The initial predictor variables used in the study were as follows: gender (D1), years on job (X1), role conflict (X2), role emblguity (X3), underutilization (X4), and job future emblguity (X5). The criterion variable was perceived job stress (Y). Gender (D1) was represented by dummy coding (i.e., 0 males, 1 females). Selected interaction terms were then created based on developed theory in the literature, that is, years on job times underalization (XB), role conflict times role emblguity (X7), and years on job times job future ambiguity (XB). Due to the fact that stress has often shown nonlinear relationships to other variables, several squared, higher order terms were included in the analysis, that is, role conflict (7272), role emblyrity (72373), under till zerian (X4X4), and job future ambiguity (XDCB). Finally, all the predicus variables with the emertim of gender (D1), were coded in order to reduce the likelihood of rounding errors in regression coefficients leading to collinearity (Mendenhall & Sincich, 1989, p. 343). Thus, to denote coded variables, "U" replaces "X" for all variables except D1 and the criterion variable Y.

#### Regults

The analysis was <u>performed</u> using SASPC and involved a number of procedures. First, the initial set of predictors was included in the general regression

procedure, PROC REG (i.e., D1, X1-X5). This analysis yielded an  $R^2 = .27$ . Next, this procedure was repeated with these variables and the additional interaction and higher order terms (i.e., D1, X1-X8, X2X2, X3X3, X4X4, X5X5). This yielded an  $R^2 =$ .34. The correlation procedure, PROC CORR, was also run at this point in order to obtain means and standard deviations for the predictor variables.

The twelve variables (excluding D1 and Y) were then coded and analyzed using the general regression procedure, PROC REG with the INFLUENCE option, (i.e., U1-U8, UZU2, U3U3, U4U4, U5U5). This analysis yielded an  $R^2 = .31$ . The subsequent inclusion of D1 (gender) raised the  $R^2$  value to .34. The DFFITS values were then examined in order to identify possible influential observations. The <u>SAS</u> <u>User's Guide: Statistics</u> (1985) describes the DFFITS statistic as "a scaled measure of the change in the predicted value of the ith observation (which is) calculated by deleting the ith observation" (p. 677). The difference,  $y_1 - y_8$ , has been divided by its standard error so that the differences can be more easily compared. The investigator is interested in values that are considerably larger relative to the other differences in predicted values. For most purposes, a value of 1.0 is considered to be sufficiently large to warrant strentim. In the present study, influence diagnostics revealed five DFFITS values greater than 1.0. These were subsequently deleted from the analysis leaving a remaining sample of n=108.

The regression procedure, PROC STEPWISE, was then utilized, specifically, the FORWARD, BACKWARD, and MAXR options. The PROC STEPWISE procedure is a good choice when there are a number of independent variables to consider. The various options do not always isolate the model with the highest R<sup>2</sup> but rather seek

the best one-variable model, two-variable model, and so forth (SAS User's Guide: Statistics, 1985). The FORWARD option requests the forward selection technique, BACKWARD requests the backward elimination technique, and MAXR requests the maximum  $R^2$  improvement technique. MAXR looks at all possible regression equations, however, as with the other options it outputs only the best models, for example, the best ten-variable, nine-variable, eight-variable models, and so forth.

After examining the output from the PROC STEPWISE analyses it was decided that the following was the best modal:  $\mathbb{R}^2 = .38$ , D1, U1, U2, U3, U4, U8, U3U3, U4U4,  $C_p = 7.87$ , with all variables significant at the 0.10 level. The  $C_p$  criterion is gleaned from the FORWARD and BACKWARD procedures (rather than MAXR) and is used to select the best subset model with a small total mean equare error ( $C_p$ ), and a value of  $C_p$  near p + 1, which indicates that elight or no bias exists  $[E(C_p) \approx p + 1]$ . In this case, the  $C_p$  value was slightly less than the number of parameters in the model ( i.e., eight).

