



---

Defining Johnson-Neyman Regions of Significance in the Three-Covariate ANCOVA Using  
Mathematica

Author(s): Steve Hunka and Jacqueline Leighton

Source: *Journal of Educational and Behavioral Statistics*, Vol. 22, No. 4 (Winter, 1997), pp. 361-387

Published by: [American Educational Research Association](#) and [American Statistical Association](#)

Stable URL: <http://www.jstor.org/stable/1165228>

Accessed: 01-09-2015 18:03 UTC

---

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



*American Educational Research Association and American Statistical Association* are collaborating with JSTOR to digitize, preserve and extend access to *Journal of Educational and Behavioral Statistics*.

<http://www.jstor.org>

## **Defining Johnson-Neyman Regions of Significance in the Three-Covariate ANCOVA Using Mathematica**

**Steve Hunka and Jacqueline Leighton**  
*University of Alberta*

**Keywords:** *Johnson-Neyman ANCOVA, Mathematica, nonhomogeneous regression*

*The Johnson-Neyman ANCOVA, used when nonhomogeneity of within-group regressions is present, poses special computational and plotting problems when three covariates are used. These problems can be overcome by using (a) an appropriate design and contrast matrix for the general linear model and (b) the Mathematica software system of computation to handle the symbolic and graphical processing requirements. Four-dimensional graphical representation of the polynomials which result are contour plotted in a three-dimensional space in order to define the regions of significance for contrasts. It is also shown that for some values of the covariate orthogonal contrasts are produced.*

Fisher's (1932, 1935) purpose in developing the analysis of covariance (ANCOVA) was based on the notion that the introduction of a covariate into the analysis of variance (ANOVA) would increase the precision of experimental results by reducing the mean square error term and by adjusting the treatment effects for a confounding variable. Most researchers today view the ANCOVA as an analysis that allows testing of treatment effects after an adjustment has been made for group differences in the covariate. In this context the covariate is not necessarily viewed as a confounding variable the effects of which are to be controlled for statistically, but rather as a component of the dependent variable. For example, reading readiness of kindergarten children in different schools might be adjusted for differences in socioeconomic background.

In addition to the assumptions made in the ANOVA (Kirk, 1982), the usual procedure for ANCOVA assumes homogeneity among the groups in the slopes of the regression lines relating the covariate and the dependent variable (Huitema, 1980). If this assumption is not tenable, the ANCOVA procedure can produce misleading results. The basis for such misleading results is explained with reference to Figure 1, in which the lines plotted represent within-group regressions lines. If a projection of the group means is made onto each axis, a difference in the observed means for the two groups on the dependent and the covariate variables is observed. If the dependent variable means are adjusted to a covariate value of 0, then these means will be  $V$  for Group 1 and  $S$  for Group 2, with  $S > V$  (i.e., Group 2 is superior to Group 1), the order being the reverse of the order of the observed mean values for  $Y$ . If an adjustment is made to a

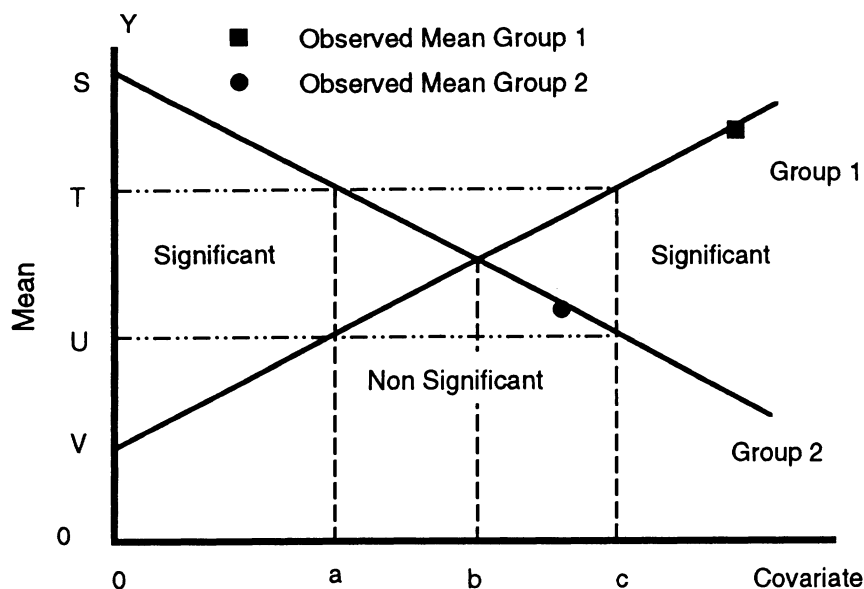


FIGURE 1. *Nonhomogeneous regressions for two groups*

covariate value of  $b$ , no differences exist between the adjusted group means. Experimental results that suggest the sort of group differences represented in Figure 1 can vary from nonsignificant to significant and can favor different groups, depending upon the value of the covariate to which the groups are adjusted. If the regression lines are parallel—that is, if the assumption of homogeneity of the regression slopes is tenable—then it makes no difference to what value of the covariate the means are adjusted, because the adjustment is identical for all values of the covariate. Many computer program packages which carry out an ANCOVA implicitly adjust the group means to a covariate value of 0 or do not tell the researcher to what value the adjustment was made.

In the case of nonhomogeneous regression slopes, the problem to be solved is to find the values of  $a$  and  $c$  in Figure 1 such that they separate the boundaries of statistical significance and nonsignificance for a nominal value of  $\alpha$  in testing the differences between the adjusted group means of the dependent variable. The values of  $a$  and  $c$  are said to delimit *regions of significance*. For a single covariate the values of  $a$  and  $c$  can be found by solving a quadratic in one unknown. For two or more covariates the problem becomes more difficult, because the regions of significance are defined by point sets of covariate values. A modification to the usual ANCOVA to accommodate nonhomogeneity of regression slopes was first detailed by Johnson and Neyman (1936) and is usually referred to as the Johnson-Neyman ANCOVA or the Johnson-Neyman procedure.

