
Multiple Linear Regression Viewpoints

A Publication sponsored by the American Educational
Research Association's Special Interest Group on
Multiple Linear Regression: The General Linear Model

MLRV

Volume 28 • Number 1 • Spring 2002

Table of Contents

On Analyzing Repeated Measures Designs With Both Univariate and Multivariate Methods: A Primer with Examples	1
Kevin M. Kieffer, Saint Leo University James A. Haley - VA Medical Center, Tampa	
A Monte Carlo Simulation Comparing Parameter Estimates from Multiple Linear Regression and Hierarchical Linear Modeling	18
Daniel J. Mundfrom, University of Northern Colorado Mark R. Schultz, University of Northern Colorado	
Scoring Above the International Average: A Logistic Regression Model of the TIMSS Advanced Mathematics Exam	22
James B. Schreiber, Southern Illinois University Carbondale	
A Discussion of an Alternative Method for Modeling Cyclical Phenomena	31
Russell Brown, Summa Health Care Isadore Newman, University of Akron	

Multiple Linear Regression Viewpoints

Editorial Board

T. Mark Beasley, Editor
University of Alabama at Birmingham

Robin K. Henson, Associate Editor
University of North Texas

Wendy Dickinson (1998-2002) University of South Florida
Jeffrey B. Hecht (2001-2005) Northern Illinois University
Robin K. Henson (2001-2004) University of North Texas
Janet K. Holt (2000-2004) Northern Illinois University
Daniel J. Mundfrom (1999-2003) Northern Colorado University
Bruce G. Rogers (2001-2005) University of Northern Iowa
Kenneth Strand (1998-2002) Illinois State University
Dash Weerasinghe (2001-2004) Dallas Independent School District

Multiple Linear Regression Viewpoints (ISSN 0195-7171) is published by the AERA Special Interest Group on Multiple Linear Regression: General Linear Model through the **University of Alabama at Birmingham**.

Subscription and SIG membership information can be obtained from:
Jeffrey B. Hecht, MLR:GLM/SIG Executive Secretary
Department of Educational Technology, Research & Assessment
Northern Illinois University
DeKalb, IL 60115-2854.
jbhecht@niu.edu

MLRV abstracts appear in CIJE, the ERIC system, and microform copies are available from University Microfilms International, 300 North Zeeb Road, Ann Arbor, MI 48106. *MLRV* is listed in the *EBSCO Librarians Handbook*.

On Analyzing Repeated Measures Designs With Both Univariate and Multivariate Methods: A Primer with Examples

Kevin M. Kieffer

Saint Leo University

James A. Haley VA Medical Center, Tampa

The present paper provides an introductory exposure to different approaches currently available for analyzing data generated through repeated measurements of a phenomenon of interest. Even univariate repeated measures data can be analyzed by employing univariate or multivariate data analytic strategies. Both univariate and multivariate methods can be valuable under certain conditions and when various assumptions are met. The present paper examined both univariate and multivariate approaches to analyzing repeated measures data and compared the results of these methods with classical ANOVA and multiple regression analyses. A small heuristic data set was utilized to make the discussion concrete and to facilitate conceptual understanding of the material.

The primary objective of all scientific research is to gather new information about a phenomenon of interest and to convey to other interested parties conclusions in an effort to further the accumulation of scientific knowledge. Since science is concerned primarily with repeatable and replicable experiments, researchers are required to examine and reduce the influence of random effects that might contaminate results and consequently distort conclusions. One manner in which researchers have sought to reduce the influence of random effects is to utilize a separate set of individuals for different levels of a treatment condition (termed a between-subjects design). The logic employed is that any individual differences manifested prior to the implementation of the treatment can be alleviated by randomly assigning individuals to only one of the various conditions. This approach has the advantage of being less cumbersome to participants but frequently places a strain on the researcher, as it is necessary to amass larger numbers of participants for each treatment and control condition so that the random assignment mechanism can work effectively to minimize the influence of outliers (Girden, 1992).

Another way to control for the influence of individual differences is through the use of the same group of individuals measured repeatedly over time, occasions, or conditions (often termed repeated measures or within-subjects designs). In this manner each individual serves as his or her own control group, and the aberrant influences garnered by one individual on a single measurement will likely be manifested by that individual on subsequent measurements as well. Consequently, the influence of different levels of the treatment condition, as well as the influence of intrapersonal (within-subjects) and interpersonal (between subjects) differences, can be examined. Advantages afforded by this approach include a decreased number of participants required to achieve reasonable power against Type II error and distorted effects due to outliers (Stevens, 1996).

The differences between the two methodologies previously described can be illustrated in a brief example. Suppose a pharmaceutical company has developed a new drug for adult depression and desires to test the drug on a group of human participants for any potentially harmful side effects. The company, having just developed the drug treatment and thus nearly exhausting the financial resources reserved for the project, needs to find a cost effective manner in which to test the treatment. One viable alternative, the between-subjects approach, would be to gather a large group of individuals and randomly assign participants to six different dosage conditions ranging from a placebo treatment (sugar pill) to a 450mg treatment of the medication (spaced in 150mg increments). Supposing a minimum number of 20 participants are required per condition, the drug company would need to recruit a minimum of 120 individuals to participate in the study.

Another potential option would invoke a repeated measures design in exploring the differential effects of the drug treatment. In this approach a group of individuals would also be gathered, but each individual would receive each of the six levels of the drug treatment in a random order. Again supposing a minimum of 20 participants per condition, a total of only 20 participants would be required to complete the study. Consequently, the demands on participants is greater in the repeated measures design, but the benefits of employing such a design can be both statistically elegant and satisfying for the researcher.

Because employing repeated measures designs can be both practical and powerful in social science research (Girden, 1992), it is lamentable that more researchers have not opted to utilize this type of design. One reason for the under utilization of this approach is disagreement regarding the appropriate analytic technique to use to evaluate results, as either a univariate or multivariate approach can be invoked (Algina & Keselman, 1997; Girden, 1992; Keselman, Keselman, & Lix, 1995; Keselman, Lix, & Keselman, 1996; Maxwell & Delaney, 1990). As noted by several authors (Algina & Keselman; Edwards, 1985; Girden; Maxwell & Delaney), each analysis has distinct advantages and disadvantages, and each type of analysis will provide a more powerful result under certain conditions and when certain statistical assumptions are satisfied.

The purpose of the present paper is to examine the similarities and differences in the univariate and multivariate analysis of univariate repeated measures data involving repeated measurements of a single dependent variable. The discussion will include comparisons with classical ANOVA and multiple regression approaches to repeated measures analysis. Power and practicality of the various approaches will be illustrated by utilizing a small heuristic data set throughout the paper.

Univariate Approach to Repeated Measures Analysis

The univariate approach to repeated measures analysis, often termed a repeated measures ANOVA, invokes many concepts also employed in a classical ANOVA analysis in that the sums of squares (SS) is partitioned into various constituent components. In repeated measures ANOVA, however, it is possible to further partition variation that classical ANOVA simply terms ‘error’ variance. Thus, in an effort to understand the advantages afforded by invoking repeated measures designs, it is first necessary to review some pertinent components of classical ANOVA.

Mechanics and Logic Underlying Classical ANOVA

ANOVA (Fisher, 1925) is a statistical procedure utilized to compare means ($k \geq 2$) in an effort to determine if the means differ from one another (Edwards, 1979). In the most traditional application of ANOVA, researchers typically utilize the technique to test the statistical significance of the differences among or between groups of means, but ANOVA can also be used to generate variance-accounted-for effect sizes as well (Wilcox, 1987). As stated by Shavelson (1988), “The purpose of... ANOVA is to compare the means of two or more groups in order to decide whether the observed differences between them represent a chance occurrence or a systematic effect” (p. 342).

Classical ANOVA with One Between Factor

A construct crucial to understanding ANOVA is the most basic unit of all statistical analysis: variance. As stated by Haase and Thompson (1992), “variance is the ‘stuff’ on which all analysis is based” (p. 3). Variance is the degree to which the scores are spread out or dissimilar in a data set. As the name implies, analysis of variance is concerned primarily with determining the sources of variation present within a set of data.

The most simplistic manner in which to examine the shared variance or the differences manifested in group means is by examining distinct groups on a singular independent variable, commonly referred to as a one-way ANOVA. In this context the word “way” (also termed factor) refers to the singular independent variable under investigation (it is assumed in ANOVA that only one dependent variable will be utilized). Researchers employing the one-way ANOVA technique typically have a group of individuals, which has been either naturally or systematically divided into two or more subgroups on a singular classification scheme. The objective of using the one-way ANOVA is to determine if the groups under investigation are appreciably different from one another by ascertaining what proportion of the total dependent variable variance is accounted for by each group.

The process of executing a one-way ANOVA involves consulting several statistics utilized in the computational process. Among the first statistics derived in this process is (a) the sum of squares total (SS_T) - the total variation in the scores on the dependent variable; (b) the sum of squares between groups (SS_B) - the variation in the scores on the dependent variable attributed to the independent variable; and (c) the sum of squares within groups (SS_w) - the variation in the scores on the dependent variable that cannot be accounted for by the independent variable. An interesting property of the three SS components is that in a balanced design SS_B and SS_w always sum to SS_T , as SS_B and SS_w are partitioned areas of SS_T . A useful

metaphor in conceptually understanding SS is to equate SS_T with an entire pie. The pie is constituted of two individual slices, one representing SS_B and one representing SS_W . Thus, SS_B and SS_W are always perfectly uncorrelated and are completely separate and unique entities. The size of each slice of pie corresponds to the proportion of variance accounted for by each SS component.

ANOVA Test Statistic and Decision Process

The test statistic used in ANOVA to determine if differences are present in the analyzed data is called the F_{ratio} or $F_{\text{calculated}}$ (F_{calc}). This statistic is generated by dividing the mean square between (SS_B divided by the degrees of freedom between (df_B)) by the mean square within (SS_W divided by degrees of freedom within (df_W)). The resultant F_{calc} can be compared to the critical values of the \mathcal{F} -distribution (F_{crit}) contained in most statistical texts, and if the F_{calc} supersedes the F_{crit} value (located by using the degrees of freedom between and the degrees of freedom within) at the specified α level, then the result of the analysis is said to be “statistically significant.” A result that is statistically significant implies that given the sample size and the specified α level, the researcher has decided to reject the null hypothesis. A statistically significant result, however, does not imply that the result is inherently important (Thompson, 1993; 1996a; 1999a; 1999b). Conversely, an F_{calc} that fails to exceed the F_{crit} value is subsequently considered not statistically significant but is not necessarily regarded as inherently unimportant, as statistical significance does not furnish any information to a researcher in regard to result importance (Thompson, 1999a, 1999b).

Another manner in which to render a decision in ANOVA is through the computation of a statistic called eta^2 , a variance accounted for effect size functionally equivalent to R^2 . Eta^2 is an uncorrected effect size (see Snyder & Lawson, 1993 for a full treatment of effect sizes) and is computed by dividing SS_B by SS_T . The resultant statistic informs the researcher as to what portion of the variance can be explained with knowledge of to which group the participants belonged.

Multiple Between-Subjects Factors in Classical ANOVA

In a two-way ANOVA, two independent variables are examined in relation to a singular dependent variable. By employing this analysis, it is possible to examine the effects of one or both of the ways on the dependent variable as well as the effects of combinations of the ways. By affording researchers the ability to examine the effects of a combination of the ways on the dependent variable, ANOVA becomes a very powerful tool in discerning the relationships between different variables.

As stated previously, ANOVA partitions the variance on the dependent variable into two distinct components, SS_B and SS_W . Recall from the earlier discussion of the one-way case that each partition is a unique and separate portion of the SS_T . In the two-way case, the SS_B portion is further partitioned into two components, the SS associated with the main effects and the SS associated with the interaction. For a balanced design, each portion of SS_W is still considered a separate and unique portion of the SS_T and will *never* overlap with any other SS partition. A final SS partition splits the SS allotted to the main effects into one portion associated with the “A” way main effect and one portion associated with the “B” way main effect. Thus, in the factorial two-way case, there are four SS partitions resulting in the components of $SS_{A\text{-way}}$, $SS_{B\text{-way}}$, $SS_{\text{INTERACTION}}$, and SS_W . Consequently, in a balanced design, any information generated by one SS component does not render any information about another SS component, as all the SS partitions are *perfectly uncorrelated* (i.e., orthogonal). The only exception to this dynamic is when one partition equals SS_T , because the other partitions must then equal zero.

The resultant effect of performing a multi-way classical ANOVA analysis can be either positive or negative (Benton, 1991). Two dynamics are at work when discerning the effect of multiple ways on a given analysis: the extent to which the SS_W and df_W are reduced. Reducing the SS_W without altering the associated df_W reduces the MS_W component and provides a larger F_{calc} value. Unfortunately, partitioning more variance in classical ANOVA costs the researcher, and the remuneration for partitioning more variance is a reduction in the df_W . The effect of reducing the df_W can counteract a reduction of the SS_W , as reducing the df_W will increase the MS_W and will subsequently reduce the F_{calc} value. Thus, as in any analysis, there is a cost-benefits decision to be rendered when including more than one factor in a classical ANOVA. If the newly included factor consumes a substantial portion of SS_W without appreciably reducing the df_W , the resultant F_{calc} value is more likely to be statistically significant.

Single Within-Subjects Factor Designs

Invoking a single within-subjects factor (e.g., repeated measures) design allows the researcher to partition variance on the dependent variable in much the same way as in classical ANOVA. The difference is that in a single factor repeated measures design, the total variation is partitioned into that component associated with variation due to individual differences in the participants (SS_{SUB}) and variation due to differences in the levels of the treatment condition (SS_{TRT}). The remaining variation (typically termed SS_{RES}) is often considered an interaction component because it represents variation due to the unique combination of the participants and the treatment levels. As noted previously, in classical ANOVA with a single between-subjects factor, it was only possible to partition variance into between-groups (variation due to differences in groups of participants) and within-groups variation (error). Because repeated measures analysis explains more of what classical ANOVA simply terms error, it is possible to reduce the MS_{RES} value and subsequently generate a larger F_{calc} value with fewer participants (Girden, 1992). Consequently, repeated measures ANOVA designs tend to be more efficient because fewer participants are required to conduct an analysis that may generate more statistical power than classical ANOVA (i.e., between-subjects) designs.

Classical ANOVA includes individual differences in performance as error (SS_w) because with only one measurement of each person the variance attributable to each person as a unique individual cannot be estimated. In repeated measures analyses, however, the variance due to differences in people can be estimated. Thus, in a single within-subjects factor situation, variance is partitioned into between-subjects (differences across the different treatment effects), within-subjects (differences due to intrapersonal variation), and residual variation.

Carry-Over, Latency Effects and Counterbalancing

As stated previously, repeated measures designs require the same group of participants to be measured on multiple occasions. It has been recognized in the social sciences that the order of presentation of stimuli can sometimes have a differential effect on participants' responses as some experiments involve repeated exposure to the same task (Girden, 1992; Keppel, 1991; Keppel & Saufley, 1980; Keppel & Zedeck, 1989; Stevens, 1996). As such, carry-over and latency effects are common problems in repeated measures research. A latency effect refers to a situation in which the effect of a treatment is not evident until a subsequent level of the treatment is introduced. A latency effect may predispose a researcher to erroneously contend that the administered treatment had little to no effect on the monitored behavior when, in actuality, the effect of the treatment was not evidenced until an additional condition had been implemented. Similarly, a carry-over effect refers to the influence of a previous level of treatment on the observed behavior in a subsequent level of the same treatment condition.

