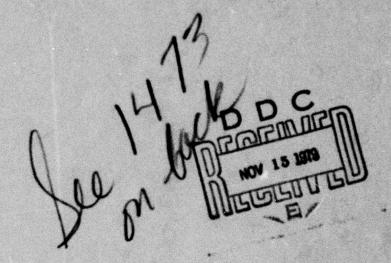




SPECIFIC INTERACTION CONTRASTS:

A STATISTICAL TOOL FOR REPEATED-MEASURES DESIGNS

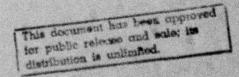
V. K. THARP, JR. ARDIE LUBIN



**REPORT NO. 76-78** 

DOC FILE COPY





# **NAVAL HEALTH RESEARCH CENTER**

P. O. BOX 85122 SAN DIEGO, CALIFORNIA 92138

NAVAL MEDICAL RESEARCH AND DEVELOPMENT COMMAND BETHESDA, MARYLAND

79 11 15 168



Specific Interaction Contrasts:

A Statistical Tool for Repeated-Measures Designs

Van K. Tharp, Jr. and Ardie Lubin

Naval Health Research Center

San Diego, California 92152

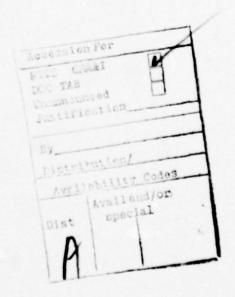


This document has been approved for public release and sale; its distribution is unlimited.

Running head: Specific Interaction Contrasts

#### Abstract

Experimental psychologists often use multifactor repeated-measure designs in which interactions are the most important effects to be assessed. An experimenter has at least five ways to evaluate such interactions: (1) a univariate repeated-measures analysis of variance (ANOVA), with (probably) inflated estimates of the degrees of freedom; (2) a univariate repeated-measures ANOVA with the Greenhouse-Geisser conservative estimate of the degrees of freedom; (3) the Greenhouse-Geisser stepwise analysis; (4) a multivariate ANOVA; and (5) specific interaction contrasts. We show that no matter which of the above paths is chosen, the careful experimenter must compute specific interaction contrasts (i.e., <u>t</u>-tests). A worked example is given.



# Specific Interaction Contrasts: A Statistical Tool for Repeated-Measures Designs

R. A. Fisher (Fisher & MacKenzie, 1923) used the concepts of interaction and additivity when he invented the Analysis of Variance (ANOVA) method for dealing with multifactor experiments. Until recently, these concepts did not play a primary role in any theoretical models commonly used in experimental psychology. However, in 1969, Sternberg incorporated both concepts directly into his additive-factor method for decomposing reaction time (RT) into processing stages.

Assuming that RT is composed of a number of additive stages in a known order, Sternberg proposed that each stage be studied by influencing its duration with various treatments. Two treatments which influence one or more stages in common should have an interaction effect on RT. But if the two treatments influence different stages, then each should have an additive effect on RT.

Taylor (1976) recently extended this methodology to conditions in which some dependence may occur between processing stages. His primary restrictive assumption is that stage dependence must be expressable as a linear function of the stage times involved. Under these conditions, Sternberg's hypothesis concerning additivity does not hold. That is, the absence of a significant interaction between the effects of two variables does not imply that these variables influence different processing stages. However, two treatments influencing one or more processing stages in common still should have an interaction effect on RT. Thus, the important question becomes: How do we test for interactions?

The traditional significance test of interactions is the F-ratio in the analysis of variance (ANOVA). We find two major difficulties in using this test with the form of experimentation advocated by Sternberg (1969) and Taylor (1976). First, multifactor repeated-measure designs, such as those required by Sternberg (1969) and Taylor (1976), do not meet a critical assumption of the ANOVA--independent scores. Second, the F-ratio is a vague test, telling the experimenter almost nothing.

R. A. Fisher's first use (1923) of the ANOVA was to study the effect of treatments on plots of ground. His most important assumption was that the criterion score (yield) for any particular plot was independent of the score for any other plot. The F-ratio, consequently, has some problems in its application to a repeated-measures experiment. For example, Sternberg (1969) and Taylor (1976) advocate exposing a subject to all possible conditions in a multifactor experiment. Even if we assume that the carryover effects of every treatment on all subsequent treatments are negligible, we must deal with the fact that any two scores measured on the same organism, can have (and usually do have) non-zero correlations. This makes the usual univariate repeated-measure F-ratio misleading or uninterpretable (e.g., see Lana & Lubin, 1963 and the justification section of the present paper).

An even greater difficulty is that the calculation of  $\underline{F}$ -ratios in conjunction with a completely general (i.e., unspecified) ANOVA model tells the experimenter almost nothing about a hypothesized interaction effect. A significant  $\underline{F}$ -ratio tells us nothing about the direction, amount, or location of the underlying interaction; while a nonsignificant  $\underline{F}$ -test may simply lack the power to detect a single large inter-

action effect. J. A. Nelder (see Plackett, 1960, p. 213) has criticized such general ANOVA models on the grounds of vagueness. He suggested, instead, that the model be specialized to fit the particular application, thus gaining power and supplying more information to the experimenter.