Recent critical discussion of the ANCOVA including the Johnson-Neyman procedure (Maxwell & Delaney, 1990; Maxwell, O'Callaghan, & Delaney, 1993; Rutherford, 1992) suggests that this method of analysis continues to be of

importance to researchers in the behavioral sciences. Maxwell et al. have suggested that one reason for a limited use of the procedure might be that researchers generalize the restrictive assumptions and problems associated with experiments in which observations are not assigned randomly to those experiments in which observations are truly randomized. This is undoubtedly true; however, an examination of the research publications describing the Johnson-Neyman procedure suggests that complexity of computation is an additional constraint. Early research publications (Johnson & Fay, 1950; Johnson & Hoyt, 1947; Johnson & Neyman, 1936; Koenker & Hansen, 1942), as well as early texts (Johnson & Jackson, 1959; Walker & Lev, 1953) and more recent texts (Huitema, 1980; Pedhazur, 1982; Searle, 1987), also indicate by their examples that the computations required for the Johnson-Neyman procedure would be difficult for most researchers to carry out.

Early attempts to ease the computational difficulties (Butsch, 1944; Koenker & Hansen, 1942) simply transformed the equations originally given by Johnson and Neyman into a series of computational steps, usually only for the two-group and one- or two-covariate case. It is important to note that Carroll and Wilson (1970) indicate that the Johnson and Fay (1950) and Koenker and Hansen (1942) articles "exhibit more or less serious errors" (p. 122), which have likely been perpetuated in subsequent papers.

To ameliorate the computational difficulties, specialized computer programs have been developed (Borich, Godbout, & Wunderlich, 1976; Carroll & Wilson, 1969, 1970; Ceurvorst, 1979; Karpman, 1980; Kush, 1986; Lautenschlager, 1987). These programs, written in programming languages such as APL, Telecomp, and FORTRAN, handle the Johnson-Neyman one-way ANCOVA with one or two covariates only, appear to follow the algebraic solutions provided by earlier papers, and have very limited plotting capabilities (if any) for defining regions of significance. These programs do not appear to have found widespread use. Using GAUSS, a more advanced programming language and system of computation, Schafer and Wang (1991) provide a solution for the one- and two-covariate case for two groups, simplifying the computations by using the root solving capabilities in GAUSS. The separate programs for the one- and two-covariate case do not appear easily generalizable to more complex cases. In addition, they do not appear to provide plotting of the regions of significance directly, and require as input the means, regression parameters, standard deviations, and correlations. Recognizing the widespread use of the SPSS, BMDP, and SAS statistical program packages, Karpman (1983, 1986) shows how these packages can be used to carry out the Johnson-Neyman procedure for the two-group and one- or two-covariate case by using the special computational options provided by these packages.

More recently, very powerful languages and computational systems such as Mathematica (Wolfram, 1991) and Maples (Abell & Braselton, 1994) have been developed which explicitly provide the numerical, symbolic, and graphical processing capabilities required for the Johnson-Neyman procedure. Using

Mathematica and the general linear model (Searle, 1971), Hunka (1994, 1995) and Leighton (1995) have illustrated how the computational and graphing problems encountered in defining the regions of significance in the one- and two-covariate case can be solved. Although these illustrations are made using Mathematica's interactive computational mode, a "canned" program is available from the Wolfram Research Web site: <http://www.wolfram.com/cgi-bin/MathSource/Enhancements/Statistics/0207-953>.

### **The Three-Covariate Johnson-Neyman Case**

A solution for the three-covariate case has been illustrated by Johnson and Hoyt (1947) for the case of two groups. As might be expected, the solution given by Johnson and Hoyt parallels the original algebraic, desk calculator-like computational procedures of Johnson and Neyman with additional complexities introduced by the requirement to plot regions of significance in a three-dimensional space. Johnson and Hoyt's solution is to define a likelihood ratio function for testing adjusted group effects for a specific level of alpha ( $\alpha = .01$ ) with reference to the incomplete beta function (Mood, 1950). Contour plots are made of the polynomial provided by the likelihood ratio function for selected values of one of the covariates. Contour plotting—that is, plotting for a specific value of the polynomial based on a fixed level of  $\alpha$ —and setting one of the covariates to a constant value reduces the plot of the region of significance to two dimensions.

Huitema (1980) illustrates a matrix solution for a two-group and three-covariate ANCOVA which provides a considerable simplification over the algebraic approach provided by Johnson and Hoyt (1947), but rather than defining regions of significance, the example shows only how to evaluate the significance of group effects adjusted for a specific set of covariate values. Since the latter computations are not done in the full context of the general linear model (although a set-to-zero design matrix is used [Searle, Speed, & Henderson, 1981]), Huitema suggests that pairwise comparisons for a specific covariate point be carried out when more than two groups exist. Pedhazur (1982) does not show how the Johnson-Neyman ANCOVA can be handled when more than two covariates are available. Other than the original Johnson and Hoyt publication, no other published reports could be found that (a) give a direct solution to the three-covariate case which includes the definition of the regions of significance, (b) use the computational approaches and efficiencies of the general linear model, and (c) reduce the computational labor required by using the more advanced computational software systems available today. We shall show how these three points can be accommodated by using the software system of computation called Mathematica.

### **The General Linear Model**

It has been shown previously (Hunka, 1995) for one and two covariates that the solution for the one-way Johnson-Neyman ANCOVA can be cast entirely in

the form of the general linear model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{Y}$  is the dependent variable of order  $(N, 1)$ , with  $N$  being the number of observations;  $\mathbf{X}$  is the design matrix of order  $(N, c)$ , where  $c$  is the number of parameters to be estimated;  $\boldsymbol{\beta}$  represents the population parameters to be estimated, of order  $(c, 1)$ ; and  $\boldsymbol{\epsilon}$  is the error vector of order  $(N, 1)$ . The usual least square estimate of  $\boldsymbol{\beta}$  can be obtained by  $\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$  under the condition that the columns of the design matrix  $\mathbf{X}$  are linearly independent. A thorough exposition of the general linear model can be obtained from Searle (1971, 1987). Casting the solution in the form of the general linear model requires that the researcher understand the nature of the parameters being estimated with the design matrix used and be able to construct an appropriate contrast matrix consistent with the nature of the design matrix for testing group effects adjusted for covariate effects. In the solution proposed here, the contrast matrix will be constructed in symbolic form and used to derive a polynomial that defines the region of significance.