Carry-over and latency effects tend to skew results by influencing the responses of participants and can be both positive or negative in nature. A strategy termed *counterbalancing* is frequently employed in repeated measures research to help combat carry-over and latency effects. Counterbalancing involves presenting levels of a treatment condition so that each level occurs equally often at each stage of practice and so that each level precedes another level as many times as it follows the level. When treatment administrations are counterbalanced, it is possible for a researcher to discern if certain combinations of treatment levels adversely affected the observed results.

Counterbalancing the presentation of treatment stimuli in repeated measures designs is critical to the generation of data that accurately represents the effect of a given treatment on the behavior of interest. The following paradigm can be invoked when determining the presentation order of the stimuli provided there are an even number of treatment levels and the number of participants is some multiple of the number of conditions:

1, 2, n, 3, n-1, 4, n-2, 5, n-3, 6, n-4, etc.

Table 1. Counterbalancing Order for Design with Six Treatment Conditions

Person	Trial Number					
	One	Two	Three	Four	Five	Six
A	1	2	6	3	5	4
B	2	3	1	4	6	5
C	3	4	2	5	1	6
D	4	5	3	6	2	1
E	5	6	4	1	3	2
F	6	1	5	2	4	3

Note. Each number (1-6) indicates which level of the treatment the participant would receive on each corresponding trial.

If six treatment levels were to be administered in a given study, the first participant would be presented the treatments levels in the order of 1, 2, 6, 3, 5, 4. To derive the order of presentation for the second participant, it would be necessary to add a 1 to each of the numbers in the preceding order:

$$2=(1+1), 3=(1+2), 1\equiv (1+6 \text{ reduces to } 1), 4=(1+3), 6=(1+5) 5=(1+4).$$

Using this generation rule, the completed order of presentation for five participants when administered five treatment conditions is presented in Table 1.

Notice that each treatment condition is presented before and after every other treatment condition. In the case that there is an odd number of levels of the treatment condition, the first order of presentation is derived as previously illustrated but the second order is computed by reversing the sequence of the first presentation. For example, if there were three levels of treatment to be administered, the first order of presentation would be 1, 2, 3 and the second order of presentation 3, 2, 1. The third order of stimuli presentation would be $1\equiv (1+3, \text{ reduces to } 1)$, $3=(1+2)$ and $2=(1+1)$. This procedure would be repeated until all of the participants were assigned stimuli presentation orders. The examples presented here conform to the criteria delineated by Girden (1992), as each treatment level occurs equally often at each stage of practice and precedes as many times as it follows a level.

The concept of counterbalancing, although important and necessary in repeated measures analyses, cannot always remedy some of the problems experienced in employing repeated measures designs. In some instances the order of presentation of the stimuli is essentially irrelevant in that after the stimuli are presented on the first occasion the participant might be able to demonstrate a substantial practice effect on later trials, which skews the responses on the dependent variable. Consider, for example, a researcher teaching undergraduate students the names of five important psychologists, and then examining the effects of differential drug treatments on memory recall. Many, if not all, of the participants may successfully commit the names to memory on the first trial and will then be able to recite them after each of the treatments regardless of the order of presentation. Counterbalancing is important in generating accurate data in repeated measures designs, but often even a good correction method cannot repair a poor research design.

SS Partitions and the Test Statistic

The remainder of the present paper uses a singular heuristic example to examine the various statistical and conceptual properties of univariate and multivariate approaches to repeated measures analyses. Consider the following situation: A pharmaceutical company has developed a new prototypical drug that is purported to alleviate depressive symptomatology in human adults. In a bid to gain FDA approval of the drug so that the treatment can be disseminated in the public sector, the company tested the medication on a small sample of human participants. The drug was presented in each of four dosage levels (0mg, 150mg, 300mg and 450mg) to each of five individuals. The different levels of treatment were counterbalanced across participants, and each additional level was given only after a sufficient time had passed to allow traces of the previous treatment to exit the participants' systems. After each dosage had an opportunity to take effect, the participants were administered a depression inventory (i.e., dependent variable) to assess

Table 2. Hypothetical Data Matrix for Medication Dosage Study

Subject	Trial Number				Sum	Mean(Y_j)
	0mg	150mg	300mg	450mg		
1	1	10	14	18	43	10.75
2	3	6	15	19	43	10.75
3	3	8	11	20	42	10.50
4	4	8	13	16	41	10.25
5	5	9	13	18	45	11.25
ΣY_j	16	41	66	91	214	
Mean(Y_j)	3.2	8.2	13.2	18.2		10.70

Note. Scores on the hypothetical depression inventory range from 0 ‘critically depressed’ to 25 ‘not depressed’.

their level of depression (ranging from 0 ‘critically depressed’ to 25 ‘not depressed’). After the conclusion of the last treatment, a data matrix was compiled and examined using a single within-subject factor (drug treatment) repeated measures analysis. These hypothetical data are presented in Table 2.

The first step in completing the repeated measures ANOVA is to partition the variance on the dependent variable. As mentioned earlier, the variance will be decomposed into three main components: A portion associated with the differences in the participants (SS_{SUB}); a component associated with differences in the treatment conditions or intervals (SS_{TRT}); and a portion associated with error variance and the effects of individual participant differences across treatment conditions (SS_{RES}). Thus, the decomposition of the total variation on the dependent variable can be described in the equation:

$$SS_T = SS_{SUB} + SS_{TRT} + SS_{RES} .$$

The SS_T component represents the sum of the squared deviation of each dependent variable score from the grand mean. SS_{SUB} is computed by summing the squared deviations of each subject mean from the grand mean. SS_{TRT} represents the sum of the squared deviations of the treatment means from the grand mean. Finally, SS_{RES} is computed by subtracting SS_{SUB} and SS_{TRT} from SS_T . Thus, for the present example, the equation can be written:

$$660.20 = 2.20 + 625.00 + 33.00.$$

The next step is to compute the appropriate degrees of freedom for each source of variation. The formulas for the df calculations are presented in Table 3. As noted earlier, the residual term is considered an interaction effect (subjects by treatment intervals) and thus the df for this component is calculated by multiplying the df for subjects (df_{SUB}) by the df for treatment intervals (df_{TRT}). After calculating the appropriate df , the mean squares, F_{calc} , eta^2 , and $omega^2$ values are calculated. Notice that the only F_{calc} value that is computed is the value corresponding to the treatment interval source of variation because it is usually of primary interest.

The results of the repeated measures ANOVA on the example data are very favorable. The different treatment conditions accounted for 94.67% of the total variation on the depression scores and rendered a statistically significant result at the $\alpha = .05$ level. Additionally, there was very little variation in the performance of individual participants when averaged across conditions as illustrated by the very small SS value for the between subjects source of variation (2.20). Based on the effect size and statistical significance of these results, the pharmaceutical company can report that different dosages of the anti-depression medication produced a decrease in the overall depression scores of the participants involved in the study.

A comparison of the results in the repeated measures ANOVA with the results generated by a classical ANOVA of the same set of data produces an interesting topic for discussion. In order to invoke a classical ANOVA with this data, it would be necessary to measure each individual only once as classical ANOVA examines effects between subjects; consequently, 20 participants would be required to complete the exact same study. In the classical ANOVA approach rather than each participant receiving all four

Table 3. Summary Table for Repeated Measures ANOVA with Table 2 Data.

Source	SS	df	MS	F_{calc}	η^2	ω^2
Subjects	2.20	$n-1 = 4$	0.55			
Intervals (TRT)	625.00	$k-1 = 3$	208.3	75.76	0.9467	0.9144
Residual	33.00	$(n-1)(k-1) = 12$	2.75			
Total	660.20	$(n)(k)-1 = 19$				

Note. n = number of Subjects, k = number of Treatment Intervals.

$$\omega^2 = (SS_{\text{TRT}} - ((k-1)MS_{\text{RES}})) / (SS_{\text{T}} + MS_{\text{SUB}} + (nMS_{\text{RES}})).$$

levels of treatment, each individual would be randomly assigned one level of treatment. Consequently, it is only possible to examine the differences between the levels of treatment and not the differences between individual participants. Thus, the variance on the dependent variable is partitioned into one less component in classical ANOVA than in the repeated measures analysis.

The results of the four level one-way ANOVA on the example data are presented in Table 4. Notice that the SS_{RES} (error) is larger in the classical ANOVA analysis (i.e., SS_{W}) because there is one less variance partition. Also note that the df_{RES} is smaller in the repeated measures ANOVA than in the classical ANOVA. Because only 2.20 SS units were accounted for by the inclusion of this source of variation in the repeated measures analysis and the cost of examining that piece of information was 3 df from a small number of 20 total observations, employing a repeated measures analysis in this case could have been detrimental to the outcome of the study in terms of statistical significance due. Even though both results are statistically significant at the $\alpha = .05$ level, the F_{calc} value is much larger in the classical ANOVA analysis. All things being equal, larger F_{calc} values will lead to a better chance of obtaining statistical significance, if the researcher is concerned with doing so.

Notice, however, that the variance accounted for by the treatment component (94.67%) remains identical in each analysis. In a real world setting with actual data and a larger sample size, the distinction between these two approaches would be more pronounced, and the power of utilizing a repeated measures approach could be illustrated more convincingly. Remember, however, that very similar results were generated in both analyses, but the repeated measures approach only utilized five participants total compared to the 20 individuals required to complete the classical ANOVA analysis.

Assumptions in Univariate Repeated Measures Analysis

Invoking a repeated measures ANOVA, however, does not come without several difficulties and considerations. As stated by Girden (1992, p. 13),

. . . several assumptions . . . were recognized by R.A. Fisher in the 1940's but were not demonstrated until the 1950's . . . with the result that many earlier studies involving repeated measures may have reached erroneous conclusions regarding the effect(s) of the independent variable(s).

It is difficult to understand why researchers and statisticians did not better explore and develop the statistical assumptions of repeated measures designs (i.e., compound symmetry and sphericity) earlier in the history of their use. Repeated measures were being used shortly after the development of ANOVA in 1925, and countermeasures were developed to contend with treatment carry-over, latency and practice effects long before Box (1954) articulated the problem of departures from compound symmetry.

Compound Symmetry

The term, *repeated measurements*, implies that each individual is measured on more than one occasion. Consequently, Box (1954) contended that the assumption of compound symmetry, which states that it is necessary for the variances and covariances of different treatment levels to be equal, must be satisfied for the results of the repeated measures analysis to be valid. The observations generated by each individual are independent of the observations produced by any other participant, but each individual's score on a treatment level is often linearly dependent on or correlated with the previous scores. The degree of correlation between a given individual's scores over the treatment conditions affects the

Table 4. Results of Four-Level One-Way ANOVA with Table 2 Data.

Source	<i>SS</i>	<i>df</i>	<i>MS</i>	F_{calc}	η^2
Between	625.00	$k-1 = 3$	208.3	94.70	.9467
Residual (Within)	35.20	$k(n-1) = 16$	2.20		
Total	660.20	$(n)(k)-1 = 19$			

statistical significance of the repeated measures result because the MS_{RES} will decrease as the degree of correlation between the observations increases (thus leading to a smaller MS_{TRT} and a smaller F_{calc} value).

It is important, therefore, to examine the degree of covariation among each pair of treatment scores (Algina & Keselman, 1997; Girden, 1992; Maxwell & Delaney, 1990), as these scores must be equal when invoking an analysis that computes an F_{calc} value ($MS_{\text{TRT}} / MS_{\text{RES}}$) with $(k-1)$ and $(k-1)(n-1)$ *df*. Prior to Box (1954), researchers and statisticians focused solely on the equality of the variances (homogeneity of variance assumption) and ignored any influence that unequal covariances might render on the results of a repeated measures analysis. If the covariance between two sets of scores is defined as

$$\text{COV}_{xy} = \sum (X - \bar{X})(Y - \bar{Y}) / (n-1),$$

then the formula for the covariation of a set of scores with itself could be written as

$$\sum (X - \bar{X})(X - \bar{X}) / (n-1).$$

The latter of the two formulas is the formula for the variance of a set of scores, thus indicating that variance is nothing more than the degree of covariation of a set of scores with itself. Thus, Box (1954) noted that examining the equality of variances in repeated measures analyses was only half of the larger issue; because variance can be defined as the covariation of a set of scores with itself, it is necessary to also examine the equality of covariances of all pairs of the treatment levels as well.

The assumption of compound symmetry derives its name from a matrix that contains information about both variances and covariances of a set of scores. For a given set of scores, the variance of each set of scores and the covariances between all possible pairs of the scores can be arranged in a single matrix termed a variance-covariance matrix. This is a special matrix of rank k , where k indicates the number of treatment levels in the single within factor repeated measures analysis. The variances of the k treatment levels are presented on the main diagonal of the matrix and the covariances between each pair of treatment levels are presented on the off diagonals. A variance covariance matrix that exhibits equal variances and equal covariances is said to demonstrate compound symmetry.

Violating Compound Symmetry

In addition to indicating that the equality of covariances between pairs of the treatment conditions must be considered, Box (1954) noted other important considerations as well. Box effectively demonstrated that if an F_{calc} value does not originate from an \mathcal{F} -distribution with $(k-1)$ and $(k-1)(n-1)$ *df*, then it is not a part of that \mathcal{F} -distribution. Box contended, therefore, that such an F_{calc} value belongs to an \mathcal{F} -distribution that is corrected by a factor called epsilon (ϵ) and which results in $\epsilon(k-1)$ and $\epsilon(k-1)(n-1)$ degrees of freedom. He further noted that the correction factor invoked by the epsilon value is more severe as the covariances become more unequal and is approximately equal to 1.0 when the equality of covariances is demonstrated.

The implication on the statistical significance of a given F_{calc} value can be profound when the epsilon correction factor is severe. Consider the following example: A researcher performs the univariate ANOVA for a single within-subjects factor and generates an F_{calc} value of 3.15 after performing all relevant calculations. The researcher then consults a table of critical values from the \mathcal{F} -distribution to determine if the result is statistically significant. Upon examination, the researcher learns that the F_{crit} value for the study, $F(4, 19)$ at the $\alpha = .05$ level, is 2.90. The researcher then revels in the limelight of a statistically significant result generated by varying the treatment conditions.

This result, however, is only valid provided that the variances and covariances of the treatment levels are equal. If the variances and covariances between treatment levels in the example study are not equal, then the researcher must alter the *df* with which to evaluate the statistical significance of the F_{calc} value by

a factor of epsilon. If the epsilon value is equal to .5, the new df for the F_{calc} would be .5(4) and .5(19); thus, when accounting for the correction factor, the df of 2 and 10 must be used. When the F_{crit} value using these df is examined, the researcher is dismayed to learn that the corrected F_{crit} value is 4.10, thus rendering the F_{calc} value of 3.15 no longer statistically significant at $\alpha = .05$. The conceptual point here is that when epsilon provides a severe correction for the inequality of variances and covariances, the df for the F_{crit} value can be radically affected. Consequently, researchers may use inappropriate df that generate a F_{crit} value that is smaller than it should be and by which a researcher may reject a null hypothesis when it should not be rejected.