Professor G. A. Barnard, also commenting on the paper by Plackett (1960), noted that the ANOVA essentially reduces to a set of independent contrasts and that we are free to select groups of contrasts in any manner we choose (see Plackett, 1960, p. 215). Since efficient use of the ANOVA involves selecting a specific model to fit the chosen application, why not devise specific contrasts to test one's hypotheses? As Geisser has said: "When there are contrasts of scientific importance the omnibus F is irrelevant" (personal communication). Contrast weights specify the size, sign, and location of the putative effects. A contrast can be tested by a t-ratio. In this way we avoid the vagueness and lack of power of the F-ratio. In particular, contrasts do not require the many assumptions necessary for a repeated-measure ANOVA.

Sternberg (1969), in fact, recommended that interactions be evaluated with specific interaction contrasts. Unfortunately, he provided no detailed worked examples of his procedure. Furthermore, the interaction in his one example, in which computation is briefly discussed, has only one possible degree of freedom, so it does not reveal the true potential of the procedure. We have seen no other use of this technique, so we can only assume that Sternberg's description was insufficient for most of his readers. The purpose of this paper, therefore, is to detail the computational procedures involved and give some theoretical justification for the use of contrasts.

## Specific Interaction Contrasts: Computation

A specific interaction contrast is an estimate of the variation due to interaction based on an explicit set of contrast weights. In a multifactor, repeated-measures design, the contrast procedure requires that one first obtain a set of contrast weights. These may be obtained by prior reasoning or from prior empirical knowledge. Each of the p conditions in the interaction must have a contrast weight, and the sum of these weights will be zero, by definition. The next step involves calculating the residuals from additivity (i.e., the interaction effect) for each subject, obtained by subtracting the overall mean and the main effects from a subject's score for each of the p conditions. Again, by definition, the residual from additivity should be a set of numbers which sum to zero. The contrast weights are then applied to the interaction effect of each subject to obtain a contrast score. Given n subjects, the set of n contrast scores can be used in a routine t-test of the null hypothesis, Ho, that the interaction effect is zero. If Ho is accepted, then additivity of the main effects is implied only for this particular weighted combination of treatment levels. Generally, the maximum number of independent, specific interaction effects is limited by the degrees of freedom for the interaction term in the ANOVA model.

Let Y be the contrast score for the kth subject. Then

$$\underline{Y}_{k} = \sum_{i j} \sum_{i j} u_{ij} \underline{q}_{ijk} \qquad \qquad i = 1 \text{ to } r$$

$$j = 1 \text{ to } c \qquad (1)$$

where  $\omega_{ij}$  is the contrast weight for the <u>ij</u>th cell, and  $\underline{q}_{ijk}$  is the residual from additivity for the <u>ij</u>th cell within the <u>k</u>th subject. First, we will give the standard method of computing  $\underline{q}_{ij}$ , the residual

from additivity. Then we will discuss the more difficult problem of estimating the contrast weights,  $\omega_{ii}$ .

Let us look at an example (Tharp, 1975) in which the time required to name pictures is measured under alcohol and baseline conditions.

The designated correct names to the stimuli came from five different word-frequency (WF) categories. The deleterious effect of alcohol on verbal reaction time is expected to increase as word-frequency decreases—an interaction effect.

In this example<sup>2</sup> there are five levels of WF and two levels of drug, giving us (5-1) (2-1) = 4 degrees of freedom for the interaction. Thus, there are four possible independent specific interaction effects. Table 1 illustrates summary data for one subject. There are ten scores per subject.

Table 1 about here

## Residuals from Additivity

Let  $\overline{X}_{ij}$  be the average score in row <u>i</u> and column <u>j</u>. First, the effect of the <u>i</u>th row,  $R_i$ , is equal to:

$$\underline{R}_{i} = \frac{\sum_{j} \overline{X}_{ij}}{2} - \overline{X}.. \qquad j = 1,2$$
 (2)

where  $\overline{X}$ . is the grand mean. Second, the effect of the <u>j</u>th column,  $\underline{C}_{j}$ , is equal to:

$$\underline{c_j} = \frac{\sum_{i=1}^{N} \overline{x_{ij}}}{5} - \overline{x}.. \qquad i = 1 \text{ to 5} \qquad (3)$$

Finally, the residual from additivity  $\underline{q}_{ij}$ , in row i and column j is equal to:

$$\underline{q}_{ij} = \overline{\underline{x}}_{ij} - \underline{R}_i - \underline{C}_j - \overline{\underline{x}}.$$
 (4)

Note that  $\underline{R_i}$  and  $\underline{C_j}$  are row and column <u>effects</u>, respectively, <u>not</u> means. For example, for the score in row 1 and column 1 of Table 1, the row effect ( $\underline{R_l}$ ) is equal to the row mean minus the grand mean (i.e., 859.5 - 1035.3 = -175.8). The column effect ( $\underline{C_l}$ ) is equal to the column mean minus the grand mean (i.e., 1124.8 - 1035.3 = 89.5). Finally, to obtain the residual from additivity ( $\underline{q_{11}}$ ), the row effect, the column effect, and the grand mean are all subtracted from the cell score (i.e., 901 + 175.8 - 89.5 - 1035.3 = 48.0). All the residuals from additivity, computed from Table 1, are shown in Table 2.

#### Table 2 about here

Each residual is equal to a random error term  $(\underline{e}_{ijk})$ , plus a putative interaction effect,  $\tau_{ij}$ .

#### Determination of Contrast Weights

Assume that the true interaction effect in the <u>ij</u>th cell is  $\tau_{ij}$ . Then the best estimate of the contrast weight,  $\omega_{ij}$ , is  $\tau_{ij}$ . If the set of ten contrast weights is equal to (or proportional to) the set of ten interaction effects, then the contrast score of equation 1 is maximized.