### Design Matrices

Unfortunately, there is little consistency in the textbooks used in the social sciences, and in the manuals of statistical program packages, in the form and nomenclature used to describe design matrices. Such inconsistencies make it more difficult for researchers to specify an appropriate contrast matrix, which is essential to the method described here for solving the Johnson-Neyman ANCOVA problem.

In Figure 2 four design matrices are illustrated for a three-group one-way ANOVA. Design matrix  $\mathbf{X}_1$  allows estimation of the parameters  $\mu$ ,  $\mu_1$ ,  $\mu_2$ , and  $\mu_3$ ; however, because the columns are linearly dependent and the inverse of  $\mathbf{X}'\mathbf{X}$  does not exist, the usual least squares solution for the parameter estimates given by  $\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$  cannot be used. Three different procedures for removing the linear dependency of columns in  $\mathbf{X}_1$  produce the design matrices  $\mathbf{X}_2$ ,  $\mathbf{X}_3$ , and  $\mathbf{X}_4$ .  $\mathbf{X}_2$  can be obtained by dropping the first column of  $\mathbf{X}_1$  to provide for the estimates  $\mu_1$ ,  $\mu_2$ , and  $\mu_3$ .  $\mathbf{X}_2$  is frequently referred to as a  $\mu$ -model or cell means design matrix.  $\mathbf{X}_3$  is a  $\Sigma$ -restricted design matrix and allows estimation of the parameters  $\mu$ ,  $\alpha_1$ , and  $\alpha_2$ , such that  $\Sigma\alpha_j = 0$ .  $\mathbf{X}_4$  is a set-to-zero design matrix that provides the estimates of  $\mu$ ,  $\delta_1$ , and  $\delta_2$ , with  $\delta_1 = \mu_1 - \mu_3$  and  $\delta_2 = \mu_2 - \mu_3$ ;

$$\mathbf{X}_1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{X}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{X}_3 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \end{bmatrix} \quad \mathbf{X}_4 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

FIGURE 2. Four equivalent ANOVA design matrices

that is, the effects for Groups 1 and 2 are their group means deviated about the mean of the last group.  $X_4$  is used frequently when only two groups are present, since it provides directly the group contrast of interest. Design matrices  $X_2$ ,  $X_3$ , and  $X_4$  all produce the same results in the sense that the sums or squares attributed to the model are the same; however, the parameter estimates will have different values and interpretations. There are many other design matrices which could be formed. Any set of three column vectors spanning the same space as  $X_1$  would provide an equivalent design matrix.

In the case of the ANCOVA having one covariate, a design matrix could be formed by appending the covariate data to the right of any one of the design matrices  $X_2$ ,  $X_3$ , and  $X_4$ . The covariate data may take the same general form as the matrices in Figure 2 and are illustrated in Figure 3 using hypothetical values.

Matrix  $Z_1$  contains a linear dependency of the same nature as in  $X_1$ . This linear dependency can be removed in the same manner as done for  $X_1$ , producing  $Z_2$ ,  $Z_3$ , and  $Z_4$ . Any one of the design matrices  $Z_2$ ,  $Z_3$ , and  $Z_4$  can be appended to any one of the design matrices  $X_2$ ,  $X_3$ , and  $X_4$  to form an appropriate design matrix for the ANCOVA. For example  $X = [X_3 | Z_3]$  has the form

$$X = \begin{bmatrix} 1 & 1 & 0 & 2 & 2 & 0 \\ 1 & 1 & 0 & 3 & 3 & 0 \\ 1 & 0 & 1 & 7 & 0 & 7 \\ 1 & 0 & 1 & 5 & 0 & 5 \\ 1 & -1 & -1 & 4 & -4 & -4 \\ 1 & -1 & -1 & 6 & -6 & -6 \end{bmatrix}$$

and provides for parameter estimates  $\mu$ ,  $\alpha_1$ ,  $\alpha_2$ ,  $\beta$ ,  $\delta_1$ , and  $\delta_2$  for which  $\alpha_3 = -\alpha_1 - \alpha_2$  and  $\delta_3 = -\delta_1 - \delta_2$  and the usual within-group regression parameters are given by  $\beta_1 = \beta + \delta_1$ ,  $\beta_2 = \beta + \delta_2$ , and  $\beta_3 = \beta + \delta_3$ . Columns 5 and 6 of  $X$  are frequently referred to as providing group-by-covariate parameter estimates. That is, these parameters are treated as interaction terms in the same sense as interaction terms in a two-factor ANOVA.

$$Z_1 = \begin{bmatrix} 2 & 2 & 0 & 0 \\ 3 & 3 & 0 & 0 \\ 7 & 0 & 7 & 0 \\ 5 & 0 & 5 & 0 \\ 4 & 0 & 0 & 4 \\ 6 & 0 & 0 & 6 \end{bmatrix} \quad Z_2 = \begin{bmatrix} 2 & 0 & 0 \\ 3 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 4 \\ 0 & 0 & 6 \end{bmatrix} \quad Z_3 = \begin{bmatrix} 2 & 2 & 0 \\ 3 & 3 & 0 \\ 7 & 0 & 7 \\ 5 & 0 & 5 \\ 4 & -4 & -4 \\ 6 & -6 & -6 \end{bmatrix} \quad Z_4 = \begin{bmatrix} 2 & 2 & 0 \\ 3 & 3 & 0 \\ 7 & 0 & 7 \\ 5 & 0 & 5 \\ 4 & 0 & 0 \\ 6 & 0 & 0 \end{bmatrix}$$

FIGURE 3. Four equivalent design matrices for covariate data



For the purpose of illustrating the solution for the three-covariate Johnson-Neyman ANCOVA, we shall employ a design matrix  $\mathbf{X}$  of the form  $\mathbf{X} = [\mathbf{X}_2 | \mathbf{Z}_2]$ , but with columns reordered so that the estimated parameters are consecutive within groups. The construction of  $\mathbf{X}$  of this form and that of any associated contrast matrix are straightforward, and the parameter estimates are easily interpreted, because they provide estimates of the adjusted group means and regression coefficients directly.