Box (1954) demonstrated that the upper bound for the epsilon correction is 1.0 which indicates that the epsilon corrected df are exactly equal to the \mathcal{F} -distribution df . At this extreme a correction factor is not invoked and erroneous conclusions are avoided. Geisser and Greenhouse (1958) determined the lower bound for the epsilon correction factor by positing that the lower bound of epsilon was equal to $1/(k-1)$. In the previous example with five treatment levels administered to five subjects, the lower bound of epsilon would have been $1/(4-1)$ or .33. At this extreme the most severe correction would be invoked, and the df for F_{crit} would be exactly one third of their original magnitude.

Sphericity

Research conducted after Box's (1954) work on compound symmetry indicated that although compound symmetry is important to consider in repeated measures analyses, it is not a necessary condition for conducting them (Girden, 1992). Sphericity (sometimes called circularity) is considered a necessary and sufficient condition to conduct repeated measures analyses (Huynh & Feldt, 1970; Rouanet & Lepine, 1970). Sphericity is the degree to which variances in the differences between pairs of treatment scores are equal. The notion of sphericity is a more flexible assumption and subsumes compound symmetry as a special case. That is, the variances of differences between treatment levels would be equal when the variances and covariances were all equal.

To satisfy the assumption of sphericity, the variances of the differences in all pairs of treatment combinations must be homogeneous (e.g., the variance of level 1 and 2 must equal the variance of level 2 and 3, etc). The variance of a difference between two treatment conditions (1 and 2) can be defined as

$$\sigma^2(Y_1 - Y_2) = \sigma^2_1 + \sigma^2_2 - 2\sigma_{12}$$

where σ^2_1 is the variance of on a set of scores, σ^2_2 is the variance of another set of scores, and σ_{12} is the covariance of the two sets of scores.

Assessing the equality of variances of differences in the treatment levels can be illustrated using the heuristic example data presented in Table 2. Table 5 contains the variances and covariances for each of the treatment levels. By invoking the formula described above for the variance of a difference between treatment levels, it is possible to compute the variance of the difference between treatment levels 1 (0mg) and 2 (150mg). The variance of the difference would be

$$\sigma^2(Y_1 - Y_2) = 2.1998 + 2.1998 + 2(-.5500) = 3.2996.$$

Because the variances and covariances of each of the treatment levels are exactly equal in this example, all of the variances of difference would be equal to 3.2996. The assumption of sphericity, therefore, would be satisfied in this example. This example has the additional feature of being compound symmetrical, as all of the variances and covariances are exactly equal.

Another way to assess the sphericity of a data set is to examine the matrix of orthonormal contrasts (Stevens, 1996). If the multivariate identity

$$\mathbf{C}'\Sigma\mathbf{C} = \sigma^2\mathbf{I}$$

where \mathbf{C} is a matrix of $(k-1)$ orthonormal contrasts, \mathbf{C}' is the transpose of \mathbf{C} , Σ is the variance-covariance matrix and $\sigma^2\mathbf{I}$ is an identity matrix (with equal variances on the main diagonal and zeros on the off diagonal) is true, the assumption of sphericity is satisfied. The first step in assessing sphericity in this manner is to create a set of orthogonal contrast variables. One set of orthogonal contrast variables for the example data are presented in Table 6. Notice that there are three contrast variables present corresponding to the $(k-1)$ contrasts possible. The contrast variables are then normalized by invoking a multiplicative constant such that the sum of the squared transformed coefficients in a given contrast is equal to 1.0. This

Table 5. Correlation and Variance-Covariance Matrix for Table 2 Data

Tx	0mg	150mg	300mg	450mg
1	2.20	-0.55	-0.55	-0.55
2	-0.55	2.20	-0.55	-0.55
3	-0.25	-0.25	2.20	-0.55
4	-0.25	-0.25	-0.25	2.20

Note. Variances are on the main diagonal. Correlations are below and Covariances are above the main diagonal.

Table 6. Orthogonal Contrast Variables.

Tx	C ₁	C ₂	C ₃
1	1	1	1
2	-1	1	1
3	0	-2	1
4	0	0	-3

Note. Tx = treatment level; C₁ = contrast variable 1; C₂ = contrast variable 2; C₃ = contrast variable 3

Table 7. Orthonormal Contrast Matrix C.

Tx	C ₁	C ₂	C ₃
1	.707	.408	.289
2	-.707	.408	.289
3	0	-.816	.289
4	0	0	-.866

is accomplished by first squaring the coefficients in the contrast and dividing by the derived number. For example, the second contrast variable contains the contrast coefficients of 1, 1, and -2. It is necessary to first square each coefficient, $(1)^2 + (1)^2 + (-2)^2 = 6.00$, and then to divide each coefficient by the result, $1/\sqrt{6}$, $1/\sqrt{6}$, and $(-2/\sqrt{6})$. This procedure would be repeated until all contrasts are transformed. The orthonormalized variables are placed in a matrix as presented in Table 7. The product of the equation, $C'\Sigma C = \sigma^2\mathbf{I}$, is then computed to determine if the assumption of sphericity is met (Stevens, 1996).

Violating the Assumption of Sphericity.

Violating the assumption of sphericity can have the same grievous consequences as violating the assumption of compound symmetry. If the variances of the differences in levels of the treatment conditions are not equal, the df for the F_{calc} value will be smaller than normal by a value of ε . Thus, just as violating compound symmetry caused the F_{crit} value to increase in magnitude, so too, does violating the assumption of sphericity. If sphericity is not examined, the unwary researcher will tend to reject the null hypothesis more often than it should be rejected resulting in higher Type I error rates (Stevens, 1996).

Correcting for Violations in the Sphericity Assumption

Adjusted Degrees of Freedom. Correcting for a violation in the sphericity assumption involves the computation of the epsilon parameter previously described. As noted earlier, Box (1954) defined the upperbound of epsilon to 1.0, and Geisser and Greenhouse (1958) calculated the lower bound of epsilon to be $1/(k-1)$, where k is the number of treatments utilized in the study. To find the actual value of epsilon, the following formula generated by Greenhouse and Geisser (1959) must be invoked,

$$\varepsilon = \frac{k^2(\bar{s}_{ii}^2 - \bar{s}_{**})^2}{(k-1)(\sum \sum s_{ij}^2 - 2k \sum \bar{s}_i^2 + k^2 \bar{s}_{**}^2)}$$

where \bar{s}_{ii}^2 is the mean of entries on the main diagonal of the variance-covariance matrix, \bar{s}_{**} is the mean of all entries in the variance-covariance matrix, s_{ij} is the ij^{th} entry of the variance-covariance matrix and \bar{s}_i is the mean of all entries in the row i . The calculation of epsilon can be illustrated using the example data in Table 2. The calculations

$$\varepsilon = \frac{4^2(2.2000 - 0.1375)^2}{(3)(22.9900 - (8)(0.075625) + (16)(0.1375^2))}$$

result in an ε value of 1.0. This value indicates that the sphericity assumption is perfectly met in the example data and that no corrections to the degrees of freedom are required. To determine if the result is statistically significant, the researcher would simply use the same degrees of freedom as calculated when constructing the summary table.

Table 8. Hypothetical Data Matrix for Worst Case of Violating Sphericity Assumption.

Subject	Trial Number				Sum	Mean(Y_i)
	0mg	150mg	300mg	450mg		
1	1	10	15	20	46	10.50
2	2	6	11	16	35	8.75
3	3	7	12	17	39	9.75
4	4	8	13	18	43	10.75
5	5	9	14	19	47	11.75
ΣY_j	15	40	65	90	210	
Mean(Y_j)	3.0	8.0	13.0	18.0		10.70

Consider an example in which the epsilon value invokes the greatest possible correction factor on the degrees of freedom used to determine F_{crit} . If the example data in Table 2 are changed only slightly to resemble the data in Table 8, the results of the repeated measures analysis are radically different. The correction factor invoked in the Table 8 data is $\epsilon = .33$ or the minimum value possible in a design with four treatment levels ($1/k-1$) because the scores are linearly dependent and completely violate the sphericity assumption. Rather than using $(k-1)=3$ and $(k-1)(n-1)=12$ degrees of freedom to arrive at an F_{crit} value of 3.49 ($\alpha = .05$), each value would be corrected by the value of ϵ . Thus, the new degrees of freedom for the F_{crit} value would be 1 and 4 and would render an F_{crit} value of 7.71 ($\alpha = .05$). The results in the present study would still be statistically significant, but in situations where the F_{calc} value is only slightly larger than the uncorrected F_{crit} value, even a small correction by epsilon can alter the results from statistically significant to not statistically significant. This comparison allows the power of a repeated measures design to be illustrated: If ϵ is equal to 1.0, the repeated measures ANOVA demonstrates the power equivalent to (nxk) participants. Consequently, fewer participants can be utilized without compromising sufficient power to reject the null hypothesis.

Estimates of Epsilon. There are several estimates of ϵ currently available, and each has its own distinct advantages and disadvantages (Maxwell & Delaney, 1990; Stevens, 1996). The most commonly employed indices of ϵ are the Greenhouse-Geisser ϵ (Greenhouse & Geisser, 1959) and the Huynh-Feldt ϵ (Huynh & Feldt, 1970). Because both ϵ values are computed in most statistical software packages, it is important to understand the particular bias that each estimate produces.

When there is only one sample, as with the present example, the Greenhouse-Geisser and Huynh-Feldt estimates are identical. However, for more two or more independent groups (i.e., split-plot design), the Greenhouse-Geisser ϵ tends to underestimate the true value of epsilon across the range of values, but the underestimation is even more pronounced as epsilon approaches 1.0. Consequently, the Greenhouse-Geisser ϵ will produce a very conservative estimate of the df utilized to obtain the F_{crit} value which may result in not rejecting the null hypothesis as often as might be indicated by the data. Conversely, the Huynh-Feldt ϵ produces an overestimation of the true value of ϵ and may result in a smaller F_{crit} value than is indicated by the data. Thus, when using the Huynh-Feldt ϵ , researchers may reject the null hypothesis more often than they should. Because the two most popular estimates of ϵ produce biased results, authors have recommended averaging the two indices as a more accurate estimate of ϵ (Barcikowski & Robey, 1984; Girden, 1992; Stevens, 1996). If one must be chosen over the other, however, it is always somewhat safer to utilize the Greenhouse-Geisser ϵ as it produces a more conservative correction factor.

Guidelines. Guidelines have been presented to facilitate the correct interpretation of univariate repeated measures analyses (Greenhouse & Geisser, 1959). Due to the advent of statistical software that readily and painlessly computes many different estimates of ϵ , these guidelines have limited pragmatic value but warrant brief consideration. As provided by Girden (1992, p. 21) the steps are as follows:

1. Compare the obtained [F_{calc}] with the tabled value corresponding to [$k-1$] and [$k-1$][$n-1$] df . If it is not greater than this most liberal value, stop at this point. It will not be significant when degrees of freedom are reduced.

2. If the obtained $[F_{\text{calc}}]$ is significantly higher than the most liberal value, enter the table with 1 and $[n-1]$ df . If the obtained F is greater than this most conservative value, it is significant. Stop at this point.
3. If the obtained $[F_{\text{calc}}]$ is higher than the tabled value for $df = [k-1]$ and $[k-1][n-1]$, but lower than the tabled value for $df = 1$ and $[n-1]$, then the ε adjustment should be applied.

It is not possible in the present example to illustrate this dynamic because the ε value of the data set is 1.0, but interested readers can examine the analysis of the data set in Girden (1992) for a good example of a situation where these guidelines might prove helpful.

Another Univariate Approach to Repeated Measures Analysis

As mentioned previously, it is possible to employ multiple regression analysis to examine the effects of repeated measures data. A brief treatment of this analytic technique is presented here, and interested readers are referred to Edwards (1985) for a more detailed discussion.

To utilize this approach it is necessary to first construct k contrasts (where k represents the number of treatment levels) in which each contrast variable represents a separate effect. In the example data set four treatment conditions were utilized; thus, four orthogonal contrasts should be created. The first three contrasts represent linear, quadratic and cubic trends in the data. [Any orthogonal trends could be used, but polynomial orthogonal trends are arbitrarily used here.] The final contrast employed is a sum vector that adds the responses of each participant over all four treatments. The orthogonal contrast matrix is presented in Table 9.

After the contrast matrix is generated, each contrast variable can be entered into the multiple regression equation in a hierarchical manner to determine the unique variance accounted for by each contrast. The results of the multiple regression analysis on the example data are presented in Table 10. Notice that the linear contrast variable (C_1) accounts for the same variance that is associated with the effects of the treatments (the interval component of the repeated measures ANOVA) in that the SS units and the R^2 values in the multiple regression analysis are exactly equal to the SS and η^2 values in the summary tables for the repeated measures ANOVA presented in Table 3 and the classical ANOVA results presented in Table 4. The only other contrast variable to generate noteworthy consideration is the contrast variable C_4 , as the SS and R^2 values for this contrast directly correspond to the SS and η^2 values for the between subjects source of variation in the repeated measures ANOVA. Thus, the multiple regression approach generates results identical to the repeated measures ANOVA by utilizing a series of orthogonal contrast variables. Because all three of these methods (classical ANOVA, repeated measures ANOVA and multiple regression) generated the exact same SS partitions with the example data, the conceptual unity of the three approaches has been illustrated.

Multivariate Approach to Repeated Measures Analysis

The preceding portion of the paper was spent solely on examining the different univariate approaches to repeated measures analysis. A single factor repeated measures design, however, can be analyzed by using a multivariate approach. The multivariate approach, Multivariate Analysis of Variance (MANOVA), invokes a different sort of logic in completing the analysis. Rather than treating several measurements over time as a single dependent variable repeatedly measured, the multivariate approach treats the repeated measurements as separate dependent variables generated by one individual. Thus, in the example that has been used consistently throughout the paper, the multivariate approach would conceptually consider each of the four measures of depression as a separate dependent variable.

There are both advantages and disadvantages to conducting repeated measures analyses via the multivariate approach. One advantage of conducting this type of analysis is that the measurements are allowed to have any correlational structure (unlike in repeated measures ANOVA). That is, the sphericity assumption becomes unnecessary. This approach may more closely honor a given researcher's reality, provided that the researcher believes measurements to be correlated in a real world situation. In repeated measures ANOVA, the dependence among measures was considered to be of fixed form, and the researcher was penalized for analyzing data that did not exhibit equal correlations between measurements. In the multivariate approach, sphericity is not a consideration because each measurement is deemed a separate and unique dependent variable.

Table 9. Orthogonal Contrasts for Regression Analysis of Example Data in Table 2.

Subject	Dose	C_1	C_2	C_3	C_4	Y
1	0	-3	1	-1	43	1
2	0	-3	1	-1	43	3
3	0	-3	1	-1	42	3
4	0	-3	1	-1	41	4
5	0	-3	1	-1	45	5
1	150	-1	-1	3	43	10
2	150	-1	-1	3	43	6
3	150	-1	-1	3	42	8
4	150	-1	-1	3	41	8
5	150	-1	-1	3	45	9
1	300	1	-1	-3	43	14
2	300	1	-1	-3	43	15
3	300	1	-1	-3	42	11
4	300	1	-1	-3	41	13
5	300	1	-1	-3	45	13
1	450	3	1	1	43	18
2	450	3	1	1	43	19
3	450	3	1	1	42	20
4	450	3	1	1	41	16
5	450	3	1	1	45	18

Note. C_1 = contrast variable 1 (linear trend), C_2 = contrast variable 2 (quadratic trend), C_3 = contrast variable 3 (cubic trend), C_4 = sum vector, and Y = the dependent variable score.