Theoretical Weights. To the extent that an experiment is based on experience, sound judgment and prior scientific knowledge, one should have little difficulty predicting the amount and direction of any

hypothesized interactions. In our example, the effect of alcohol was expected to remain relatively constant from WF<sub>1</sub> to WF<sub>3</sub> and to increase dramatically at WF<sub>4</sub> and WF<sub>5</sub>. These results were expected because (1) an exponential increase in a priori stimulus-response uncertainty occurs from WF<sub>1</sub> to WF<sub>5</sub> and (2) the effect of alcohol increases as a function of experimentally defined stimulus-response uncertainty (Tharp, Rundell, Lester, & Williams, 1974). Thus, based on prior knowledge, we might postulate the set of weights,  $\omega_{ij}$ , shown in Table 3.

Table 3 about here

There are two restrictions on these interaction contrast weights. First,

$$\int_{j}^{\omega} ij = 0 \text{ for every row} \qquad j = 1,2$$
 (5)

and second

$$\int_{3}^{2} \omega_{ij} = 0 \text{ for every column} \qquad i = 1 \text{ to 5} \qquad (6)$$

One can adjust any set of weights to fit these rules by subtracting from each weight the appropriate row mean, finding column means of the row-adjusted scores, and then subtracting the appropriate column mean from each cell.<sup>3</sup>

Contrast weights are usually given as single digits lying between -9 and +9. Single digit weights might be obtained by smoothing and dividing all scores by their lowest common denominator. Two-digit accuracy might be justified with 50 or more subjects.

Cross-Validation Procedure. If one cannot ask questions about interaction effects in terms of prior contrast weights, then empirical

post-hoc contrast weights can be estimated from the cell means for half the experimental subjects (analysis group) by setting  $\omega_{ij} = \underline{q}_{ij}$ . The remaining subjects (cross-validation group) then can be used to get an unbiased estimate of each putative interaction effect and to test the obtained effect for significance. This cross-validation method is standard operating procedure for psychometricians in multiple regression, test construction, etc., (Mosier, 1951). Cross-validation is rarely used as such by statisticians. However, the "jack-knife method", popularized by Tukey, as well as the Geisser "predictive sample reuse method", can be viewed as a generalization of the usual two-group cross-validation procedure (see Mosteller & Tukey, 1968; Geisser, 1975).

Our example consisted of 24 subjects, so 12 of them (i.e., the analysis or training group) are used to estimate the weights while the remaining 12 are the cross-validation (testing) group. Table 4 gives the set of weights derived from the analysis group.

Table 4 about here

Notice that these values are quite similar to our theoretical values in that most of the alcohol effect occurs at  ${
m WF}_5$ .

# Significance Tests

The significance test will have maximum power when the interaction effect,  $\tau_{ij}$ , for each cell is equal to, or proportional to, the contrast weight for that cell. For simplicity in what follows, we omit the case of proportionality. Thus, the experimenter's hypothesis is:

$$\underline{H}_{1}: \quad \tau_{ij} = \omega_{ij}, \tag{7}$$

and the null hypothesis is:

$$\underline{H}_0: \quad \tau_{ij} = 0 \tag{8}$$

Both hypotheses assume that  $\tau_{ij}$  equals the expected value of  $\underline{q}_{ij}$ , so  $\underline{H}_0$  is equivalent to  $\underline{q}_{ij}$  = 0, where

$$\overline{q}_{ij} = \frac{\sum_{k=1}^{n} q_{ijk}}{n}$$
 (9)

and  $\underline{n}$  refers to the number of subjects. The contrast score for the  $\underline{k}$ th subject is defined by equation 1. The expected value of  $\underline{Y}$ , computed over the  $\underline{n}$  subjects, will equal zero if  $\underline{H}_0$  is true. If  $\underline{H}_1$  is true, then the expected value is:

$$E(\overline{Y}) = \sum_{i,j} \sum_{i,j}^{2} \omega_{i,j}^{2} = \sum_{i,j} \sum_{i,j}^{2} \tau_{i,j}^{2}$$
 (10)

We now have a between-subjects <u>t</u>-test, eliminating the repeated-measures problems.

Theoretically, under  $\underline{H}_1$  the  $\underline{Y}_k$  scores for the cross-validation group will range from zero to  $\sum\limits_{i=j}^{n}\omega_{ij}^2$ . Under  $\underline{H}_0$ , the  $\underline{Y}_k$  scores can be positive or negative, centering on zero. The conventional  $\underline{t}$ -test is simple to obtain.

$$\underline{t} = \underline{\underline{Y}} \sqrt{\underline{n}}$$
where  $\underline{H}_0$  is  $\underline{E}(\underline{\underline{Y}}) = 0$ 

$$\underline{H}_1 \text{ is } \underline{E}(\underline{Y}) > 0$$

and the t-ratio has (n-1) degrees of freedom.

To complete our example, Table 5 illustrates the  $\frac{Y}{k}$  scores for the cross-validation group using the weights obtained from the analysis group (right) or the theoretical weights (left).

Table 5 about here

For Sternberg (1969) additive effects between two treatments are as meaningful as interactive effects. Thus, he suggests that

... one might present findings in terms of mean interaction contrasts of theoretically interesting magnitudes, and adjust tests so that errors of Types 1 and 2 have equal probabilities with respect to such alternatives. (p. 310)

One could, for example, use the appropriate <u>t</u>-ratio at the 50% level as a rejection point. This adjustment cannot be used to infer additivity, however, within the context of the test we recommend for interactions. That is, rejection of a specific interaction hypothesis does not imply additivity. In view of Taylor's (1976) cautions against interpreting additivity, this inability to infer an additive relationship from a specific test does not appear to be a serious drawback.