### Sum of Squares for Contrasts

Searle (1971) shows that the sum of squares  $SS_k$  associated with any contrast matrix  $\mathbf{K}'$  of order  $(r, c)$  and rank  $r$ , for testing the hypothesis  $H_0: \mathbf{K}'\boldsymbol{\beta} = \mathbf{0}$ , can be obtained by

$$SS_k = (\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}.$$

In the case of three covariates ( $P$ ,  $Q$ , and  $S$ ), in order to test the means of three groups adjusted to the covariate values of  $p$ ,  $q$ , and  $s$  for an ANCOVA in which the parameter estimates  $\boldsymbol{\beta}$  are given as  $\mu_1, \beta_{11}, \beta_{12}, \beta_{13}, \mu_2, \beta_{21}, \beta_{22}, \beta_{23}, \mu_3, \beta_{31}, \beta_{32}$ , and  $\beta_{33}$ —that is, the basic form of the design matrix is  $\mathbf{X} = [\mathbf{X}_2 | \mathbf{Z}_2]$ —the following contrast matrix is used:

$$\mathbf{K}' = \begin{bmatrix} 1 & p & q & s & -1 & -p & -q & -s & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & p & q & s & -1 & -p & -q & -s \end{bmatrix}.$$

The hypothesis  $\mathbf{K}'\boldsymbol{\beta} = \mathbf{0}$  using this contrast matrix is

$$\begin{bmatrix} (\mu_1 + p\beta_{11} + q\beta_{12} + s\beta_{13}) - (\mu_2 + p\beta_{21} + q\beta_{22} + s\beta_{23}) \\ (\mu_2 + p\beta_{21} + q\beta_{22} + s\beta_{23}) - (\mu_3 + p\beta_{31} + q\beta_{32} + s\beta_{33}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Furthermore, any elementary row operation to  $\mathbf{K}'$  results in testing an equivalent hypothesis and produces an identical sum of squares. If the number of groups is greater than three, then  $\mathbf{K}'$  has more columns and rows, but its general form remains the same as illustrated. If the design matrix is of the form  $\mathbf{X} = [\mathbf{X}_3 | \mathbf{Z}_2]$ , but with columns reordered to provide estimates  $\mu, \alpha_1, \beta_{11}, \beta_{12}, \beta_{13}, \alpha_2, \beta_{21}, \beta_{22}, \beta_{23}, \beta_{31}, \beta_{32}$ , and  $\beta_{33}$ , then an acceptable contrast matrix is

$$\mathbf{K}'_1 = \begin{bmatrix} 0 & 1 & p & q & s & -1 & -p & -q & -s & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 2 & p & q & s & -p & -q & -s \end{bmatrix}.$$

In order for a test of the hypothesis to be significant, the following relationship must be true:

$$\frac{[(\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}]/df_k}{MS_e} \geq F_{\alpha; df_k, df_e}, \quad (1)$$

where  $df_k$  is the degrees of freedom ( $r$ ) for the contrast, and  $df_e$  the degrees of freedom for the mean square error term  $MS_e$ . Since the right-hand term may



involve a simple modification to accommodate variations such as those required for simultaneous statements about the covariate values or the use of the Bonferroni approach (Huitema, 1980; Rogosa, 1981), and has no material importance to the method of defining the regions of significance proposed here, the regular value for  $F$  will be employed.

Equation 1 can be rearranged in the following form:

$$[(\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}] = df_k MS_e F_{\alpha; df_k; df_e} \quad (2)$$

which, when true for specific values of covariate point sets  $(p, q, s)$ , represents significance just being achieved at a nominal level of  $\alpha$ . The matrix terms of (2) generate a polynomial with coefficients in symbolic form, for which the unknown values  $p, q$ , and  $s$  of the covariates, embedded in the contrast matrix  $\mathbf{K}'$ , are sought such that (2) is true, that is, equal to the product  $(df_k MS_e F_{\alpha; df_k; df_e})$  for some specified  $\alpha$  level. Since this polynomial has three unknowns, a plot of the expression for a range of values of  $p, q$ , and  $s$  in addition to the evaluated result would require a four-dimensional space. Using the Mathematica computer software, the expression is plotted for a value equal to the product  $(df_k MS_e F_{\alpha; df_k; df_e})$  for a selected  $\alpha$  level. This contour plot occupies three dimensions with each axis representing a covariate. Any coordinate point  $(p, q, s)$  that falls on the surface of this representation, or within its enclosing volume, will be significant at the  $\alpha$  level selected or less.

## Mathematica: A System of Computation

Operators in Mathematica can be executed step-by-step interactively, or they can be grouped together as functions and executed as a single unit. Mathematica statements are for the most part self-explanatory, because they use the usual symbols for arithmetic operators, and numerical processes are usually specified using customary mathematical names. For example, `Inverse[x]` provides the inverse of the matrix  $\mathbf{x}$ , whether in numeric or symbolic form. For the purpose of illustrating the computational steps required for the three-covariate case, step-by-step calculations using Mathematica are shown for the first example only. Where an operation may not be obvious, an explanation is given either in the introductory text or as a comment contained between the symbols  $(*)$  and  $(*)$  with the operation itself. In the examples that follow, operations given to Mathematica appear boldface in the Courier font, and the results lightface in the Courier font.

### Example 1: Johnson and Hoyt Data

The Johnson and Hoyt (1947) data consist of two different college physics classes with 111 and 257 observations, respectively. The dependent variable was the Mechanics section of the Cooperative Physics Test for College Students, and the three covariates were the scores on the American Council on Education Psychological Examination (ACEPE), proficiency in mathematics, and an honor point variable. The unadjusted group means are given in Table 1.