Although discarding the sphericity assumption may sound enticing to even the most seasoned researcher, there are disadvantages to employing the multivariate approach to analyze repeated measures data. The primary disadvantage to utilizing the repeated measures approach is that statistical significance is difficult to obtain. Because the multivariate approach presumes that each measurement by an individual participant is a separate dependent variable, the advantage of repeatedly measuring participants is forfeited (i.e., a smaller number of participants necessary to generate results similar to classical ANOVA). Consequently, sample size issues become a paramount consideration, and a sample size that was sufficient to demonstrate an effect in a repeated measures ANOVA may be too small to demonstrate a similar effect using the multivariate approach.

The multivariate approach can be utilized in one of two ways: (a) by either transforming the score set into $(k-1)$ difference variables and then analyzing the new variable set or (b) by creating a matrix of orthogonal or orthonormal coefficients, weighting each score by the corresponding coefficient, and analyzing the new matrix. Either method will generate identical results (Stevens, 1996), and there is no clear advantage to employing either method. The only stipulation in using either method is that only $(k-1)$ new variables are created. The new matrix (either the differenced matrix or the transformed orthonormal matrix) is then analyzed using Hotelling's T^2 calculated by:

$$T^2 = n/(n-1)[(C'M)'(C'\Sigma C)^{-1}(C'M)],$$

where Σ^{-1} is the inverse of the variance-covariance matrix, \mathbf{M} is a column vector of means and \mathbf{C} is a matrix of $(k-1)$ orthogonal, orthonormal, or difference contrasts, \mathbf{C}' is the transpose of \mathbf{C} . The calculated value of T^2 for the example data is 56.8181. This value can easily be transformed into an F_{calc} by invoking the formula:

$$F_{\text{calc}} = T^2 (n - k + 1)/(k - 1),$$

where n is the number of participants and k is the number of treatment conditions. In this example, the

Table 10. Results of Multiple Regression Analysis with Table 2 Data.

Source	SS	df	MS	F_{calc}	R	R^2
C_1	625.00	1	625.00	94.70	.9730	.9467
C_2	0.00	1	0.00	0	0	0
C_3	0.00	1	0.00	0	0	0
Residual	35.20	18	1.96			
C_1	625.00	1	625.00	227.27	.9730	.9467
C_2	0.00	1	0.00	0.00	0	0
C_3	0.00	1	0.00	0.00	0	0
C_4	2.20	4	0.55	0.36	.0580	.0033
Residual	33.00	12	2.75			
Total	660.20					

computation is performed [$F_{\text{calc}} = 56.8181(5 - 4 + 1)/(4 - 1)$] and the transformed value of F_{calc} is equal to 37.879. The statistical significance of the F_{calc} value can be evaluated by using $(k-1)$ and $(n-k+1)$ degrees of freedom. For this example, the F_{crit} value is $F(3,2) = 99.164$ at the $\alpha = .01$ level. In contrast to the results of the univariate ANOVA, this result is not statistically significant primarily due to a radically smaller df_{RES} than in the univariate repeated measures ANOVA analysis.

Another reason statistical significance was not obtained by using the multivariate approach is that the sample size appears smaller. Because the ε value was 1.0 in the example data, the repeated measures ANOVA demonstrated the power equivalent of 20 participants (five participants measured on four occasions) whereas the multivariate approach presumed only five individuals measured on four occasions with four different dependent variables. The divergence in these two outcomes illustrates that the multivariate approach to repeated measures analysis is a more conservative method of analyzing effects and may not indicate statistically significant results even when other analytic methods do.

Are Univariate or Multivariate Methods Superior?

Repeated measures designs offer a number of advantages to social science researchers, one of which is the economy of research participants included in the study. As noted previously, the total size of the sample is decreased in the repeated measures case because the logic employed in the research design is different. That is, in a repeated measures design, participants serve as their own controls to eliminate the effects of individual participant error, thus necessitating a fewer number of participants. However, the advantages of the using the repeated measures design does not come without a cost. Disadvantages of this strategy include the introduction of carry-over, latency, and general practice effects, none of which can be fully removed from a given study (Keppel, 1991).

In regards to the analysis of repeated measures data, statisticians and researchers have debated whether univariate or multivariate approaches to repeated measures analyses are preferable with no clear consensus emerging from the discussions (Algina & Keselman, 1997; Barcikowski & Robey, 1984; Stevens, 1996). There are situations in which a univariate repeated measures ANOVA would be the most effective method of analyzing the data (e.g., if the sphericity assumption is satisfied and/or the between subjects and interval variance partitions account for the majority of the variance on the dependent variable) and other situations in which a multivariate approach would provide more favorable results (e.g., when the sphericity assumption is severely violated and/or when there is more variation in the treatment intervals than between the participants). Authors have recommended that both analyses be performed, because one or the other may be more powerful depending on the characteristics of the data (Barcikowski & Robey, 1984). This raises the question of which of the two analyses should be reported, and common sense indicates that if both analyses are completed, then this fact should be acknowledged to the consumers of the research. Most recently, however, the logic of employing both analytic strategies in the same analysis has been questioned (Keselman, Keselman, & Lix, 1995).

An important caveat must be noted. Multivariate methods were described earlier as offering an alternative to the univariate approach to analyzing repeated measures data and by which the sphericity assumption was not relevant. Although this comment is literally true, it is only accurate if the assumption

of multivariate normality is met. Multivariate normality is the degree to which all linear combinations of several variables are normal (Henson, 1999). Further, simply ascertaining that variables are univariate and bivariate normal is not adequate to ensure that the system of variables is multivariate normally distributed. Rather, it is necessary to determine that each variable is normally distributed about fixed variables on all other variables.

Failing to assess multivariate normality can have deleterious consequences on statistical analyses, including those conducted with repeated measures designs. As noted by Marascuilo and Levin (1983, p. 203), “multivariate normality’s impact and role . . . are basic to the inference procedures of multivariate analysis.” Similarly, Thompson (1996b, p. 4) noted that, “Although multivariate normality is not required to estimate most multivariate parameters (e.g., function coefficients, structure coefficients), even in these cases the distributions of the variables must be reasonably comparable.” Consequently, multivariate normality is an assumption that must be satisfied to ensure the accuracy and correct interpretation of multivariate results.

When multivariate normality has been met, however, the actual α level of the study is mathematically guaranteed to equal the preset α at the beginning of the study (Maxwell & Delaney, 1990). Consequently, if the multivariate normality assumption is not met, the α rate for the entire study will rarely equal the preset α level and may run as high as 10-20% (Tanguma, 1999). It is on this basis that Maxwell and Delaney recommended the use of the multivariate approach to repeated measures analysis if it is the researcher’s desire is to avoid falsely rejecting the null hypothesis.

Higher-than-preset α rates are often a concern when the sphericity assumption is violated as well because many researchers do not adjust the degrees of freedom accordingly and completely overlook this very important statistical assumption. Even when sphericity is examined, there are only two options for adjusting for its violation: (a) adjust the degrees of freedom by a value of epsilon; or (b) use the multivariate approach. If the former is chosen as the method of choice, the results are only approximate because the epsilon value is an estimated correction factor.

Thus, the following general rules are posited based on the conjecture presented in previous research: (a) if sphericity is not violated, the univariate methods tend to be more powerful than the multivariate methods; (b) if sphericity is violated, neither method (univariate versus multivariate) tends to be preferable to the other, although a multivariate method can be used providing the variables are multivariate normal; and (c) if the size of the sample is substantially greater than the number of levels in the repeated variable, then the multivariate methods are preferred and tend to produce better results.

Summary and Conclusions

The present paper explored the analysis of univariate repeated measures designs by using both univariate and multivariate approaches. The univariate tests examined generated desirable results when certain statistical assumptions in the data were satisfied. The multivariate approach to repeated measures analysis was found to generate results generally free from the statistical assumptions in the univariate case but which tended to be a more conservative estimate of the effects. Both types of analyses were found to be valuable when certain situations were presented and when various assumptions were met.

If the assumption of sphericity is violated, the researcher is always able to employ the multivariate analysis as sphericity is not a required assumption in this approach. However, the multivariate method presumes that multivariate normality has been met, which is often not the case. The univariate approach is often more powerful than the multivariate approach when score variances are homogeneous as the df_{RES} is larger for the univariate test than for Hotelling's T^2 . If the data are characterized by extreme variability, small effects of the treatment conditions may be hidden in the univariate data but elucidated by the multivariate approach. In situations where both the sphericity and normality assumptions are violated, multivariate rank-based procedures have been shown to control Type I errors and provide more statistical power than parametric procedures (Agresti & Pendergast, 1986; Beasley, in press). The bottom line, therefore, is that common analytic sense should guide practice, and researchers should critically examine their data to determine the most appropriate analytic strategy.

References

- Agresti, A., & Pendergast, J. (1986) Comparing mean ranks for repeated measures data. *Communications in Statistics: Theory & Method*, 15, 1417-1433.
- Algina, J., & Keselman, H.J. (1997). Detecting repeated measures effects with univariate and multivariate statistics. *Psychological Methods*, 2, 208-218.
- Barcikowski, R.S., & Robey, R.R. (1984). Decisions in single group repeated measures analysis: Statistical tests and three computer packages. *American Statistician*, 38, 148-150.
- Beasley, T. M. (in press) Multivariate aligned rank test for interactions in multiple group repeated measures designs. *Multivariate Behavioral Research*.
- Benton, R.L. (1991). Statistical power considerations in ANOVA. In B. Thompson (Ed.), *Advances in educational research: Substantive findings, methodological developments* (Vol. 1, pp. 119-132). Greenwich, CT: JAI Press.
- Box, G.E.P. (1954). Some theorems on quadratic forms in the study of analysis of variance problems: II. Effects of inequality of variance and of correlation between errors in the two-way classification. *Annals of Mathematical Statistics*, 25, 484-498.
- Edwards, A.L. (1979). *Multiple regression and the analysis of variance and covariance*. San Francisco: Freeman and Company.
- Edwards, A.L. (1985). *Experimental design in psychological research* (5th ed.). New York: Harper & Row.
- Fisher, R.A. (1925). *Statistical methods for research workers*. Edinburgh, England: Oliver and Boyd.
- Girden, E.R. (1992). *ANOVA: Repeated measures*. Newbury Park, CA: Sage.
- Geisser, S., & Greenhouse, S.W. (1958). An extension of Box's results on the use of the F distribution in multivariate analysis. *Annals of Mathematical Statistics*, 29, 885-891.
- Greenhouse, S.W., & Geisser, S. (1959). On methods in the analysis of profile data. *Psychometrika*, 24, 95-112.
- Haase, T., & Thompson, B. (1992, January). *The homogeneity of variance assumption in ANOVA: What it is and why it is required*. Paper presented at the annual meeting of the Southwest Educational Research Association, Houston, TX.
- Henson, R.K. (1999). Multivariate normality: What is it and how is it assessed? In B. Thompson (Ed.), *Advances in social science methodology* (vol. 5, pp. 193-211). Stamford, CT: JAI Press.
- Huynh, H., & Feldt, L.S. (1970). Conditions under which mean square ratios in repeated measurements designs have exact *F*-distributions. *Journal of the American Statistical Association*, 65, 1582-1589.
- Keppel, G. (1991). *Design and analysis: A researcher's handbook* (3rd ed.). Englewood Cliffs, NJ: Prentice-Hall.
- Keppel, G., & Saufley, W.H. (1980). *Introduction to design and analysis: A student's handbook*. San Francisco: W.H. Freeman & Company.
- Keppel, G., & Zedeck, S. (1989). *Data analysis for research designs*. New York: W.H. Freeman & Company.
- Keselman, H.J., Keselman, J.C., & Lix, L.M. (1995). The analysis of repeated measures measurements: Univariate tests, multivariate tests, or both? *British Journal of Mathematical and Statistical Psychology*, 48, 319-338.
- Keselman, J.C., Lix, L.M., & Keselman, H.J. (1996). The analysis of repeated measurements: A quantitative research synthesis. *British Journal of Mathematical and Statistical Psychology*, 49, 275-298.
- Marascuilo, L.A., & Levin, J.R. (1983). *Multivariate statistics in the social sciences: A researcher's guide*. Monterey, CA: Brooks/Cole.
- Maxwell, S.E., & Delaney, H.D. (1990). *Designing experiments and analyzing data: A model comparison perspective*. Belmont, CA: Wadsworth.
- Rouanet, H., & Lepine, D. (1970). Comparison between treatments in a repeated-measurement design: ANOVA and multivariate methods. *British Journal of Mathematical and Statistical Psychology*, 23, 147-163.
- Shavelson, R.J. (1988). *Statistical reasoning for the behavioral sciences* (2nd ed). Boston: Allyn & Bacon.
- Snyder, P.A., & Lawson, S. (1993). Evaluating results using corrected and uncorrected effect size

- estimates. *Journal of Experimental Education*, 61, 334-349.
- Stevens, J. (1996). *Applied multivariate statistics for the social sciences*. Mahwah, NJ: Lawrence Erlbaum.
- Tanguma, J. (1999). Analyzing repeated measures designs using univariate and multivariate methods. In B. Thompson (Ed.), *Advances in social science methodology* (vol. 5, pp. 233-250). Stamford, CT: JAI Press.
- Thompson, B. (1993). The use of statistical significance tests in research: Bootstrap and other alternatives. *Journal of Experimental Education*, 61, 361-377.
- Thompson, B. (1996a). AERA editorial policies regarding statistical significance testing: Three suggested reforms. *Educational Researcher*, 25(2), 26-30.
- Thompson, B. (1996b). *Problems with multivariate normality: Can the multivariate bootstrap help?* Paper presented at the annual meeting of the Society for Applied Multivariate Research, Houston. (ERIC Document Reproduction Service No. ED 420 154)
- Thompson, B. (1999a). Journal editorial policies regarding statistical significance tests: Heat is to fire as *p* is to importance. *Educational Psychology Review*, 11, 157-169.
- Thompson, B. (1999b). If statistical significance tests are broken/misused, what practices should supplement or replace them? *Theory and Psychology*, 9, 167-183.
- Wilcox, R.R. (1987). *New statistical procedures for the social sciences: Modern solutions to basic problems*. Hillsdale, NJ: Lawrence Erlbaum.

Send correspondence to: Kevin M. Kieffer, Department of Psychology MC2127,
Saint Leo University, P.O. Box 6665, Saint Leo, FL 33574.
Email: kmkieffer@earthlink.net or Kevin.Kieffer@saintleo.edu.

A Monte Carlo Simulation Comparing Parameter Estimates from Multiple Linear Regression and Hierarchical Linear Modeling

Daniel J. Mundfrom

Mark R. Schultz

University of Northern Colorado

In this simulation study, the parameter estimates obtained from hierarchical linear modeling (HLM) and multiple linear regression (MLR) were examined for differences under different values of the intraclass correlation. 15,000 data sets were generated for each of ten different ranges of intraclass correlations. The resulting vectors of parameter estimates from both HLM and MLR were subtracted, averaged across 50 data sets and compared to a null vector of zeros using Hotelling's T^2 statistic. Little difference was found between the vectors of parameter estimates in any of the intraclass correlation ranges.