# Several Interaction Hypotheses

The major interest of the study may not be the comparison of  $\underline{H}_0$  with  $\underline{H}_1$ . In some cases, divergent theories might lead to divergent hypotheses as to the nature and direction of the interaction. For example, suppose that some evidence predicted a "golden-mean" theory

in which the optimal effect of alcohol occurs at the median word frequency level,  $WF_3$ , while decrement occurs at the extreme word frequency levels,  $WF_1$  and  $WF_5$ . Thus, the theoretical weights might be those shown in Table 6. We now have two experimental hypotheses,

#### Table 6 about here

 $\underline{H}_1$  and  $\underline{H}_2$ . Let the vector of weights given in Table 4,  $\underline{V}_1$ , represent  $\underline{H}_1$ ; and the vector of weights given in Table 6,  $\underline{V}_2$ , represent  $\underline{H}_2$ . Each vector will yield a <u>t</u>-ratio--the higher the <u>t</u>-ratio, the better the hypothesis fits the data. To make an accurate estimate of the significance of the <u>t</u>-ratios, the rejection levels must be adjusted to take account of the fact that two similar significance tests have been obtained from the same data.

Although there are many solutions to this multiple-comparison problem (e.g., Miller, 1966), we prefer the Dunn-Bonferroni method (Dunn, 1959). Assume that you want to hold your experimentwise Type 1 error at .05 (i.e., when Ho is true then either the trational or both the trational will be significant in five percent of all comparisons). In the simplest version of the Dunn-Bonferroni, the chosen alpha level is divided by the number of comparisons to obtain a critical level of significance for each comparison. Thus, in this example we would divide .05 by two to obtain .025—the level at which each the test would be evaluated. This is a very conservative test which guarantees that the Type 1 error will be .05 or greater, so it lacks some of the power of other tests. In general, it will be more conservative (i.e., the Type 1 error will

be larger) as the correlation between  $\underline{V}_1$  and  $\underline{V}_2$  increases. Thus, the experimenter should be careful to avoid testing redundant hypotheses. When the weights of  $\underline{V}_1$  correspond exactly to the true interaction effects,  $\omega_{ij}$ , then the contrast scores for  $\underline{V}_1$  will account for all interaction variance (i.e., any vector which is orthogonal to  $\underline{V}_1$  will have a  $\underline{Y}_k$  of zero).

Most experiments contain several families of statistical hypotheses. The Dunn-Bonferroni adjustment may apply separately to each family (Miller, 1966). For example, one might be interested in the effect of a new treatment on various information processing stages, and thus introduce that treatment into a multifactor design with several treatments whose locus of effect has already been "established" by means of the Sternberg-Taylor procedure. Such an experiment would involve a two-step analysis. Step one, constituting one family of statistical tests, would involve confirming the interactions between the established treatments in data which did not include the new treatment. The second step, constituting the second family, would be a search for interactions between the effects of the new treatment and the established ones.

#### Justification

Let us review the remaining four ways to evaluate interactions in a repeated measures design: (1) the univariate repeated-measures

ANOVA using the (probably) inflated estimate of degrees of freedom;

(2) multivariate ANOVA: (3) conservative degrees of freedom with a univariate ANOVA: and (4) the Greenhouse-Geisser stepwise analysis. We will demonstrate that no matter what the outcome of the given procedures, the careful experimenter must use specific interaction contrasts.

#### Univariate Repeated-Measures ANOVA

For the conscientious experimenter few problems are so exasperating as the analysis of a repeated-measures design. In 1954, Box showed that the exact analysis of a repeated-measures design is a multivariate analysis and that the basic answers must be given in terms of multivariate <u>F</u>-ratios (see, also, Winer, 1971, sections 4.4 and 4.9). The usual univariate <u>F</u>-ratio approach, set forth in almost every psychological statistics text, is strictly valid only under a set of necessary and sufficient conditions known variously as (a) the "circularity property" (Rouanet and Lepine, 1970); (b) "equality of variances of differences" (Huynh and Feldt, 1970); or (c) "homogeneity of interaction variances" (McNemar, 1962, pp. 315-316).

All E-ratios with only one degree of freedom for the numerator are valid under the standard multivariate normal assumptions (see Appendix). With three or more measures per subject, one could test for the "circularity property" and drop the univariate approach when it does not hold. When the circularity property does hold, one need only worry about interpreting the meaning of a significant E-ratio.

# Multivariate Analysis of Variance

One possible solution for analyzing repeated-measures data is to avoid the univariate approach--do a multivariate ANOVA as soon as we have three or more measures per subject. Unfortunately, the multivariate ANOVA presents many problems unless one has a large number of subjects.

Given  $\underline{rc}$  scores per subject, the number of subjects must be equal to or greater than  $(\underline{rc}-1)$  in order to compute a multivariate ANOVA. This absolute bar is present because every multivariate ANOVA demands the computation of the inverse of a within-group covariance matrix. When

 $\underline{\mathbf{n+1}}$  is less than the number of measures per subject, the determinant of the covariance matrix must be zero, and the inverse does not exist.