TABLE 1  
Means for Johnson-Hoyt data

Group	Mechanics	ACEPE	Mathematics	Honor point
1	11.928	87.640	31.081	11.586
2	26.393	92.397	56.074	12.689

From the summary data provided by Johnson and Hoyt (1947), and using a  $\mu$ -model design matrix ( $\mathbf{X}$ ) to estimate the parameters  $\mu_1$ ,  $\beta_{11}$ ,  $\beta_{12}$ ,  $\beta_{13}$ ,  $\mu_2$ ,  $\beta_{21}$ ,  $\beta_{22}$ , and  $\beta_{23}$ , with the parameters in  $\beta$  being the within-group regression coefficients, the matrix  $\mathbf{X}'\mathbf{X}$  is

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 111 & 9,728 & 3,450 & 1,286 & 0 & 0 & 0 & 0 \\ 9,728 & 905,694 & 307,220 & 117,560 & 0 & 0 & 0 & 0 \\ 3,450 & 307,220 & 118,846 & 42,408 & 0 & 0 & 0 & 0 \\ 1,286 & 117,560 & 42,408 & 20,084 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 257 & 23,746 & 14,411 & 3,261 \\ 0 & 0 & 0 & 0 & 23,746 & 2,338,412 & 1,349,410 & 319,746 \\ 0 & 0 & 0 & 0 & 14,411 & 1,349,410 & 823,945 & 191,223 \\ 0 & 0 & 0 & 0 & 3,261 & 319,746 & 191,223 & 53,417 \end{bmatrix}$$

and  $\mathbf{Y}'\mathbf{X}$  is

$$\mathbf{Y}'\mathbf{X} = [1,324 \quad 118,635 \quad 44,200 \quad 17,532 \quad 6,783 \quad 645,402 \quad 386,261 \quad 92,843].$$

The following steps are used to analyze this data:

- (1) A list structure is defined in Mathematica for the matrix  $\mathbf{X}'\mathbf{X}$  and for the vector  $\mathbf{Y}'\mathbf{X}$ .
- (2) The OLS regression parameters are calculated to confirm those reported by Johnson and Hoyt (1947).
- (3) The sum of squares is calculated, for testing the homogeneity of the regressions slopes.
- (4) A contour plot of the regression hyperplanes is made, to show the nature of the nonhomogeneity of the regression slopes.
- (5) The polynomial provided by (2) is calculated, and its contour plot made, to provide the region of significance.
- (6) The Johnson-Hoyt plots of the region of significance are reproduced.

#### Step 1: Defining the Data

The data  $\mathbf{X}'\mathbf{X}$  are defined as the list structure  $\mathbf{xx}$ , and  $\mathbf{X}'\mathbf{Y}$  as  $\mathbf{xy}$ . Although not shown, the normal response from Mathematica is to echo back the list structure if it is properly defined, or to indicate errors if it is not.

```
xx = {{111., 9728., 3450., 1286., .0, .0, .0, .0},
      {9728., 905694., 307220., 117560., .0, .0, .0, .0},
      {3450., 307220., 118846., 42408., .0, .0, .0, .0},
```

```
{1286.,117560.,42408.,20084.,.0,.0,.0,.0},
{.0,.0,.0,.0,257.,23746.,14411.,3261.},
{.0,.0,.0,.0,23746.,2338412.,1349410.,319746.},
{.0,.0,.0,.0,14411.,1349410.,823945.,191223.},
{.0,.0,.0,.0,3261.,319746.,191223.,53417.}]
xy = {1324.,118635.,44200.,17532.,6783.,645402.,
386261.,92843.]
```

### Step 2: OLS Regression Estimates

The OLS estimates of the regression parameters are calculated using  $\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ , with the result being an eight-element vector of the parameter estimates  $\mu_1, \beta_{11}, \beta_{12}, \beta_{13}, \mu_2, \beta_{21}, \beta_{22}$ , and  $\beta_{23}$ .

```
b = Inverse[xx].xy
```

```
{2.01566, 0.00103219, 0.192459, 0.331443, 10.2658,
0.0677165, 0.0856049, 0.399582}
```

The last results correctly reproduce those provided by Johnson and Hoyt (1947, p. 350).

### Step 3: Obtaining the Sum of Squares

The total sum of squares for the data is provided by Johnson and Hoyt (1947) as  $194,001 + 19,862 = 213,863$ . The sum of squares for the model is obtained using  $\mathbf{\beta}'\mathbf{X}'\mathbf{X}\mathbf{\beta}$ ,

```
b.xx.b
```

```
200610
```

and, by subtraction, the error sum of squares is  $213,863 - 200,610 = 13,253$ , which provides a mean square error of 36.814 with  $df_e = 360$ .

In order to calculate the sum of squares for testing the hypothesis that the difference in within-group regression slopes is zero in the population, a one-line Mathematica function defined as `ssk` is created. It has the following form

```
ssk[k_,xx_,pb_] :=
(k.pb).(Inverse[k.Inverse[xx].Transpose[k]]).(k.pb);
```

and takes the dummy arguments `k_`, the contrast matrix; `xx_`, the matrix  $\mathbf{X}'\mathbf{X}$ ; and `pb_`, a vector of the estimated parameters. This function calculates the result for  $SS_k = (\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}$ . Next, the contrast matrix to test  $H_0: \beta_{11} - \beta_{21} = 0$  is formed and included with  $\mathbf{X}'\mathbf{X}$  and the estimated parameters as the argument list for the function `ssk` to obtain the sum of squares for the contrast:

```
k1={{0,1,0,0,0,-1,0,0}}; (*matrix of order 1,8 *)
ssk[k1,xx,b]
```

```
149.927
```

The sum of squares of 149.927 has 1 degree of freedom. Although not illustrated using Mathematica,  $F = 4.073$ , and  $p \leq .044$ . In a similar fashion the test of  $H_0: \beta_{12} - \beta_{22} = 0$  produces  $F = 1.562$  and  $p \leq .212$ , and the test of  $H_0: \beta_{13} - \beta_{23} = 0$  produces  $F = .334$  and  $p \leq .560$ . Mathematica can provide the probabilities associated with the  $F$  distribution.