Multiple Linear Regression (MLR) and Hierarchical Linear Modeling (HLM) are two statistical procedures that can be used to model the relationship between a numerical dependent (i.e., response) variable and two or more numerical independent (i.e., predictor) variables. For an algebraic description and comparison of the HLM and MLR models see Mundfrom & Schultz (2001). Although similar in many ways, these two procedures are not identical in how they analyze the data and consequently may not produce the same results on any specific set of data. Raudenbush & Bryk (1986) and Goldstein (1987) have suggested that HLM is useful for data-analytic situations that may not be adequately handled using the general linear model, of which MLR is a specific case. Specifically, HLM is believed to be, and in fact was designed to be, more accurate in situations involving multi-level data, i.e., situations in which data are measured at more than one level.

Multi-level data situations are not uncommon, particularly in educational research. A common example involves data measured at both the student level and also at the teacher and/or the school level. This hierarchical or nested structure does not appear to be adequately modeled using the general linear model framework or more specifically, multiple linear regression. However, Mundfrom and Schultz (2001) found that little difference existed between the predicted values generated using HLM and those obtained from MLR when an appropriate MLR model was utilized. They did find some differences in parameter estimates between the two procedures, although in most cases those differences were small. Although their findings were obtained from comparing a relatively few actual data sets, the results would seem to indicate that MLR may be an appropriate alternative for analyzing multi-level data.

Purpose

This study was designed to examine more closely differences among the parameter estimates between HLM and MLR. Littel, Milliken, Stroup, & Wolfinger (1996) show examples of analyses using the general linear model produce identical results to ones using an HLM model. Bryk & Raudenbush (1992) on the other hand cite examples in which the analyses using the two procedures produce similar, but different results. One possibility for explaining why in some cases HLM and MLR produce parameter estimates that are the same whereas in other instances these estimates differ could be differing correlational structures in the data. Specifically, perhaps the size of the intraclass correlation could be affecting the parameter estimates.

The intraclass correlation is often described as the proportion of variability in the dependent variable that is explained by the group membership (Montgomery, 1997). In the typical multi-level data structure, one or more characteristics are measured on individuals in each of several groups, and one or more characteristics are also obtained on these same individuals at the group level. That is, each individual in the same group will have the same value for the group level characteristic(s). Hence, differences among the group-level characteristic(s) can account for some of the variation in the responses, and this variation is referred to as the intraclass correlation. Murray (1998) also refers to this quantity as a clustering effect.

It is conceivable that multi-level data in which little or no differences exist among the group-level characteristic(s) (i.e., little or no intraclass correlation) will produce parameter estimates that are very similar, if not identical, when analyzed using HLM and MLR. But when the intraclass correlation is greater, i.e., larger differences among group-level characteristics, the two analyses will produce parameter estimates that exhibit larger differences among them. The purpose of this study is to compare the parameter estimates obtained from HLM with those obtained from MLR on simulated data in which the intraclass correlation is systematically varied from small values to larger ones.

Method

In this simulation study the outcome of interest was the difference between the vector of parameter estimates obtained when data are analyzed using both MLR and HLM. The independent variable was the size of the intraclass correlation among the groups. Our objective was to use ten different values of the intraclass correlation: 0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, and 0.9. However, because we were unable to generate directly data sets with a given intraclass correlation, we were forced to “back into” these values by generating data with given correlations and then checking the intraclass correlation. Trial and error and the ability to learn while doing allowed us to become more proficient as the process progressed. However, we were unable to fix the intraclass correlation at these specified values and instead were forced to settle for intraclass correlations within a small, specified range. Therefore, the ranges of intraclass correlations examined in this study were as follows: < .05, between .05 and .15, between .15 and .25, and so on through between .85 and .95. Within each of these intraclass correlation ranges, 15,000 data sets were generated.

The HLM model used in each of these data sets was a simple two-level model with one individual-level variable and one group-level variable. This model can be expressed as:

$$Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + r_{ij},$$

where Y_{ij} represents the response for the i^{th} individual in the j^{th} group

β_{0j} represents the y -intercept of the regression line for the j^{th} group

β_{1j} represents the slope of the regression line for the j^{th} group

X_{ij} is the measurement on the individual-level variable for the i^{th} individual in the j^{th} group,

r_{ij} represents random error associated with the response for the i^{th} individual in the j^{th} group, and j ranges from 1 to J , the number of groups in the data set.

In the HLM model, however, the group parameters, β_{0j} and β_{1j} , are not estimated individually from the raw individual data, but instead are estimated from a second-level model using group-level data. This model can be expressed as:

$$\beta_{kj} = \gamma_{k0} + \gamma_{k1}W_j + u_{kj},$$

where β_{kj} represents the k^{th} parameter for the j^{th} group

γ_{k0} represents the y -intercept of the regression line for the k^{th} parameter

γ_{k1} represents the slope of the regression line for the k^{th} parameter

W_j represents the measurement on the group-level variable for the j^{th} group, and

u_{kj} represents random error associated with the k^{th} parameter for the j^{th} group.

In this model, there are two second-level models, one for the y -intercept, β_{0j} , and one for the slope, β_{1j} :

$$\beta_{0j} = \gamma_{00} + \gamma_{01}W_j + u_{0j} \quad \text{and} \quad \beta_{1j} = \gamma_{10} + \gamma_{11}W_j + u_{1j}.$$

The MLR model examined in this study includes both the individual-level variable and the group-level variable along with their interaction. This model, with two independent variables and their interaction term, can be expressed as:

$$Y_{ij} = \beta_0 + \beta_1 X_{1ij} + \beta_2 X_{2ij} + \beta_3 X_{1ij}X_{2ij} + e_{ij},$$

The number of groups was set to ten with the number of individuals in each group allowed to vary, but the mean group size was approximately 50. Each data set contains 500 observations on the response variable, 500 observations on the individual-level variable, and approximately 50 observations in each of the ten groups on the group-level variable. Hence, the parameter estimate vectors for multiple regression contained four values, an intercept, a coefficient for the individual-level variable, a coefficient for the group-level variable, and a coefficient for the “interaction” between the individual-level variable and the group-level variable, $[\beta_0 \beta_1 \beta_2 \beta_3]'$. For the HLM model, the parameter estimate vectors contained the estimates of the second-level fixed effects parameters (often referred to as “gammas”), $[\gamma_{00} \gamma_{10} \gamma_{01} \gamma_{11}]'$.

Once the 15,000 data sets were produced within each intraclass correlation range, each set was analyzed using both a hierarchical linear model and a multiple linear regression model to produce a vector of parameter estimates. All of the analyses were run using SAS, with the MLR analyses performed with PROC GLM and PROC MIXED used for the HLM analyses. Wang (1997) demonstrated the similarities in results between multi-level analyses performed in SAS PROC MIXED and the HLM software (e.g., Raudenbush, Bryk, and Congdon, 1999). For each data set, the parameter estimate vectors from MLR and HLM were subtracted to obtain a difference vector. Consecutive sets of 50 difference vectors were grouped together to form mean difference vectors and Hotelling's T^2 statistic was used to determine the proportion of mean difference vectors that exceeded an F -critical value at the 5% significance level.

Results

The parameter estimate vectors from both HLM and MLR analyses on each of the 15,000 data sets showed remarkable similarity when compared with each other. Once the parameter estimate vectors were obtained and subtracted to form difference vectors, the average difference vectors were compared to a null hypothesis of a zero vector which would be indicative of no difference between HLM and MLR in terms of the parameter estimates. These mean difference vectors were tested using Hotelling's T^2 statistic. If this null hypothesis were true, we would expect to find about 5% of the mean difference vectors differing from the zero vector, simply by chance variation. Larger percentages of the mean difference vectors differing from the zero vector would indicate that HLM and MLR produce parameter estimates, which differ significantly from one another.

For each of the ten intraclass correlation ranges, the results are displayed separately in Table 1. Notice that none of the intraclass correlation ranges had F -approximations that approach the expected 5% Type I error criteria. Only when the intraclass correlation exceeded 0.65 did the percentage of F 's exceeding the 5% critical value surpass even 1%. In fact, these pairs of parameter estimates are so similar to each other that in more than half of the intraclass correlation ranges studied, less than 1% of the mean difference vectors differed significantly from each other. There appears to be little evidence in these data that the size of the intraclass correlation has any influence upon the difference between the parameter estimates from HLM and MLR.

Discussion

In this study, numerous data sets were generated and analyzed to investigate the effect, if any, that the intraclass correlation has on the parameter estimates of multiple linear regression as compared to those from hierarchical linear modeling. Previous work had indicated that in some cases these estimates were very similar, if not identical, and that in other cases, the parameter estimates from these two procedures were quite different. It was conjectured that in those data that produced different estimates of the parameters, perhaps it was a larger intraclass correlation that could account for these estimates differing. These results would seem to indicate that it is not the size of the intraclass correlation that is responsible for differences in parameter estimates between HLM and MLR. There can be no doubt that in some instances these estimates do indeed differ. Why they differ in some cases and not in others is still a question that needs answering. It would appear, however, that the intraclass correlation could be eliminated from the list of possible explanations for these differences.

Table 1. Percentages of Calculated F -Values Beyond the 5% F -critical value.

Intraclass Correlation Range	Percentage of F 's beyond the 5% F Critical Value
0.00-<0.05	0.9
0.05-<0.15	0.3
0.15-<0.25	0.9
0.25-<0.35	0.6
0.35-<0.45	0.0
0.45-<0.55	0.0
0.55-<0.65	0.0
0.65-<0.75	1.3
0.75-<0.85	1.3
0.85-<0.95	1.1

It should be noted that the MLR model that was compared to HLM in this study is the model that contains the interaction term between the individual-level variable and the group-level variable. Previous work by Mundfrom and Schultz (2001) indicated that the simpler MLR model without interactions was not an adequate competitor to HLM in terms of similar predicted values or parameter estimates. Also, in this study, only the simplest full multi-level model was investigated. It could still be the case that more complex multi-level models may show more influence due to the size of the intraclass correlation in terms of differences between parameter estimates from HLM and MLR. It should also be noted that HLM has been shown to provide better, more accurate (i.e., larger) estimates of standard errors of parameter estimates

than does multiple regression. Indeed, it is this characteristic of HLM that provides its premier advantage in the analysis of multi-level data. This study did not compare standard errors of the two techniques.

This study did not assume a balanced design (i.e., equal sample sizes per group) although further research on the effect of unbalanced designs on the respective parameter estimates may be warranted. Similarly, the role of heterogeneous variances in the comparison of these parameter estimates may also be worthy of further investigation.

Finally, it is not the intent of this study to imply that hierarchical linear modeling is in any way inadequate or inappropriate for use in analyzing data, particularly multi-level data. The intent here is simply to investigate differences and similarities between the results obtained when these two procedures are used with the same multi-level data sets, and to see if reasons can be identified as to why, or in what situations, these procedures produce different outcomes.

References

- Bryk, A. S. & Raudenbush, S. W. (1992). *Hierarchical linear models: Applications and data analysis methods*. Newbury Park, CA: Sage Publications.
- Goldstein, H. (1987). *Multilevel models in educational and social research*. London: Griffin.
- Littel, R. C., Milliken, G. A., Stroup, W. W., & Wolfinger, R. D. (1996). SAS system for mixed models. Cary, NC: SAS Institute, Inc.
- Montgomery, D. C. (1997). *Design and analysis of experiments* (4th edition). New York: Wiley.
- Mundfrom, D. J. & Schultz, M. R. (2001). A comparison between hierarchical linear modeling and multiple linear regression in selected data sets. *Multiple Linear Regression Viewpoints*, 27(1), 3-11.
- Murray, D. M. (1998). *Design and analysis of group-randomized trials*. New York: Oxford University Press.
- Raudenbush, S. W. & Bryk, A. S. (1986). A hierarchical model for studying school effects. *Sociology of Education*, 59, 1-17.
- Raudenbush, S. W., Bryk, A. S., & Congdon, R. (1999). *HLM Version 5*. Chicago IL: Scientific Software International.
- Wang, J. (1997). Using SAS PROC MIXED to demystify the Hierarchical Linear Model. *Journal of Experimental Education*, 66, 84-93.

Send correspondence to: Daniel J Mundrom, Department of Applied Statistics and Research Methods, University of Northern Colorado, Greeley, Colorado 80639. Email: Daniel.Mundfrom@unco.edu.

Scoring Above the International Average: A Logistic Regression Model of the TIMSS Advanced Mathematics Exam

James B. Schreiber

Southern Illinois University Carbondale

This study examined 2349 advanced mathematics students from the Third International Mathematics and Science Study. Students were separated into two groups: those that scored above the international average of 501 and those who scored below. Logistic regression was then utilized to model the student data. Results indicate that students whose parents had less than a high school education were one-fourth as likely to score above the international average and students who were enrolled in advanced mathematics and physics courses were three times as likely to score above the international. A probability model is also discussed.

The intense concern over the mathematics achievement levels of U. S. elementary and high-school students can be seen through the number of articles published (Baker, 1993a, 1993b; Bracey, 1991, 1992, 1993; Freudenthal, 1975, Rotberg, 1990; Stedman, 1994a, 1994b). A great deal of debate has centered on the fact that the U.S. scored lower than many countries on most of the assessments. Very little discussion has focused on U.S. high school seniors and on the simultaneous examination of factors that are associated with scoring above the international mean. The questions, which motivate this study, are: What student variables are associated with scoring above the international mean on the TIMSS advanced mathematics test? How much of an impact do these variables have?

Theoretical Framework

Over the past few decades a great deal of research has been conducted on the student level factors that impact mathematics achievement. The positive relationship between attitude (ATT) and achievement (ACH) has long been observed (Ethington & Wolfe, 1984, 1986; Lester, Garafalo, & Kroll, 1989; Ma, 1997; Suydam & Weaver, 1975). Ma and Kishnor (1997) conducted a meta-analysis concerning the effect of mathematical self-concept and achievement in mathematics (ACH) and found a significant effect size of .23 (statistically different than zero) for self-concept about mathematics and ACH. The authors concluded that a positive self-concept about mathematics is associated with higher achievement in mathematics. Ma (1997) observed that attitudes towards mathematics (e.g., importance, difficulty, enjoyment) influenced achievement. As crucial as attitude towards mathematics is for mathematics achievement, understanding a student's academic related beliefs is also an important key to comprehending the student's achievement level (Dweck, & Elliot, 1983; Felson, 1984). Academic beliefs have been observed to impact the achievement level of students (Dweck & Elliot, 1983) and Schoenfeld (1985) has pointed out that students' beliefs about mathematics may impede their ability to solve problems.

Extracurricular activities (EA) for students have been glorified and chastised over the years (Gerber, 1996). Two main views of EA exist. The first view is that extracurricular activity is similar to a zero sum game where greater activity will subvert academic achievement (Coleman, 1961; Marsh, 1992). The second view is that EA experiences "further the total development of the students," thereby enhancing non-academic goals and possibly facilitating academic goals (Holland & Andre, 1987; Marsh, 1992).

A common extracurricular activity is athletics (e.g., football, softball, baseball, soccer). A review of research by Holland and Andre (1987) focused on the examination of the relationship between athletic participation and achievement. They reported that the research demonstrated that male high school athletes received somewhat higher grade-point averages (GPA's) than did non-athletes. When one considers standardized achievement or aptitude tests, males, whose only after school activity is sports, scored lower than the national average on the Standardized Achievement Test. No significant difference in either grade-point average or standardized test scores was observed between female athletes and non-athletes.