Most experimenters probably can slip under the (<u>rc-1</u>) barrier with one or two degrees of freedom to spare. If so, then we usually run into the problems of the ill-conditioned covariance matrix and the inherently large sampling variance of correlation coefficients and variances. The idea of estimating the inverse of a 4x4 or 5x5 covariance matrix, which is based on less than 20 degrees of freedom, only appeals to those who have an overwhelming faith in small samples.

The inherent sampling instability of second-order statistics (e.g., correlations, variances, etc.) will interact with small sample size to emphasize the problems of the multivariate-normal model. The unbiased estimate of a covariance matrix demands many more assumptions than the unbiased estimate of a difference between two independent means. Some of these assumptions are given in the Appendix. Ordinarily, even when we do not have exact normality or homogeneous variances, the Central Limit Theorem guarantees that the difference between a pair of independent means goes very quickly towards a normal distribution with homogeneous variance, provided that scores are independent and that there are enough of them. But second-order statistics are very sensitive to even slight deviations from normality (e.g., the fourth moment), and the Central Limit Theorem has very little effect with small samples. Furthermore, if some fundamental assumption such as linearity or independence has been violated, then the Central Limit Theorem simply does not apply. For example, if one measure has a non-monotonic relation to another, then no increase in the sample size will linearize that relation. In summary, then, multivariate analysis of variance is an exact way to evaluate repeated-measure interactions. Nevertheless, a multivariate ANOVA is practical only (a) if one has fifty cases or more (given that rc is more than 5) and (b) if the important assumptions of linearity, a well-conditioned covariance matrix, etc., hold (see Appendix).

#### Conservative Estimate of the Degrees of Freedom

Some investigators have attempted to avoid the repeated-measures problems of the univariate ANOVA by using the conservative test advocated by Greenhouse and Geisser (1959). This approach is based on a statistic, epsilon, developed by Box (1954). When epsilon equals unity, then the usual univariate ANOVA of repeated-measures is valid. When epsilon is less than unity, multiplying the degrees of freedom for the numerator and the denominator of the univariate E-ratio by epsilon, gives the approximate degrees of freedom for a valid evaluation of the univariate E-ratio.

The Greenhouse and Geisser (1959) conservative test avoids the estimation of epsilon. It uses the fact that epsilon cannot go any lower than 1/(p-1) where p is the number of measures per subject. Consequently, by using this minimum value of epsilon, one usually underestimates the degrees of freedom and the significance level of the p-ratio. If the conservative test is significant, then the multivariate analysis would also be significant. But, the Greenhouse-Geisser conservative test will only reject the null hypothesis for rather large p-ratios. When the conservative test is not significant, the experimenter is given very little information.

## Greenhouse-Geisser Stepwise Analysis

Greenhouse and Geisser in describing their stepwise analysis (1959, p. 110) advised that one should test the univariate  $\underline{F}$ -ratio first, using the nominal degrees of freedom without the epsilon correction. When this  $\underline{F}$ -ratio is significant, one must then perform the conservative test. The epsilon correction is needed only when the  $\underline{F}$ -ratio with inflated degrees of freedom is significant and the conservative test is not. This stepwise approach, also recommended by Lana and Lubin (1963), assumes that no further analysis is needed when the univariate  $\underline{F}$ -ratio with possibly inflated degrees of freedom is not significant. This assumption is wrong. Davidson (1972) showed that one could easily obtain a significant Hotelling  $\underline{T}^2$  (i.e., a significant multivariate  $\underline{F}$ -ratio) on the same data. Thus, a non-significant  $\underline{F}$ , with possibly inflated degrees of freedom, also gives the experimenter very little information.

In summary, if a significant  $\underline{F}$  is obtained with the conservative test or by using epsilon to approximate the valid degrees of freedom, then the only problem is to interpret the meaning of the F-ratio. Otherwise, one can turn to a multivariate ANOVA if enough subjects are available (e.g., 50 or so, given that  $\underline{p}$  is about 5) and the important multivariate assumptions are met.

# Specific Interaction Contrasts Versus the Omnibus F-Ratio

Justification. When a non-significant F is obtained either with the Greenhouse-Geisser conservative test or by using epsilon to approximate the valid degrees of freedom and the conditions are not appropriate for a multivariate ANOVA, only one of our suggested solutions remains—specific interaction contrasts. Such a situation always justifies the use of specific interaction contrasts.

However, the most powerful argument for the use of interaction contrasts is a pragmatic one. Even when <u>F</u>-ratio tests (univariate or multivariate) are appropriate, only two outcomes are possible—a significant <u>F</u>-ratio or a nonsignificant <u>F</u>-ratio. Both outcomes should lead the careful experimenter to test for specific interactions.

A nonsignificant <u>F</u>-ratio may be due to a lack of power, the result of considering all possible interaction contrasts simultaneously. Therefore, a careful experimenter should always apply any prior interaction contrast (derived from hypothesis or an analysis group) to the data to see if the <u>t</u>-test is significant even though the overall <u>F</u>-ratio was not.

If the <u>F</u>-ratio is significant, then the experimenter still has to determine which cells yielded the significant interaction effect, as well as the direction and amount of each cell effect. A significant F-ratio does <u>not</u> guarantee that the experimenter's hypothesis about the interaction is correct. For example, the cell sizes may be as hypothesized, but with opposite signs; all signs may be as hypothesized, but the amounts may be wrong; or possibly the experimental hypothesis may be wrong about both the direction and magnitude of the interaction effects.