#### Step 4: Plotting the Regression Hyperplanes

Obtaining a visual indication of the regression hyperplanes is a little more complicated, because (a) the regression estimates are contained in an eight-element vector defined as the Mathematica variable **b** rather than as separate estimates for each group, and (b) each regression equation or hyperplane of the general form  $y = f(p, q, s)$  must be plotted for some specific value of  $y$ , that is, a contour plot. First the regression equation is formed by restructuring the vector of eight regression parameters into a matrix of order (2, 4) by the **Partition** operator. Then the first row is taken by the **Part** operator, and the result multiplied times the vector **{1, p, q, s}**. The result is shown below as a new variable defined as **regeqn1**

```
regeqn1=Part[Partition[b,4],1].{1,p,q,s}
```

```
2.01566 + 0.00103219 p + 0.192459 q + 0.331443 s
```

and is equivalent to the usual algebraic form  $y = 2.01566 + .00103219p + .192459q + .331443s$ .

In a similar fashion, the regression equation for the second group is defined as the variable **regeqn2**. The next step is to make a contour plot of each hyperplane for some specific value of  $y$  and then to display the plots together in a single representation. There is no problem in displaying the contour plot of each hyperplane separately for a specific value of  $y$ ; however, if the same value of  $y$  is selected for both plots, there is no guarantee for the range of the covariate values over which the plot is made that the same value is possible for both plots. In other words, for the equations of the two hyperplanes,  $y_1 = 2.01566 + .00103219p + .192459q + .331443s$  and  $y_2 = 10.2658 + .0677165p + .0856049q + .399582s$ , there is no guarantee that a value of, say, 10 can be produced by both equations varying  $p$ ,  $q$ , and  $s$  over the same integer range of reasonable values. The approach used here is to plot contour values equal to the mean of the dependent variable and then to adjust these values slightly if an intersection of the planes is desired. The Mathematica plot operation shown below plots the expression **regeqn1** for the ranges of  $p$  from 60 to 100 for covariate  $P$  (mathematics),  $q$  from 20 to 70 for covariate  $Q$  (ACEPE), and  $s$  from 5 to 15 for covariate  $S$  (honor point) for which the value of  $y_1$  is equal to 12 (**Contours**  $\rightarrow$  {12}).

```
ContourPlot3D[regeqn1,{p,60,100},{q,20,70},{s,5,15},
Contours->{12},Axes->True,
DisplayFunction->Identity,(*do not display*)
AxesLabel->{"p","q","s"}]
```

In a similar fashion, the regression hyperplane for the second group is obtained, but for a value of  $y_2 = 23$  (`Contours-> {23}`). The two plots can now be brought together and displayed as shown in Figure 4, in which the viewpoint has been modified to more clearly show the planes intersecting.

```
Show[%,%%, (*display last two results together*)
      DisplayFunction->$DisplayFunction,
      ViewPoint->{3.910,-0.150,1.760}]
```

#### Step 5: Forming and Plotting the Polynomial

In order to define the polynomial expressed by  $(\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}$ , the contrast matrix is defined in symbolic form, with the unknown values of the covariate to which an adjustment is to be made defined as  $p$ ,  $q$ , and  $s$ . The Mathematica function `ssk` used earlier to obtain the sum of squares using a contrast matrix in completely numerical form can also handle a contrast matrix defined with symbolic values. The result returned by the function `ssk` is a polynomial. Because the function `ssk` returns about 40 lines of output for the polynomial, the operation `Simplify[. . .]` asks Mathematica to use its own rules to simplify the polynomial. The polynomial expression is defined as the Mathematica variable `polyxpr`. Mathematica represents the result as the ratio of two polynomials in the three unknowns  $p$ ,  $q$ , and  $s$ , as shown below.

```
ks={{1,p,q,s,-1,-p,-q,-s}} (*symbolic contrast*)
polyxpr=Simplify[ssk[ks,xx,b]]

(149.927 (15306.6 + 247.439 p2 + 1. p - 396.496 q -
3.20479 p q + 2.56767 q2 + 252.837 s + 2.04363 p s -
3.2747 q s + 1.04411 s2))/
(15575.4 - 128.282 p + 1. p2 - 443.029 q - 0.64581
p q + 6.69302 q2 + 192.002 s - 1.85544 p s - 7.05243
q s + 12.5227 s2)
```

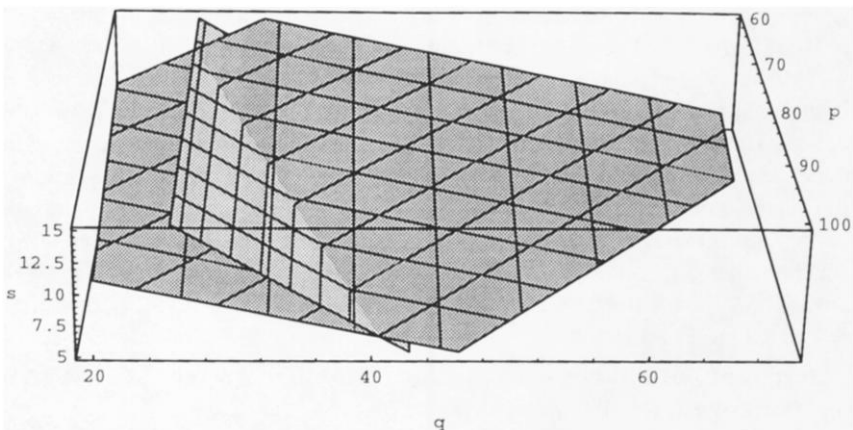


FIGURE 4. Regression hyperplanes for both groups plotted for fixed values of  $y_1$  and  $y_2$

The next step is to contour plot the polynomial for a resulting value of 246.859, formed by the product  $df_k MS_e F_{\alpha; df_k, df_e}$ , which is the minimum value required for significance at  $\alpha = .01$ , and then to project the representation onto planes bounding the figure.

```
<<Graphics`ContourPlot3D` (*load special Graphics
library*)
Shadow[ContourPlot3D[polyxpr, {p, 0, 100}, {q, -20, 60},
{s, -10, 20}, Contours->{246.859}], Axes->True,
AxesLabel->{"p", "q", "s"},
DisplayFunction->Identity,
FaceGrids->{{-1, 0, 0}, {0, 1, 0}, {0, 0, -1}},
DisplayFunction->$DisplayFunction]
```

In Figure 5, any combination of the covariate  $(p, q, s)$  that falls on the surface, or within the volume enclosed by this surface, is significant at  $\alpha = .01$  or less. Alternatively, any coordinate formed by  $(p, q, s)$  that simultaneously falls within the outer boundaries of the three projected regions will be significant at  $\alpha = .01$  or less.