The influence of part time employment, another common after school activity for American high school students, on achievement has shown mixed results in the research literature. Some studies indicate

that part-time employment has a negative influence on achievement (e.g., Brown & Steinberg, 1991; Cooper et al., 1999); whereas, other studies show a positive influence (D'Amico, 1984), or no influence (Green & Jacquess, 1987). Cooper et al. (1999) reported that the number of hours per week had a significant inverse association with standardized test scores ($r=-.29$) and with teacher assigned grades ($r=-.17$). With 51.6% of high school senior female students and 47.5 percent of high school senior male students working, the examination of the effects of employment on achievement is critical (Green, Dugoni, Ingels, & Camburn, 1995).

With regard to homework, Cooper (1989) observed a positive linear relationship between hours per week spent on homework (0 to 10 hours) and achievement. Television has traditionally been assumed to lessen achievement (Comstock, 1991; Keith et al., 1986). Simply, television viewing displaces academic activities and reduces the amount of time available for completing homework and other academic activities, thereby reducing achievement. Keith et al. (1986) observed a small but negative relationship between the amount of television watched and achievement. In a review of research, Williams, Haertel, Haertel, and Walberg (1982) observed a strong negative effect after ten hours per week. The negative effect increases and is extreme after thirty hours per week. A more recent study by Cooper et al. (1999) observed a significant negative association between achievement and television viewing (mean viewing was 1-2 hours per night).

Differences in mathematics achievement for male and female students have been observed (Anderson, 1989; Sherman, 1987). Currently there is still debate about the observed differences. Some researchers have observed a large difference (Benbow & Stanley, 1980) and some no difference (Radhawa, Beamer, & Lundberg 1993). The more consistent finding in the research appears to be that some differences still exist but not in all areas of mathematics (Fennema & Carpenter, 1981).

Finally, numerous research articles have been written concerning family background variables such as parental socio-economic status (SES) or parent education (Keith & Cool, 1992; Lockheed & Komenan, 1989; McConeghy, 1987; Santiago & Okley, 1992). The positive relationship between parental SES or parent education level and achievement has been consistently observed (Goleman, 1988; Heyns, 1978). Similar results for SES have been observed in multilevel modeling (Lee & Bryk, 1989).

Methodology

Sample

The sample for this study consists of U.S. students from the Third International Mathematics and Science Study (TIMSS) Population 3 final year of secondary school cohort (High School Seniors) who were administered the advanced mathematics test. All students who were administered the advanced mathematics instrument were designated advanced mathematics students or advanced mathematics and physics students. Total sample size is 2349 with 1158 girls and 1191 boys.

Variables

A listing of each variable, TIMSS label, and coding information is provided in Appendix A. The dichotomous outcome variable is based on the first plausible value of the advanced mathematics test. Students who scored over the international average of 501 were coded as one, and students who scored below the international average were coded as zero. Though dichotomizing typically reduces statistical power, with all the variables being statistically significant in the logistic regression model this concern is reduced. Scoring above the international average is not the primary goal of mathematics education in the United States, just a component of this study. The justification for this dichotomous split is due to the fact many of the articles concerning international data discuss those who score above and below the international average (e.g., Carson, Heulkamp, & Woodall, 1993, Garden, 1989; Mullis, Martin, Beaton, Gonzales, Kelly, & Smith, 1998) but do not discuss in detail the combination of factors that are associated with this separation. Secondly, the author was interested in creating a probability model, which is possible with logistic regression.

Table 1. Frequency Tabulations of Dichotomous Variables

Variable	Below International Average	Above International Average
Gender		
Boy (0)	816	375
Girl (1)	918	240
Advanced Math/Physics		
Advanced Math (0)	914	150
Advanced Math and Physics (1)	820	465
^a Parent Education 1 (High School)		
No = 0	840	458
Yes = 1	753	135
Parent Education 2 (Elem. School)		
No = 0	1446	573
Yes = 1	147	20

Note: *a.* 1131 students had parents who graduated from a university, 888 from high school, 167 only elementary school.

Table 2. Descriptive Statistics for Continuous Variables

Variable	Below International Average	Above International Average
Attitude	10.53 (2.94)	8.89 (2.45)
Natural Talent	2.47 (0.75)	2.20 (0.76)
Hard Work	1.58 (0.66)	1.86 (0.78)
Television	2.76 (0.96)	2.58 (0.94)
Employment	2.66 (1.68)	2.20 (1.46)
Sports	2.36 (1.21)	2.28 (1.14)
Studying Math	2.14 (0.79)	2.26 (0.73)

Note. Standard Deviations (SD) are in parantheses.

The independent variables are categorized into three areas:

Background of Student consisted of: Parent Education (dummy coded), Gender (Girls), and Advanced Math/Physics;

After School Time consisted of: Television Viewing, Employment, Sports, and Studying Math;

Affective Factor was composed of: Attitude, Natural Talent Belief, Hard Work Belief.

Analysis

To answer the research questions posed, a logistic regression analysis was conducted and the variables of interest were run in blocks. The first block contained background variables: Parent Education, Gender (Girls), Math/Physics Expert. The second block contained the first block and affective factors: Attitude, Natural Talent, Hard Work. The final block included the previous two blocks and after school time variables: Television, Employment, Sports, Studying Math.

Results

Tables 1 and 2 provide descriptive statistics of the variables of interest separated by scoring above or below the international average. Table 3 provides results from the from the logistic regression analysis. Block 1 results indicate that parental education, gender (girls) and advanced mathematics and physics were significantly related to scoring above international mean. Taking the exponential of the coefficient for advanced mathematics and physics ($e^{1.105}$) provides an odds ratio of 3.019, which indicates that students enrolled in advanced mathematics and physics were three times as likely to score above the international average than students only enrolled in advanced mathematics. Students whose parents had only completed elementary school ($e^{-1.368}$) were only one-fourth as likely to score above the international mean. The Hosmer and Lemeshow statistic indicates that the model was a good fit. The Cox and Snell and Nagelkerke R^2 values indicate less than 20 percent of the “variance” or likelihood was accounted with the Block 1 variables. These values are analogous to the traditional R^2 in multiple regression, but are specific to logistic regression due to the dichotomous dependent variable.

Table 3. Logistic Regression Results for the Three Blocks

Factor	Block 1			Block 2			Block 3		
	<i>b</i>	<i>SE</i>	<i>p</i>	<i>b</i>	<i>SE</i>	<i>p</i>	<i>b</i>	<i>SE</i>	<i>p</i>
High School	-1.17	0.17	<.001	-1.31	0.12	<.001	-1.25	0.12	<.001
Elem. School	-1.36	0.25	<.001	-1.69	0.28	<.001	-1.75	0.28	<.001
Gender (Girl)	-0.37	0.10	<.001	-0.20	0.11	.063	-0.39	0.12	.001
Advanced Math & Physics	1.10	0.11	<.001	0.97	0.12	<.001	0.91	0.12	<.001
Attitude				-0.26	0.02	<.001	-0.25	0.02	<.001
Natural Talent				-0.45	0.08	<.001	-0.45	0.08	<.001
Hard Work				0.39	0.08	<.001	0.49	0.08	<.001
Televisions							-0.16	0.06	.006
Job							-0.17	0.05	<.001
Sports							-0.16	0.04	.001
Math Study							0.27	0.08	.001
-2log likelihood	2260.85			2007.13			1938.44		
Hosmer-Lemeshow χ^2	5.82 <i>df</i> = 6			8.66 <i>df</i> = 8			5.968 <i>df</i> = 8		
Cox & Snell R ²	.126			.214			.231		
Nagelkerke R ²	.183			.309			.335		

The results from the Block 2 analysis indicate that all the variables, except Gender, were significantly related to scoring above the international mean. This indicated the possibility of an interaction between Gender and the affective variables (Edward St. John, Personal Communication). Since the interaction of affective variables and gender are well documented (McLeod, 1992), intermediate models were run (see Table 4) to examine the association of interactions on the model (Gender*Attitude, Gender*HardWork, Gender*Natural Talent). The interactions were not significant, the $-2\log$ likelihood did not significantly change, and the R² estimates were not significantly altered. According to Hosmer and Lemeshow (1989) this indicates that affect may be a confounder but not an effect mediator.

Also, after including affect variables – in essence statistically controlling for the association of affect – changes in the negative association of formal parent education changes were evident. The coefficients for both variables increased. The impact of being enrolled in advanced mathematics and physics had a decrease in association from Block 1 to Block 2. For the continuous variable Attitude, the results indicate the worse the students attitude the less likely the student scored over the international average. The coefficient for Natural Talent indicates that the more a student believed that natural talent was a key to mathematics success the more likely the student scored above the international average. Finally, the less a student believed in hard work as a key to mathematics success the more likely the student scored over the international average.

The $-2\log$ likelihood value dropped from Block 1 to Block 2 and the Hosmer-Lemeshow statistic indicates that the Block 2 model is a better fit for the data. Finally, The R² estimates indicate more “variance” was accounted for by the inclusion of the Affect variables, which was expected.

Block 3 results indicate that the variables from Blocks 1 and 2 and the new variables are significantly associated with scoring above the international mean. Taking the exponential of the Studying Math coefficient ($e^{.27}$) shows that for every one unit increase on the scale the odds of scoring above the international average increases by 1.3. Television viewing, employment and sports participation decreased the likelihood of scoring above the international mean. For example, for every unit increase in sports participation the probability of scoring over the international average decreases by a factor of 0.85 ($e^{-.16}$). In essence decreasing the probability of scoring over the international average.

Table 4. Examination of Gender and Affect Interaction

Interaction	<i>b</i>	<i>SE</i>	-2Log Likelihood	Cox & Snell R ²	Nagelkerke R ²
Gender*Attitude	0.03	0.04	2006.66	.214	.310
Gender*Natural Talent	-0.17	0.15	2005.79	.214	.310
Gender* Hard Work	-0.18	0.16	2005.82	.214	.310

Table 5. Interaction Effects of Gender by Extracurricular Activity.

Interaction	<i>b</i>	<i>SE</i>	-2Log Likelihood	Cox & Snell R ²	Nagelkerke R ²
Gender*Television	0.002	0.112	1938.44	0.231	0.335
Gender*Job	0.050	0.073	1937.97	0.233	0.331
Gender*Sports	0.223	0.336	1937.32	0.232	0.336
Gender*Studying	0.028	0.154	1938.41	0.231	0.335

From Block 2 to Block 3 the variable Gender went from not significant to significant indicating a possible interaction with extracurricular activities. Therefore, these interactions (Gender*Television, Gender*Employment, Gender*Sports, Gender*Studying Math) were examined. None of the interactions were significant and the interactions did not significantly change the -2log likelihood or the R² values (see Table 5).

Though mostly completed with epidemiological research, a probability equation can be developed to examine a student's probability of scoring above the international average. The equation is:

$$\pi(X) = \exp(A) / [1 + \exp(A)], \text{ where}$$

$$A = 2.483 - 1.252(\text{HS}) - 1.756(\text{ES}) - .392(\text{Girls}) + .913(\text{AMP}) - .256(\text{Att})$$

$$- .457(\text{NT}) + .499(\text{HW}) - .166(\text{TV}) - .163(\text{Sports}) - .172(\text{Job}) + .272(\text{MStdy}).$$

A boy whose parents graduated from college, is not enrolled in advanced mathematics and physics, has a good attitude towards mathematics and average score on beliefs, watches TV, works, plays sports and studies a couple of hours per day, would have a value of 0.48 indicating a 48 percent probability of scoring above the international average. Keeping everything the same except making the student enrolled in both advanced mathematics and physics increases the value to .74. Using an equation such as this, though, is only useful when the variables can be controlled, such as working and studying.

Discussion

With the publication of the international scores and the discussion about the rank of American students average score it was important to look at factors that were associated with scoring above the international mean. The advanced mathematics and advanced mathematics and physics students were chosen because there has been less emphasis on their performance. By using such a unique group of students it was hoped by the author that traditional variables would not have a strong impact on scoring above or below the international average. Unfortunately, formal parent education level and gender were significantly associated with scoring above the international mean showing that this traditionally observed disparity still exists. This indicates that there is still a great deal of work to be done to close this gap. The results indicate a more in-depth examination of this is needed. For example, a study of classrooms of advanced mathematics students with varying levels of SES and gender followed for a specified time period may help to understand the factors associated with these disparities.

The result for advanced mathematics and physics students was expected. Students who are concurrently enrolled in mathematics and physics courses see many of the same equations and solve many of the same types of problems. The dual coverage of content, amount of time with the material, and similar activities in the courses provides the students more opportunity to learn the content. Therefore,

advanced mathematics and physics students are potentially spending twice as much time during each day with advanced mathematics problems. Another aspect that may be influencing the difference is an abstract to concrete representation of the material (Danesi, 1993). A great deal of work completed in advanced mathematics courses may appear to students to be abstract. In the physics courses those students may be seeing concrete representations of the abstract ideas from their mathematics course. The connection may have assisted the students who are enrolled in both course score above the international average. But this too needs to be researched further to examine if this could be contributing to the observed difference or something else, such as students who are concurrently enrolled in both courses have more “ability” or a higher “IQ.”

The affective factors, though important to understand and monitor are difficult to change. Being a former mathematics instructor, the author understands the difficulty in attempting to change attitudes and beliefs that have developed over many years. More importantly, though, the results demonstrate that our most advanced students are impacted in a similar way as the rest of our mathematics students (Schreiber, 2000). With regard to beliefs, the data indicate that (due to reverse coding) those students who believe that success in mathematics is due to natural talent were more likely to score above the international mean. This is in contrast to work by Dweck (1986) and Schommer (1990). Traditionally those students who have believe that ability or intelligence is fixed tend to perform lower and not persist in difficult tasks (Dweck, 1986) or, avoid difficult tasks if they think their ability is low (Dweck, 1986). Further, Schommer (1990) observed that students who believe success is tied to natural talent or ability tend to do worse academically. One explanation for this observation is the uniqueness of the cohort and the may not have experience much failure in the mathematics domain. This may be part of the reason for the development of this belief. Recently, Dai, Moon, and Feldhusen (1998) observed that gifted and talented students that have high levels of self-efficacy attribute success to both high ability and hard work. Overall, affective factors are associated with performance and need to be considered.

Time variables, overall, may be more readily modified than affective factors. Their association is smaller in the model than other variables, but still significant. As students engage in more hours per day in these activities they may positively (studying) or negatively (TV, job, sports) impact performance. Scoring above the international average in this study is associated with less time in non-academic activities and more time studying. The issue is not to remove these activities because they do provide a great deal of societally beneficial effects (time management, responsibility, working as a group, school retention), but to monitor the amount of time and energy these activities are absorbing. Even within these observations, more in-depth research needs to be conducted on the effects of different types of activities (e.g., school/non school, academic/non-academic, organized/unorganized) on achievement.

In addition to the importance of the findings, this study also demonstrates that logistic regression can be a useful tool in modeling educational data because it allows for the examination of factors related to academic achievement based on a criterion cut off point. Finally, the reader is reminded that these data are from a very select group of students and only one model or view of the factors associated with advanced mathematics achievement was examined, therefore, these results should be weighed with previous observations and other research findings.