In summary, a nonsignificant F-ratio tells one almost nothing. A significant F-ratio is merely a hunting permit, with the interaction contrast and its associated t-test as weapon.

Advantages. A specific interaction contrast is pragmatic, powerful, robust and easy to compute. Moreover, as we have shown, it is unavoidable when evaluating interactions in a repeated-measures design.

Specific interaction contrasts allow the experimenter to test his interaction hypothesis exactly. Since a specific hypothesis is tested.

the  $\underline{t}$ -test performed is one-tailed and always more powerful than the comparable  $\underline{F}$ -ratio.

Only three basic assumptions are required of the data in order to test an interaction contrast: (1) independence of the <u>n</u> contrast scores, (2) a normal distribution, and (3) homogeneous variance. Independence is guaranteed by independent selection and scoring of the <u>n</u> subjects. The latter two assumptions—normality and homogeneous variance—are not guaranteed, but can be tested by using the Wilcoxon signed rank test in tandem with the <u>t</u>-test. As <u>n</u> increases, normality and homogeneous variance become irrelevant. Thus, specific interaction contrasts, in conjunction with appropriate rank-order tests, are robust.

#### References

- Box, G. E. P. Some theorems on quadratic forms applied 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, 1954, 25, 484-498.
- Davidson, M. L. Univariate versus multivariate tests in repeatedmeasures experiments. <u>Psychological Bulletin</u>, 1972, 77, 446-452.
- Dixon, W. J. (Ed.). <u>Biomedical computer programs</u>. Berkeley: University of California Press, 1975.
- Dunn, O. J. Confidence intervals for the means of dependent normally distributed variables. <u>Journal of the American Statistical</u>
  Association, 1959, 54, 613-621.
- Fisher, R. A., & MacKenzie, W. A. Studies in crop variation. II.

  The manurial response of different potato varieties. <u>Journal of Agricultural Science</u>, 1923, 13, 311-320.
- Gaito, J. Repeated measurements, designs and counterbalancing.
  Psychological Bulletin, 1961, 58, 46-54.
- Geisser, S. The predictive sample reuse method with applications.

  Journal of the American Statistical Association, 1975, 70, 320-328.
- Geisser, S., & Greenhouse, S. W. An extension of Box's results on the use of the F distribution in multivariate analysis. <u>Annals of Mathemat-ical Statistics</u>, 1958, <u>29</u>, 885-891.
- Greenhouse, S. W., & Geisser, S. On the methods in the analysis of profile data. Psychometrika, 1959, 24, 95-112.

- Huynh, H., & Feldt, L. S. Conditions under which mean square ratios in repeated measurements designs have exact F-distribution.

  Journal of the American Statistical Association, 1970, 65, 1582-1589.
- Lana, R. E., & Lubin, A. The effect of correlation on the repeated measures design. <u>Educational and Psychological Measurement</u>, 1963, 23, 729-739.
- McNemar, Q. Psychological statistics. New York: Wiley, 1962.
- Miller, R. G., Jr. <u>Simultaneous statistical inference</u>. New York: McGraw-Hill, 1966.
- Mosier, C. I. Problems and designs of cross-validation. <u>Educational</u>
  and <u>Psychological Measurement</u>, 1951, <u>11</u>, 5-11.
- Mostelier, F., & Tukey, J. W. Data analysis, including statistics.

  In G. Lindzey & E. Aronson (Eds.), <u>Handbook of social psychology</u>

  (Vol. 2). Reading, Mass: Addison-Wesley, 1968.
- Plackett, R. L. Models in the analysis of variance. <u>Journal of the</u>
  Royal Statistical Society, 1960, 22, 195-217.
- Rouanet, H., & Lepine, D. Comparison between treatments in a repeatedmeasurement design: ANOVA and multivariate methods. <u>British Journal of</u> Mathematical and Statistical Psychology, 1970, 23, 147-163.
- Sternberg, S. The discovery of processing stages: Extension of Donder's method. In W. G. Koster (Ed.), <u>Attention and performance II</u>.

  Amsterdam: North Holland, 1969. (Reprinted from <u>Acta Psychologia</u>, 1969, 30, 276-315.
- Taylor, D. A. Stage analysis of reaction time. <u>Psychological Bulletin</u>, 1976, 83, 161-191.

- Tharp, V. K. Alcohol and response selection: An information processing analysis. (Doctoral dissertation, University of Oklahoma, 1975).

  <u>Dissertation Abstracts International</u>, 1975, 36, 6438B. (University Microfilms No. 76-14,096).
- Tharp, V. K., Rundell, O. H., Lester, B. K., & Williams, H. L. Alcohol and information processing. Psychopharmacología, 1974, 40, 33-52.
- Votaw, D. F., Jr. Testing compound symmetry in a normal multivariate distribution. Annals of Mathematical Statistics, 1947, 19, 447-473.
- Winer, B. J. <u>Statistical principles in experimental design</u>. New York: McGraw-Hill; 1st ed., 1962, 2nd ed., 1971.

#### Footnotes

This research was supported by the Department of the Navy, Bureau of Medicine and Surgery, under Work Unit Numbers MR000.01.01.6011 and MR041.01.03-0152. The views presented in this paper are those of the authors. No endorsement by the Department of the Navy has been given or should be inferred.

We thank S. Geisser, H. Williams, J. Callan, and the staff of the Psychophysiology Division at the Naval Health Research Center for their helpful comments during the preparation of this manuscript.