#### Step 6: Reproducing the Johnson-Hoyt Plots

As illustrated by Johnson and Hoyt (1947), simplified plots of the region of significance can be obtained by plotting the polynomial for a given value of the

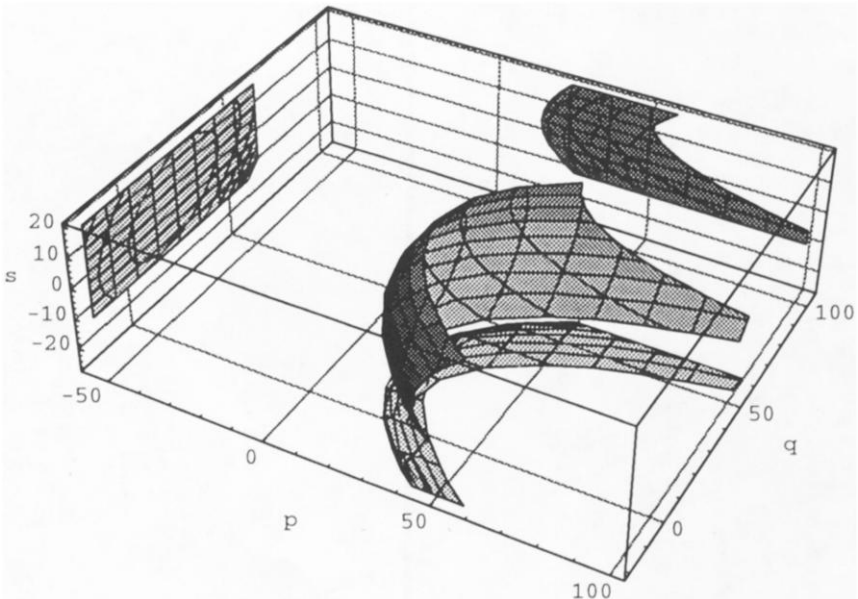


FIGURE 5. Region of significance and shadow projections

covariate  $s$ . Johnson and Hoyt first plot the region for the honor point value of  $s = 7.823426$  in raw score metric. By substituting this value into the polynomial expression `polyxpr` provided by Mathematica, the polynomial is reduced to two unknowns  $p$  and  $q$ . The polynomial can then be contour plotted in the form

$$[(\mathbf{K}'\mathbf{B})'[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B}] - df_k MS_e F_{\alpha; df_k; df_e} = 0.$$

To make the substitution (of  $s = 7.823426$  into the expression defined as `polyxpr`), use the Mathematica statement that follows. It requests that in the expression `polyxpr` the rule defined as `s->7.823426` replace (/.) all occurrences of  $s$  and that the result (which has only two unknowns,  $p$  and  $q$ ) be defined as `polyxpr1`.

```
polyxpr1=polyxpr /. s->7.823426
149.927 (17348.5 + 263.428 p2 + 1. p - 422.115 q -
3.20479 p q + 2.56767 q2)
-----
17844. - 142.798 p + 1. p2 - 498.203 q - 0.64581 p q +
6.69302 q2
```

To achieve significance at  $\alpha = .01$ , the sum of squares must be at least 246.859, so a three-dimensional plot is made of the expression `polyxpr1-246.859`. The result is given in Figure 6, in which the value of `polyxpr1-246.859` is labeled on the vertical axis as `ss-c` (where  $ss$  = sum of squares and  $c$  represents a constant, in this case 246.859).

```
Plot3D[(polyxpr1-246.859),{p,0,100},{q,-20,60},
PlotPoints->30,PlotRange->{-10,20},
AxesLabel->{"p","q","ss-c"}]
```

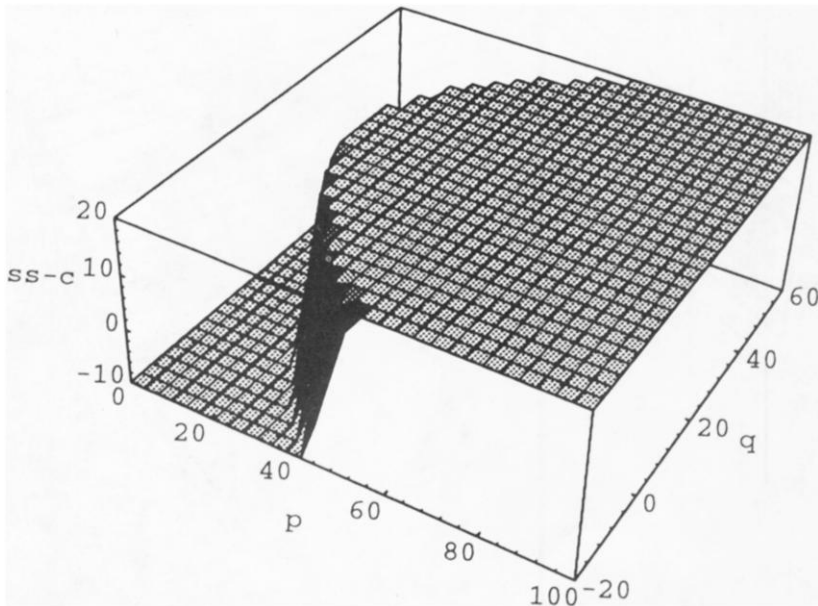


FIGURE 6. Plot of the polynomial for covariate  $s = 7.823426$  truncated at  $ss - c = 20$



In Figure 6 the plotting range on the axis  $ss - c$  was restricted to be between  $-10$  and  $20$  in order to more clearly show the nature of the polynomial in the region of interest, that is, where  $ss - c = 0$ . The upper flat surface of Figure 6 is due to this restriction. The region of significance is identified by the surface produced at the intersection of a plane parallel to the  $p - q$  plane at the height of  $ss - c = 0$ . This region of significance is shown in Figure 7 as a contour plot of Figure 6 with the contour set at a value where  $ss - c = 0$ .

```
Show[ContourGraphics[%], AspectRatio->.8,
      ContourShading->False, Contours->{0},
      FrameLabel->{p,q}]
```

In a similar manner, the region of significance equivalent to a covariate value of  $s = 17.823426$  and plotted in Johnson and Hoyt's (1947) Figure 3 can be found. Figure 8 brings the regions of significance for  $s = 7.823426$  and  $s = 17.823426$  together in one plot (as in Johnson & Hoyt's Figure 1, p. 348).