This model was then analyzed using the general regression procedure, PROC REG, with the VIF, P, R, DW, and INFLUENCE options. VIF prints variance inflation factors with the parameter estimates; variance inflation is the reciprocal of tolerance; P calculates predicted values from the estimated model and input data; R analyzes the residual and includes the Cook's D statistic which is an overall measure of influence for each observation, the standard errors of the predicted and residual values, and the studentized residual; DW calculates the Durbin-Watson statistic; INFLUENCE prints the following diagnostics used in the present study for each

observation: the residual, studentized residual, HAT or leverage value (h.), and the

negatively correlated (Mendenhall & Sincich, 1989, p. 307). However, calculation of They were subsequently deleted leaving a final cample of n=101. values were influential observations and abould be aliminated from the data set. revealed four values greater than twice the average value, suggesting that these the average leverage value, h = (k + 1)/n = .17, and examination of the HAT values  $\mathbb{R}^2 = .41$ , and a Durbin-Watson,  $\mathbb{D} = 2.21$ , suggesting the residuals were alightly general regression procedure, PROC REG, using the above best model yielded an This value was subsequently deleted leaving a sample of n=105. A rann of the revealed a value greater than +2 standard deviations, that is, a possible outlier. Emmination of the plot RESID\*PRED (readuals times prediced ecores)

U3, U4, U7, U8, U3U3, U4U4, C, = 8.85, DW = 2.21. Thus, the final model included the following variables: gender (D1), years on job (U1), role conflict (U2), role correlation. Durbin-Watson statistic was close to a value of two, suggesting minor reddual BACKWARD procedure as one dropped to the eight variable models; and 4) the model was chosen as the best model for several reasons; 1) the C, value was only significant at the 0.10 level; 3) there was a significant drop in R<sup>4</sup> using the elightly less then the number of predictore (Younger, 1986), whereas it was and MAXR revealed algorificant gains in R<sup>2</sup> values. At this point, a nine-variable aignificantly larger for other modals with similar R<sup>a</sup> magnitude; 2) all variables were Econdmetion of PROC SIEPWIER options, that is, FORWARD, BACEWARD, Therefore, the best model chosen was as follows: R<sup>2</sup> = .48, D1, U1, U2,

embiguity (U3), underutilization (U4), role conflict times role embiguity (U7), years on job times job future embiguity (UB), and the equared, higher-order terms utilizing role ambiguity (UJU3) and undervalization (U4U4). The only conflicting evidence was the value of the variance inflation factors (VIFs) for U3, U4, U3U3, and U4U4. These values were greater than 10, whereas the VIFs for all other variables in the model were approximately 10 or less. VIFs greater than 10 indicate the presence of collinearity where, (VIF) =  $1/(1-R^2)$ , i = 1, 2, ..., k (Mendenhall & Sincich, 1989. p. 237). Values greater than 10 occurred only in those variables used both singularly and squared in the higher-order terms, making them obvious candidates for collinearity. In addition, Mendenhall and Sincich (1989) discuss the need to code the dependent, as well as the independent variables, in order to properly calculate VIFs (p. 236). The criterion variable, perceived job stress (Y), was not coded in this study. Finally, the  $R^2 = .48$  representing the best model did not appear to be sufficiently large to indicate the presence of collinearity. Consistent with this finding, the standard errors of the individual beta parameters were not inflated, and the t-tests on the individual beta parameters were significant suggesting lack of evidence for collinearity (Mendenhall & Sincich, 1989, p. 236).

#### Disvertion and Recommendations

Of several hundred studies of stress examined by the authors, it appears that none have used the regression diagnostics discussed in this paper, suggesting that the results of models presented in the literature may be weaker than necessary. Results of the present study illustrate that the use of the various regression diagnostics can improve best model fit considerably. In addition, it should be

obvious that investigators cannot depend solely on regression selection options such as MAXR, FORWARD and BACKWARD when searching for the best subset regression model. Options such as MAXR will provide R<sup>3</sup> values for all generated models, however, the final decision as to which is the best model cannot be made without the C, statistic and other values, for <u>example</u>, regression diagnostics such as HAT values, Cook's D, results of jackinifing procedures such as deleted residuals, DFFITS, DFBETAS, and the Durbin-Wetson statistic which are available under the FORWARD and BACKWARD options of the PROC STEPWISE procedures.