Address requests for reprints to: Van K. Tharp, Jr., Ph.D., Naval Health Research Center, San Diego, California 92152.

<sup>1</sup>A contrast is a weighted combination of scores where the sum of the weights equals zero.

<sup>2</sup>Since there are only two levels of drug (i.e., 1 df), one can subtract baseline from alcohol scores to simplify the computational procedures for computing residuals from additivity. We have not simplified in the above example in order to show how such residuals are computed with both rows and columns. Finding residuals for the simplified scores is a straightforward generalization of the example given.

<sup>3</sup>This procedure is equivalent to the method detailed for finding residuals from additivity.

We assume that the interaction effect for the ijth cell,  $\tau_{ij}$ , is a constant for all  $\underline{n}$  subjects. Any interaction with the subjects is thrown into the error deviance.

SWhen the specific interaction hypothesis is correct, under most circumstances the contrast weights will be proportional to the actual

interaction effect. For example, we recommend smoothing the contrast weights, but not the residuals from additivity for each subject, to single digit numbers. Thus, the theoretical maximum value can be stated more accurately as being  $\underbrace{k}_{i}\sum_{j}\omega_{ij}^{2}$ , where  $\underline{k}$  is the constant of proportionality.

<sup>6</sup>Each <u>t</u>-ratio may have a different between-subjects variance in this example. The completely general ANOVA demands that every <u>F</u>-ratio and <u>t</u>-ratio for interaction must have the same error variance to comply with the assumption of homogeneous variance.

<sup>7</sup>The routine application of stepwise multiple regression (e.g., Dixon, 1975) to the matrix of vectors representing the hypotheses will guarantee linear independence and thus eliminate redundancy. If there are  $\underline{\mathbf{h}}$  degrees of freedom for the interaction deviance in the ANOVA, then one can construct  $\underline{\mathbf{h}}$  vectors that are mutually orthogonal to one another.

"compound symmetry" property (Votaw, 1948). Compound symmetry holds when, given p repeated-measures, the p variances are equal and the p(p-1)/2 correlations are identical. Rouanet and Lepine (1970) have shown that compound symmetry is sufficient but not necessary to guarantee the validity of all F-ratios in a univariate ANOVA of a repeated-measures design. To confuse the issue further, the Greenhouse and Geisser 1959 article has a slip (p. 95): the word necessary was applied to the compound symmetry model rather than sufficient, as was implied by their first paper (Geisser and Greenhouse, 1958). This slip was copied by Lana and Lubin (1963), Gaito (1961), Winer (1962) and others. Winer corrected this slip in his second edition (1971, pp. 282-283).

<sup>9</sup>Rouanet and Lepine (1970), as well as Huynh and Feldt (1970), verified that when the circularity property holds, the Box epsilon criterion equals unity.

10 If some of the measures with small variances have very high correlations with one another, and if the remaining measures have high variance and low inter-correlations with all other measures, then the multivariate ANOVA will unerringly pick out the linear compound that gives maximum weight to the differences between the means of the highly correlated measures. The univariate ANOVA gives equal weight to all differences, and so may end up with a nonsignificant result.

Table 1
Verbal RT in Milliseconds, N=1

	Alcohol	Baseline	
WF <sub>1</sub>	901	818	$\bar{x}_1 = 859.5$
WF <sub>2</sub>	932	854	$\overline{X}_2 = 893$
WF <sub>3</sub>	949	870	$\bar{x}_3 = 909.5$
WF4	1109	982	$\bar{X}_4 = 1045.5$
WF <sub>5</sub>	1733	1205	$\overline{X}_5 = 1469$
	$\bar{x}_{1} = 1124.8$	$\overline{X}_{2} = 945.8$	X = 1035.3

Table 2
Residuals From Additivity For One Subject

	Alcoho1	Baseline	Row Effects
WF <sub>1</sub>	-48.0	48.0	-175.8
WF <sub>2</sub>	-50.5	50.5	-142.3
WF <sub>3</sub>	-50.0	50.0	-125.8
WF4	-26.0	26.0	10.2
WF <sub>5</sub>	174.5	-174.5	433.7
Column Effects	89.5	-89.5	x = 1035.3

Table 3

A Set of Theoretical Contrast Weights

	Alcohol	Baseline	Sum
WF <sub>1</sub>	-1	1	0
WF <sub>2</sub>	-1	1	0
WF <sub>3</sub>	-1	1	0
WF4	0	0	0
WF <sub>5</sub>	3	- 3	0
	0	0	

Table 4

Contrast Weights Obtained Empirically

from the Analysis Group

#### Condition

	Baseline	Alcohol	
WF <sub>1</sub>	-1	1	
WF <sub>2</sub>	-5	5	
WF <sub>3</sub>	1	-1	
WF4	-4	4	
WF <sub>5</sub>	9	-9	

Table 5 Contrast Scores

Subjects	Using Hypothesized Weights	Using Analysis Group Weights
Υ1	648	2,458
Y <sub>2</sub>	3,101	8,293
Y <sub>3</sub>	2,055	6,547
Y <sub>4</sub>	4,446	13,388
Y <sub>5</sub>	438	1,096
Y <sub>6</sub>	3,616	9,800
٧,	271	2,249
Y <sub>8</sub>	580	1,888
Y9	-524	-1,594
Y <sub>10</sub>	155	-15
Y <sub>11</sub>	530	1,528
Y <sub>12</sub>	818	618
Ÿ	1,344.500	3,854.667
sy	1,574.214	4,570.038
t	2.959	2.922
P .	.006	.007
r <sup>2</sup>	.443	.437