References

- Anderson, J. (1989). Sex related differences on objective tests among undergraduates. *Educational Studies in Mathematics*, 20, 165-177.
- Baker, D. P. (1993a). Compared to Japan, the U.S. is a low achiever really: New evidence and comment on Westbury. *Educational Researcher*, 22 (3), 18-20.
- Baker, D. P. (1993b). A rejoinder. *Educational Researcher*, 22 (3), 25-26.
- Benbow, C. P., and Stanley, J. C. (1980). Sex differences in mathematical ability: Fact of artifact? *Science*, 12, 1262-1264.
- Bracey, G. W. (1991). Why can't they be like we were? *Phi Delta Kappan*, October, 105-117.
- Bracey, G. W. (1992). The second Bracey report on the condition of public education. *Phi Delta Kappan*, 74, 104-117.

- Bracey, G. W. (1993). The third Bracey report on the condition of education. *Phi Delta Kappan*, 75, 105-117.
- Brown, B. B., & Steinberg, L. (1991). *Noninstructional influences on adolescent engagement and achievement. Final Report Project 2*. Madison, WI: National Center on Effective Secondary Schools. (ERIC Document Reproduction No. ED340641).
- Carson, C. C., Heulskamp, R. M., & Woodall, R. D. (1993). Perspectives on education in American. *The Journal of Educational Research*, 86, 259-310.
- Coleman, J. S. (1961). *The adolescent society*, New York, NY: Free Press of Glencoe.
- Comstock, G. (1991). *Television and the American child*. New York: Academic Press.
- Cooper, H. (1989). *Homework*. New York: Longman.
- Cooper, H., Valentine, J. C., Nye, B., & Lindsay, J. J. (1999). Relationships between five after-school activities and academic achievement. *Journal of Educational Psychology*, 91 (2), 369-378.
- Dai, Y. D., Sidney, M. M., & Feldhusen, J. F. (1998). Achievement motivation and gifted students: A social cognitive perspective. *Educational Psychologist*, 33 (2/3), 45-63.
- D'Amico, R. (1984). Does employment during high school impair academic progress? *Sociology of Education*, 57, 152-164.
- Danesi, M. (1993). *Messages and meaning: An Introduction to semiotics*. Toronto, ON: Canadian Scholar's Press.
- Dweck, C. S. (1986). Motivational processes affecting learning. *American Psychologist*, 41 (10), 1040-1048.
- Dweck, C. S., & Elliot, E. S. (1983). Achievement motivation. In E. M. Hetherington (Ed.), P. H. Mussen (Series Ed.), *Handbook of child psychology: Vol. 3. Socialization, personality, and social development*, 4th ed, (643-691). New York, NY: Wiley.
- Ethington, C. A., & Wolfe, L. M. (1984). Sex differences in a causal model of mathematics achievement. *Journal for Research in Mathematics Education*, 15, 361-377.
- Ethington, C. A., & Wolfe, L. M. (1986). A structural model of mathematics achievement for men and women. *American Educational Research Journal*, 23, 65-75.
- Felson, R. B. (1984). The effect of self-appraisal of ability on academic performance. *Journal of Personality and Social Psychology*, 47, 944-952.
- Fennema, E., & Carpenter, T. P. (1981). Sex related differences in mathematics: Results from a national assessment. *Mathematics Teacher*, 74, 554-559.
- Freudenthal, H. (1975). Pupil's achievement internationally compared—The IEA. *Educational Studies in Mathematics*, 6, 127-186.
- Garden, R. A. (1989). Students' achievements: Population be. In D. F. Robitaille and R. A. Garden (Eds.), *The IEA study of mathematics II: Contexts and Outcomes of School Mathematics*. (126-152). New York: Pergamon Press.
- Gerber, S. B. (1996). Extracurricular activities and academic achievement. *Journal of Research and Development in Education*, 30 (1), 42-50.
- Goleman, D. (1988, April 10). An emerging theory on blacks' I.Q. scores. *New York Times* (Education Life Section), 22-24.
- Green, G., & Jacquess, S. N. (1987). The effect of part-time employment on academic achievement. *Journal of Educational Research*, 80, 325-329.
- Green, P. J., Dugoni, B. L., Ingels, S. J., & Camburn, E. (1995). *A profile of the American high school senior in 1992. National Educational Longitudinal Study of 1988. Statistical Analysis Report*. National Opinion Research Center. (ERIC Document No. ED 386 502)
- Heyns, B. (1978). *Summer learning and the effects of schooling*. New York: Academic Press.
- Holland, A., & Andre, T. (1987). Participation in extracurricular activities in secondary school: What is known, what needs to be known? *Review of Educational Research*, 57, 437-466.
- Keith, T. Z., Reimers, T.M., Fehrman, P.G., Pottebaum, S.M., & Aubey, P.G. (1986). Parental involvement in homework and TV time: Direct and indirect effects on high school achievement. *Journal of Educational Psychology*, 78, 373-380.

- Keith, T.Z. & Cool, V. A. (1992). Testing models of school learning: Effects of quality of instruction, motivation, academic coursework, and homework on academic achievement. *School Psychology Quarterly*, 7 (3), 207-226.
- Lee, V. E., & Bryk, A. S. (1989). A multilevel model of the distribution high school achievement. *Sociology of Education*, 62 (3), 172-192.
- Lester, F. K., Garofalo, J., & Kroll, D. L. (1989). Self-confidence, interest, belief, and metacognition: Key influences on problem-solving behavior. In D. B. McLeod & V.M. Adams (Eds.), *Affect and mathematical problem solving: A new perspective* (75-88). New York: Springer-Verlag.
- Lockheed, M.E. and Komenan, A. (1989). Teaching quality and student achievement in Africa: The case of Nigeria and Swaziland. *Teacher and Teacher Education*, 5 (2), 93-113.
- Ma, X. (1997). Reciprocal relationships between attitude toward mathematics and achievement in mathematics. *The Journal of Educational Research*, 90 (4), 221-229.
- Ma, X., & Kishnor, N. (1997). Attitude toward self, social factors, and achievement in mathematics: A meta-analytic review. *Educational Psychology Review*, 9 (2), 89-120.
- Marsh, H. W. (1992). Extracurricular activities: Beneficial extension of the traditional curriculum or subversion of academic goals? *Journal of Educational Psychology*, 84 (4), 553-562.
- McConeghy, J. I. (1987). Mathematics attitudes and achievement: Gender differences in a multivariate context. ERIC Services Document 284742.
- McLeod, D. B. (1992). Research on affect in mathematics education: A reconceptualization. In D. A. Grouws (Ed.), *Handbook of research on mathematics teaching and learning* (pp. 575-596). New York: Macmillan.
- Mullis, I. V. S., Martin, M. O., Beaton, A. E., Gonzales, E. J., Kelly, D. L., & Smith, T. A. (1998). *Mathematics and Science Achievement in the Final Year of Secondary School*. Chestnut Hill, MA: Boston College.
- Radhawa, B. S., Beamer, J. E. & Ingvar, L. (1993). Role of mathematics self-efficacy in the structural model of mathematics achievement. *Journal of Educational Psychology*, 85 (1), 41-48.
- Rotberg, I. (1990). I never promised you first place. *Phi Delta Kappan*, 72, 296-303.
- Santiago, R. W., & Okley, J.R. (1992). The effects of advisement and locus of control on achievement in learner controlled instruction. *Journal of Computer-Based Instruction* 19 (2), 47-53.
- Schoenfeld, A.H. (1985). *Mathematics problem solving*. Orlando, FL: Academic Press.
- Schommer, M. (1990). The effects of beliefs about the nature of knowledge on comprehension. *Journal of Educational Psychology*, 82, 498-504.
- Schreiber, J. B. (2000, April). *High-school seniors' mathematics literacy*. A poster presented at the annual meeting of the American Educational Research Association, New Orleans.
- Sherman, J. (1987). *Sex related cognitive differences*. Springfield, IL: Wiley.
- Stedman, L. C. (1994a). The sandia report and U.S. achievement: An assessment. *Journal of Educational Research*, 87 (3), 133-146.
- Stedman, L. C. (1994b). Incomplete explanations: The case of U.S. Performance in the international assessments of education. *Educational Researcher*, 23 (7), 24-32.
- Suydam, M. N., & Weaver, J. F. (1975). Research on mathematics learning. In J. N. Payne (Ed.), *Mathematics learning in early childhood: Thirty seventh yearbook* (44-67). Reston, VA: National Council of Teachers of Mathematics.
- Williams, P. A., Haertel, E. H., Haertel, G. D., & Walberg, H. J. (1982). The impact of leisure-time television on school learning. *American Educational Research Journal*, 19, (1), 19-50.

Send correspondence to: James B. Schreiber, Department of Educational Psychology and Special Education, Southern Illinois University, Mail Code 4618, Carbondale, IL 62901
 Email: jschreib@siu.edu

Appendix A

The following variables were either selected from those in the TIMSS Population 3 data set or were developed using other variables from the data set.

1. **Gender:** ITSEX. 1 = Girls; 0 = Boys
2. **Parent Education:** CSDGEDU. This variable was dummy coded into:

	Parent Education 1	Parent Education 2
University Graduate	0	0
High School Graduate	1	0
Primary School Graduate	0	1

3. **Advanced Math/Physics:** IDSUBPOP Recoded 1 = Adv. Math and Physics, 0 = Advanced Math
4. **Attitude (Composite):** CSBENJY, CSBMBORE, CSBMEASY, CSBMLIFE, CSBMLIKE. For each individual variable 1 = strongly agree to 4 = strongly disagree. CSBMBORE and CSBMLIKE were reverse coded to match the direction of the others. The variables were added together to create the composite. The lower the attitude score, the better the attitude towards mathematics in general. Alpha = .8241
5. **Natural Talent:** CSBMDOW1. 1 = strongly to agree 4 = strongly disagree. Those who chose strongly agree indicated that natural talent was the key to mathematics success.
6. **Hard Work:** CSBMDOW3. 1 = strongly agree to 4 = strongly disagree. Those who chose strongly agree indicated that hard work was the key to mathematics success.
7. **Television:** CSBGDAY1. 1 = No time per school day, 2 = Less than one hour per school day, 3 = 1-2 hours per school day, 4 = 3-5 hours per school day, 5 = more than five hours per school day.
8. **Sports:** CSBGDAY6. 1 = No time per school day, 2 = Less than one hour per school day, 3 = 1-2 hours per school day, 4 = 3-5 hours per school day, 5 = more than five hours per school day.
9. **Employment:** CSBGDAY5. 1 = No time per school day, 2 = Less than one hour per school day, 3 = 1-2 hours per school day, 4 = 3-5 hours per school day, 5 = more than five hours per school day.
10. **Studying Math:** CSBGDAY. 8 1 = No time per school day, 2 = Less than one hour per school day, 3 = 1-2 hours per school day, 4 = 3-5 hours per school day, 5 = more than five hours per school day.
11. **Television:** CSBGDAY1. 1 = No time per school day, 2 = Less than one hour per school day, 3 = 1-2 hours per school day, 4 = 3-5 hours per school day, 5 = more than five hours per school day.

A Discussion of an Alternative Method for Modeling Cyclical Phenomena

Russel Brown
Summa Health Care

Isadore Newman
University of Akron

The purpose of this paper is to provide an answer to the question of the relative effectiveness of the cosine function versus a polynomial function in the description and stability of prediction of a specific set of longitudinal data. If the data conforms to a known function (such as the cosine function), can we test for that function more effectively (that is to say by the stability of the weights upon cross-validation) than using a polynomial function for developing prediction equations?

The purpose of this paper presentation is to describe a methodological approach to facilitate the ease of the description and prediction of cyclical events. The proposed method will be compared with traditional methods for fitting curvilinear relationships, and the potential benefits of the proposed method will be outlined. Finally, an example will be provided to demonstrate the method and its advantages over traditional curve fitting methods for the description and prediction of cyclical events.

Introduction

There are many phenomena with cyclical patterns that are of interest to statisticians, psychologists, and epidemiologists. Seasonal Affective Disorder is a clear example of a cyclical phenomena that is of growing interest in psychological research for both children and adults (Glod, & Baisden, 1999; Rohan & Sigmon, 2000). In turn, Partonen, Piironen, Loennqvist, and Jouko (2000) have demonstrated that additional psychological symptoms may have a seasonal pattern. Accidental death has also been shown to be tied to an annual pattern (Coren, 1996). On the other hand, there are ready examples of phenomena that cycle on a monthly basis. For example, suicide rates have been shown to cycle on a monthly basis with suicides being more frequent during the first and second weeks of the month (Phillips & Ryan, 2000). Clearly, there are numerous additional phenomena that cycle on a monthly or daily basis.

In the regression literature, authors have traditionally suggested that these phenomena be predicted by curve fitting or log-linear techniques (Cohen & Cohen, 1983; McNeil, Newman, & Kelly, 1996; Pedhazur, Pedhazur-Schmelkin, 1991). A model of a cyclical phenomena (Cosine) using a polynomial equation will take the general form:

$$\text{Model } Y = a_0U + a_1X + a_2X^2 + \Sigma(a_{1+2i}X^{1+2i} + a_{2+2i}X^{2+2i}) + \text{Error}$$

for i cycles

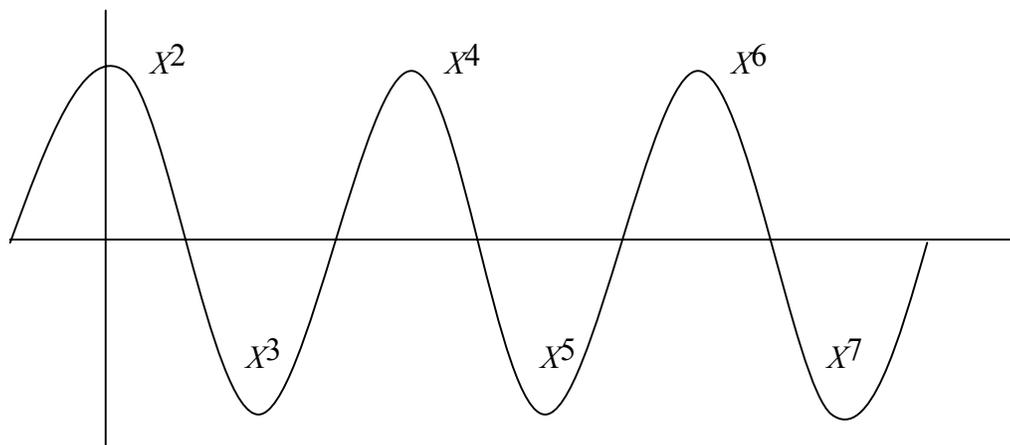


Figure 1. General polynomial representation of a cyclical function.

As can be seen in Figure 1, the first cycle requires a model with a 4th degree polynomial, and each additional cycle requires an addition of two variables to model. For example, a series of 4 cycles would require a model extending to a 10th degree polynomial. If one uses the trigonometric function (cosine), this relationship can be modeled with two predictor variables:

$$\text{Model } Y = a_0U + a_1\text{COS}(X) + a_2X + \text{Error}$$

In this presentation, we will outline a method for using trigonometric functions (cosine) to predict cyclical phenomena; although this technique is not new, few examples have been provided in the literature. A method will be discussed for utilizing the sine or cosine function to predict a cyclical phenomena. In order to facilitate this discussion, the method will be anchored to an example: modeling temperature based on the time of year. Using this example, we will outline the general steps necessary to model a cyclical phenomena. In turn, the example will be used to demonstrate some of the advantages of the use of trigonometric functions in modeling cyclical phenomena.