The FORWARD and BACKWARD options offer different best models. That is, they each output best models based on the particular programmed criteria embedded in their respective routines, with the R<sup>2</sup> as the salient criterion. However, a strong R<sup>3</sup> value is not inequivocably the last word on model fit. For example, if two models with similar R<sup>2</sup> values are examined, it may be that the model with the slightly lower R<sup>2</sup> will better satisfy the other criteria discussed above and will thus be the better choice overall. Therefore, the investigator needs to utilize the power of these routines coupled with intelligent decision making regarding the various procedures. Coding variables reduces the likelihood of collinearity, and outputting regression disgnostics enables the investigator to experiment with dropping outliers and influential observations to see how their absence affects the variance accounted for by the overall model. In summary, there is nothing automatic about the process. SASPC and other packages will provide the mathematics, but it remains the responsibility of the investigator to expendibly to arrive at truly the best model.
The analysis of the present study was theory-driven rather than exploratory in nature. In other words, because of the authors' preference for confirmancy modalling techniques, a limited number of interaction and higher-order terms were chosen based on the literature. However, the literature is replete with more complex models, that is, path analyses and LISREL models that strempt to model indirect effects. Thus, future analyses could be improved by studying the literature in more depth to arrive at other plausible variables and higher-order terms. Possible variables to be included in additional studies involve two general categories, that is, 1) job design facets such as autonomy, responsibility, feedback, task significance, task wholeness, leadership style; and 2) moderator variables consisting of personality characteristics and other demographics, such as Type-A, locus of control, and growth need strength. In addition, existing studies could be strengthened through replication and unliastion of the regression diagnostics detailed in the present study.

- Brief, A. P., & Aldag, R. J. (1976). Correlates of role indices. Journal of Applied Psychology, 61, 469-472.
- Ceplan, R. D., Cobb, S., French, J. R. P., Harrison, R. V., & Pinneeu, S. R., Jr. (1980). Job demands and worker health. Ann Arbor, MI: Institute for Social Research.
- French, J. R. P., & Caplan, R. D. (1972). Organizational stress and individual stress. In A. J. Marrow (Ed.), <u>Failure of success</u>, New York: AMACOM.
- Hamner, W. C., & Tosi, H. L. (1974). Relationship of role conflict and role embiguity to job involvement measures. Journal of Applied Esychology, 39, 497-499.
- Kahn, R. L., Wolfe, D. M., Ouinn, R. P., Snoek, J. D., & Rosenthal, R. A. (1964). Occupational stress: Studies in role conflict and ambiguity. New York: Wiley.
- Liddell, W. W., & Slocum, J. W. (1976). The effects of individual-role competibility upon group performance: An extension of Schutz's FIRO theory. <u>Academy of</u> <u>Management Journal, 19</u>, 413-426.
- Matteson, M. T., & Ivano-vich, J. M. (1987). <u>Controlling work stress: Mective</u> <u>human resource and management strategies.</u> San Francisco: Jossey-Bass. Mendenhall, W., & Sincich, T. (1989). <u>A second course in Dusiness statistics</u>:

Recreasion analysis (3rd ed.). San Francisco: Dellen Publishing Company. SAS Institute, Inc. (1985). <u>SAS User's Guide: Statistics</u> (5th ed.). Cary, NC: Author. Younger, M. S. (1985). <u>A first course in linear recreasion</u> (2nd ed.). Boston: Durbury Press.

## MULTIPLE LINEAR REGRESSION VIEWPOINTS VOLUME 18, NUMBER 1, FALL 1991

If you are submitting a research article other than notes or comments, I would like to suggest that you use the following format if possible:

Title Author and affiliation Indented abstract (entire manuscript should be single spaced): Introductions (purpose == short review of literature, etc.) Method Results Discussion (conclution): References

All imanuscripts should be cent to the editor at the above address. (All manuscripts should be camera-ready.)

It is the policy of the M.L.R. SIG-multiple linear regression and of *Viewpoints* to consider articles for publication which deal with the theory and the application of multiple linear regression. Manuscripts should be submitted to the editor as original, double-spaced, *camera-ready copy*. Citations, tables, figures and references should conform to the guidelines published in the most recent edition of the *APA Publication Manual* with the exception that figures and tables should be put into the body of the paper. A cost of \$1 per page should be sent with the submitted paper. Reprints are available to the authors from the editor. Reprints should be ordered at the time the paper is submitted, and 20 reprints will cost \$.50 per page of manuscript. Prices may be adjusted as necessary in the future.

A publication of the Multiple Linear Regression Special Interest Group of the American Educational Research Association, *Viewpoints* is published primarily to facilitate communication, authorship, creativity and exchange of ideas among the members of the group and others in the field. As such, it is not sponsored by the American Educational Research Association nor necessarily bound by the association's regulations.

"Membership in the Multiple Linear Regression Special Interest Group is renewed yearly at the time of the American Educational Research Association convention. Membership dues pay for a subscription to the *Viewpoints* and are either individual at a rate of \$5, or institutional (libraries and other agencies) at a rate of \$18. Membership dues and subscription requests should be sent to the executive secretary of the M.L.R. SIG."



The University of Akron is an Equal Education and Employment Institution ©Copyright The University of Akron 1991/991-ED-014