Table 6

A Golden-Mean Set of Contrast Weights

	Alcohol	Baseline
WF <sub>1</sub>	-1	1 .
WF <sub>2</sub>	0	0
WF <sub>3</sub>	2	-2
WF4	0	0
WF <sub>5</sub>	-1	1

#### Appendix

Some Assumptions of Multivariate Analysis of Variance

- Independence of Subjects. Each subject was selected independently of any other subject.
- (2) Normality. Each of the p measures has a marginal normal distribution.
- (3) Homogeneity of Residual Variance (Homoscedasticity). When we predict one of the p measures from any linear component of the remaining (p-1) measures, all errors of prediction must have the same variance.
- (4) <u>Linearity</u>. Each of the p measures has a linear regression on any weighted combination of the remaining (p-1) measures.
- (5) Homogeneous Universe. Each subset of subjects in the sample has the same covariance matrix as any other subset.
- (6) Well-conditioned Covariance Matrix. The determinant of the p by p covariance matrix is clearly greater than zero.

SECURITY CLASSIFICATION OF THIS PAGE (When Date Entered)

1. REPORT NUMBER 2 GOVT ACC	BEFORE COMPLETING FO
	ESSION NO. 3. RECIPIENT'S CATALOG NUMBER
76-78	
4 TITLE (and Subtitle)	S TYPE OF REPORT & PERIOD CO
Specific Interaction Contrasts: A Statist	ical (9 Interim rept.
Tool for Repeated-Measures Designs .	FEGERAMING ONG. REPORT OF
7. AUTHOR(4)	6. CONTRACT OR GRANT NUMBER
Van K./Tharp, Jr. Ardie/Lubin	61153W
9. PERFORMING ORGANIZATION NAME AND ADDRESS	10. PROGRAM ELEMENT, PROJECT,
Naval Health Research Center	
P.O. Box 85122	MRODE 27 01, 6011 and
THE SHUTA PRICE NO PICE NAME AND ADDRESS	MR041 01 03 0152
TO THE PROPERTY OF THE HAME AND ADDRESS	Dec 276
Naval Medical Research & Development Comman	nd 19 NUMBER OF PAGES
Bethesda, MD 20014	30
14 MONITORING AGENCY NAME & ADDRESS/II different from Controlle	as office) is security class for this report
Department of the Navy	UNCLASSIFIED
Bureau of Medicine and Surgery	150 DECLASSIFICATION DOWNGRA
Washington, D. C. 20372	
/ 11 \ ]	2.0
61150	
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, II	
17. DISTRIBUTION STATEMENT (of the abotract entered in Black 20, 11	different from Report)
17. DISTRIBUTION STATEMENT (of the abotract entered in Block 20, 11	different from Report)
	different from Report)
17. DISTRIBUTION STATEMENT (of the abetract entered in Black 20, 11	different from Report)
17. DISTRIBUTION STATEMENT (of the abetract entered in Black 20, 11	different from Report)
17. DISTRIBUTION STATEMENT (of the abetract entered in Block 20, 11	different from Report)
12. DISTRIBUTION STATEMENT (of the abotract entered in Block 20, II  14 NAVHLTHRS  18. SUPPLEMENTARY MOTES	CHC-76-78
12. DISTRIBUTION STATEMENT (of the abotract entered in Block 20, II  14 NAVHLTHRS  18. SUPPLEMENTARY MOTES	CHC-76-78
12. DISTRIBUTION STATEMENT (of the abotract entered in Block 20, II  14 NAVHLTHRS  18. SUPPLEMENTARY MOTES	CHC-76-78
12. DISTRIBUTION STATEMENT (of the abotract entered in Block 20, II  14 NAVHLTHRS  18. SUPPLEMENTARY MOTES	different from Report)
12. DISTRIBUTION STATEMENT (of the aborract entered in Block 20, 11  14	CHC-76-78  CHC-76-78  RØ41Ø1
12. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, 11  14	CHC-76-78  CHC-76-78  RØ41Ø1  RØ41Ø1
19. KEY HORDS (Continue on reverse side II necessary and identity by Side (U) Experimental psychologists often use	CHC-76-78  CHC-76-78  RØ41Ø1  RØ41Ø1  emultifactor repeated-measure
12. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, 11  14	CHC-76-78  CHC-76-78  RØH1Ø1  RWH1Ø1  emultifactor repeated-measure important effects to be assessed a valuate such interactions: (A)
12. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, II  14. NAVHLTHRS  15. KEY BORDS (Continue on reverse side II necessary and identity by Mineractions Contrasts  26. ABSTRACT (Continue on reverse side II necessary and identity by Mineractions (U) Experimental psychologists often use	CHC-76-78  CHC-76-78  RØH1Ø1  RWH1Ø1  emultifactor repeated-measure important effects to be assessed a riance (ANOVA), with (probably)

DD . TORM 1473

EDITION OF I NOV 45 IS OBSOLETE S/N 0102-014-6601

UNCLASSIFIED SECURITY CLASSIFICATION OF THIS PAGE (When Date Enter

391 642

CURITY CLASSIFICATION OF THIS PAGE(When Date Entered)

20. (Continued)

ANOVA; and (5) specific interaction contrasts. We show that no matter which of the above paths is chosen, the careful experimenter must compute specific interaction contrasts (i.e., <u>t</u>-tests). A worked example is given.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Date Entered)