Method

The use of a trigonometric function to model a cyclical phenomena requires the estimation of two variables: (1) the period, and (2) the amplitude of the function. The period is the time taken to make one complete cycle (a complete oscillation). For the Sine and Cosine functions, the period is measured in radians and is equal to 2π . The amplitude, on the other hand, is the displacement (the distance from the crest to the trough) of the cycle.

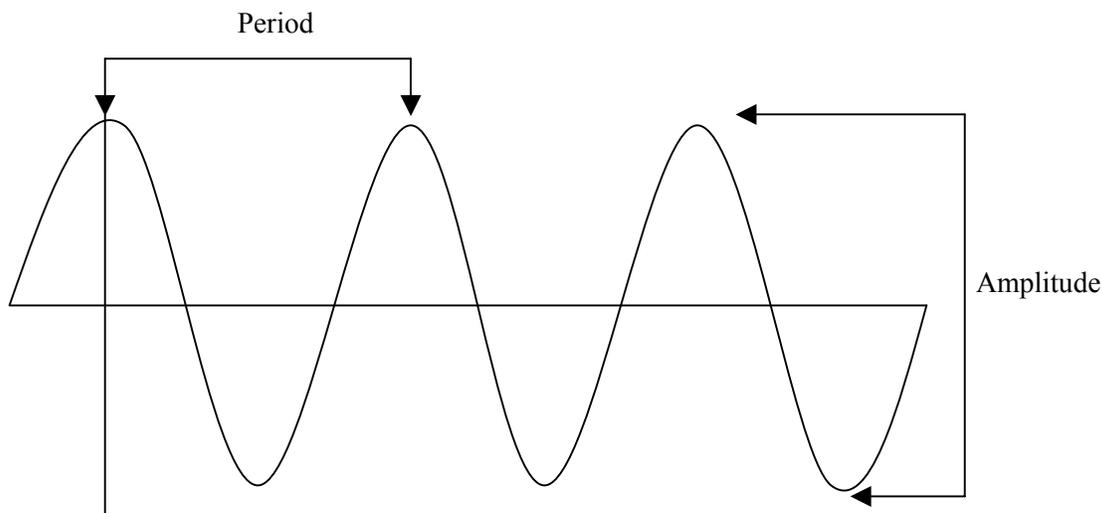


Figure 2. Parameters of a wave function.

One could estimate the period of a cyclical function through either theoretical or empirical means. The choice would be driven by the researcher's question and by whether theory is suggestive of a fixed cycle. Some psychological phenomena (e.g. seasonal affective disorder) are clearly linked to an annual cycle (period = 1 year). In other cases, one could estimate the period empirically. This can be achieved by plotting the dependent measure across time. For this presentation, we plotted temperature readings across time. As would be expected, there was a clear cyclical pattern of warmer temperatures in the summer and cooler temperatures in the winter.

At this point, we had a clear cycle with peaks occurring in July and troughs occurring in January. The period for this cycle is 12 months; therefore, all time measurements were divided by 12 and multiplied by 2π in order to convert the time measurement to a radian scale.

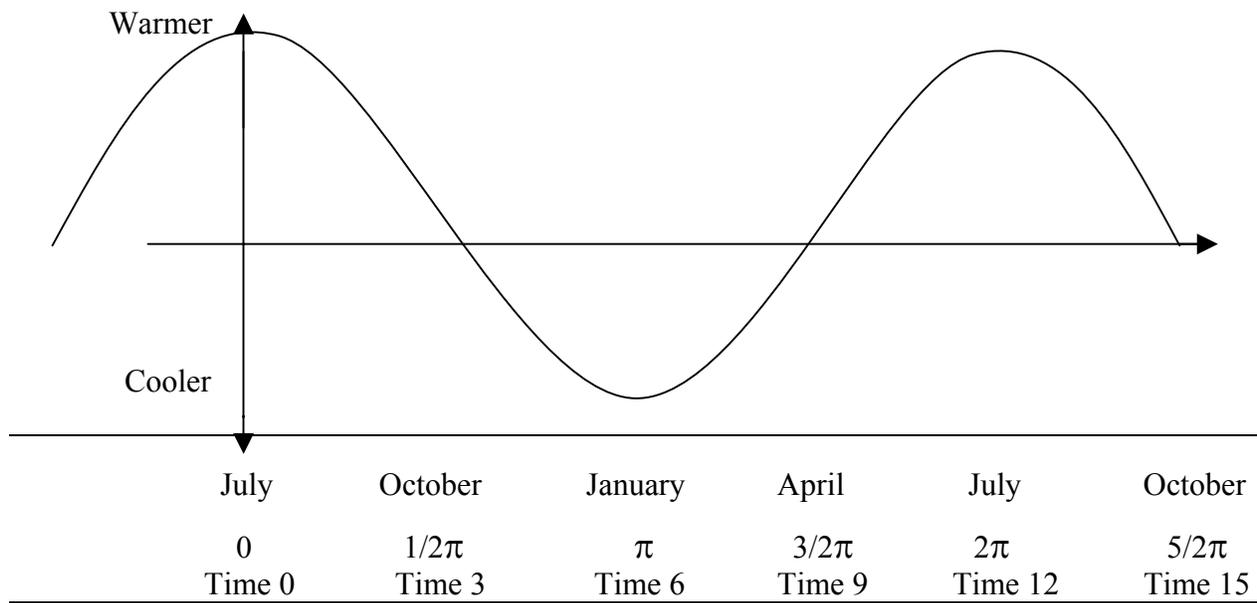


Figure 3. Plotting temperature across time using a radian representation of time.

Once the period was determined and the time measurements were converted to a radian scale, we were able to create a regression model to fit the temperature data:

$$\text{Model } Y = a_0U + a_1\text{COS}(X) + a_2X + \text{Error}$$

where Y = Temperature and X = Measurement Time in Radians

The amplitude of the wave is, then, represented by the weight (a_1) of the $\text{COS}(X)$ variable. The least squares regression solution for this model calculates the amplitude (a_1) in such a manner that the error sums of squares is minimized. This regression model can then be used with either the whole data set (i.e. ten year's data) or with smaller subsets (e.g. one year's data) within the whole data set.

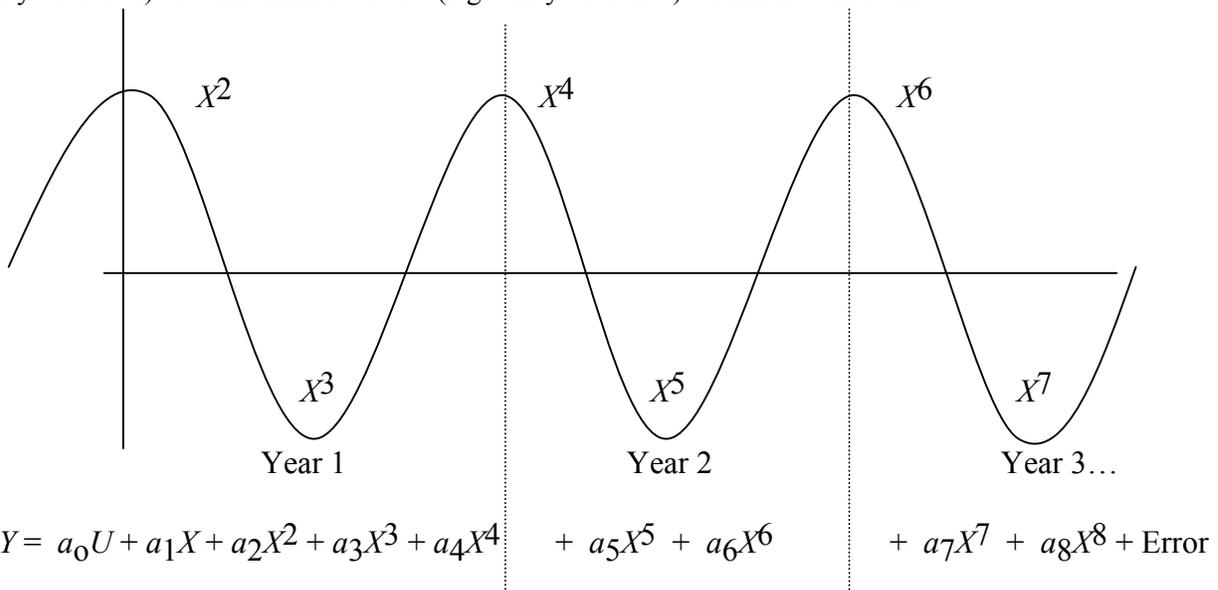


Figure 4. Polynomial representation of a cyclical function.

In order to model the same phenomena using a polynomial, it is necessary to take into account the number of cycles to be modeled when developing a regression equation to fit the data. Each inflection point requires the additional of a polynomial in order to be modeled. Therefore, the model needed to fit one year's data would be as follows:

$$\text{Model } Y = a_0U + a_1X + a_2X^2 + a_3X^3 + a_4X^4 + \text{Error}$$

where Y = Temperature and X = Measurement Time in Radians

Each additional year would require the addition of two polynomials in order to model the inflection points in the curve. Ten year's data would require a model extending to the 22nd power of the time measurement!

Results

In order to demonstrate the differences between these modeling methods, both were applied to a longitudinal temperature data set. The data for this analysis were collected from a web site maintained by Utah Climate Center, Utah State University (<http://climate.usu.edu/free/default2.htm>). This web site contains weather data for hundreds of sites nationally. The temperature data for this study were taken, specifically, from a site in Colorado (Akron – site # 5011403) and consist of daily maximum temperature readings for a ten year span beginning January 1, 1980.

The data were first modeled using a polynomial function and subsequently modeled using the cosine function. The analysis began with a sample of data from a single year; and, then, proceeded to a sample from five year's data. Each sample, ten percent of the cases, was drawn randomly from the time range sampled. Cross-validations were based on random, matched (ten percent), samples. The results of these analyses are in Table 1.

In the first year, the polynomial function accounted for slightly more variance than did the cosine function. As can be seen in the column of shrunken R-squares, both functions were stable when cross-validated. At five years, the cosine function accounted for more variance in the temperature data and was considerably more stable when cross-validated (12.2% versus 95.9% shrinkage).

The problem of shrinkage with the polynomial functions is further exacerbated when one attempts to model additional cycles (years) of data. To demonstrate this point, we attempted to model ten years of daily maximum temperatures using both polynomial and cosine functions. Again, a ten percent sample of the cases, was drawn randomly from the time range (ten years) sampled. Cross-validations were based on random, matched (ten percent), samples.

To demonstrate the explanatory power of the polynomials, the data was first modeled with only a second degree polynomial. Each subsequent model added an additional polynomial to the model until the final model included the second to eighth degree polynomials. As can be seen from Table 2, the addition of higher order polynomials incrementally adds to the variance accounted for by the model, but these models are not stable when cross-validated. In turn, the polynomial models did not account for nearly as much variance as did the cosine function, which was also stable when cross-validated.

Table 1. Modeling One and Five Years Data: Polynomials versus Cosine Function

Model	R-square	Adjusted R-square	Shrunken R-square	Percent of Shrinkage
Polynomial 1 year	.7622	.7334	.7367	3.8%
Polynomial 5 years	.3044	.2651	.0123	95.9%
Cosine 1 year	.7311	.7157	.6951	4.9%
Cosine 5 years	.6551	.6514	.5751	12.2%

Table 2. Modeling Ten Years Data: Polynomials versus Cosine Function

Model	Largest Polynomial in the Model	R-square	Shrunken R-square
Polynomial	X^2	.0018	.0016
Polynomial	X^3	.0042	.0007
Polynomial	X^4	.0072	.0003
Polynomial	X^5	.0182	.0005
Polynomial	X^6	.0498	.0007
Polynomial	X^7	.0508	.0010
Polynomial	X^8	.0883	.0007
Cosine	NA	.6718	.6537

Note: A weight of zero is applied to higher order polynomials (>8th power) when they are modeled.

Discussion

The purpose of this paper was to provide an answer to the question of the relative effectiveness of the cosine function versus a polynomial function in the description and stability of prediction of a specific set of longitudinal data. With the present data set, the cosine function clearly provided a more stable prediction of the maximum daily temperatures in Akron, Colorado between the years of 1980 and 1990. The benefit of using a cosine function to predict the temperature scores was particularly evident when more than one year of data was modeled. When only one year's data was modeled, there was a slight advantage to using a polynomial model to predict temperature. The polynomial models appear to capitalize on unique variance in the sample. On the other hand, the cosine function appears to be relatively stable across samples and time. Therefore, if the data conforms to a known function, using the function to model the data, especially when numerous cycles are to be modeled, would give one a more stable prediction and hence greater confidence when generalizing beyond the specific sample used in a given study. Currently, we are beginning to work to determine the point at which there is a clear advantage for using one method rather than the other.

There are times, however, when the two modeling methods could be used in a complementary fashion. If one is unsure as to whether a known function exists within a data set, one could model a small interval (2 or 3 cycles) of the data using a polynomial. If the polynomial model is suggestive of a known function (e.g. cosine), one could then model the following intervals using the known function. In this manner, the two modeling methods could be used in conjunction with one and other.

References

- Cohen, J., & Cohen, P. (1983). *Applied multiple regression/correlational analysis for the behavioral sciences*. (2nd ed.). Hillsdale, NJ; Erlbaum.
- Coren, S. (1996). Accidental death and the shift to daylight savings time. *Perceptual & Motor Skills*, 83(3, Pt 1), 921-922.
- Glod, C., & Baisden, N. (1999). Seasonal affective disorder in children and adolescents. *Journal of the American Psychiatric Nurses Association*, 5(1), 29-33.
- McNeil, K., Newman, I., & Kelly, F. (1996). *Testing research hypotheses with the general linear model*. Carbondale, IL: Southern Illinois Press.
- Partonen, T., Piironen, P., & Loennqvist, J. (2000). Season-dependent symptoms in consultation-liaison patients. *International Journal of Psychiatry in Clinical Practice*, 4(2), 151-154.
- Pedhazur, E., & Pedhazur-Schmelkin. (1991). *Measurement, design, and analysis: An integrated approach*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Phillips, D., & Ryan, N. (2000). An abrupt shift in U.S. suicide levels around the month boundary. In R.W. Maris (Ed.), *Review of Suicidology, 2000*. (pp. 34-44). New York; Guilford Press.
- Rohan, K., & Sigmon, S. (2000). Seasonal mood patterns in a northeastern college sample. *Journal of Affective Disorders*, 59(2), 85-96.

Send correspondence to: Russell Brown, 200 Strecker Drive, Tallmadge, OH 44278
Email: rcbrown@uakron.edu.

POSTMASTER: Send address changes to:

Jeffrey B. Hecht, MLR/GLM SIG Executive Secretary

Department of Educational Technology, Research & Assessment

Northern Illinois University

DeKalb, IL 60115-2854

Multiple Linear Regression Viewpoints (ISSN 0195-7171) is published by the
AERA Special Interest Group on Multiple Linear Regression: General Linear Model
through the University of Alabama at Birmingham.