Because Johnson and Hoyt's (1947) plots are hand drawn and of small scale, it is difficult to compare their accuracy with that of the plots given here. Using Mathematica's capability to read coordinate points from a two-dimensional plot, point sets of  $(p, q)$  falling on each contour having fixed values of  $s$  equal to  $7.823426$  and  $17.823426$  were obtained. Using the summary data provided by

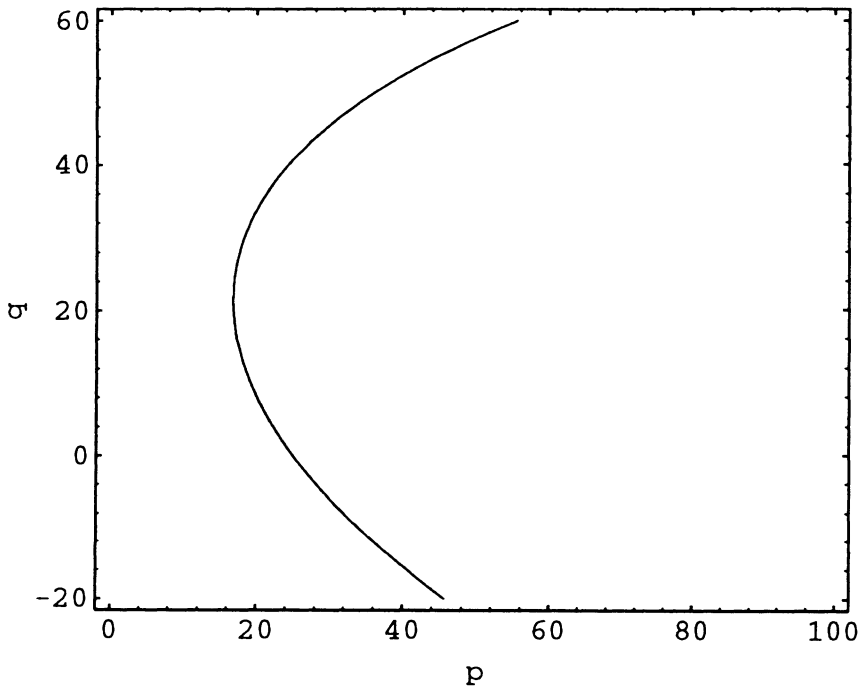


FIGURE 7. Region of significance at covariate  $s = 7.823426$

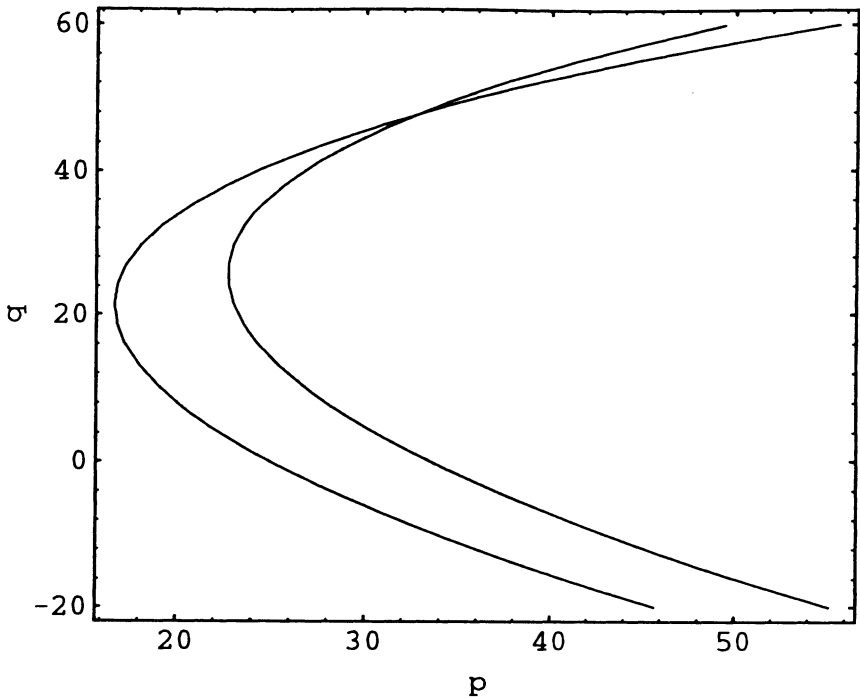


FIGURE 8. Regions of significance at covariate  $s = 7.823426$  and  $17.823426$

Johnson and Hoyt as input to an ANCOVA program, these coordinate values were then used in contrast matrices to test group effects. The probabilities associated with these tests were within the range of 0.0089 to 0.0102, which confirms the correctness of the solution for  $\alpha = .01$ .

**Example 2: Three Groups and Three Covariates**

In this example, the Mathematica statements required to carry out the analysis are not given, because a separate computer program written in Mathematica was used.<sup>1</sup> This program carries out all of the basic calculations illustrated in the last example, except those plots provided by Johnson and Hoyt (1947), and provides results for user-defined contrasts.

Huitema (1980) illustrates the Johnson-Neyman ANCOVA procedure using data for two groups and three covariates with 15 observations in each group. For the purpose of illustration, this data matrix is divided into three groups of 10 observations each. Table 2 gives the means of the variables for each group. A  $\mu$ -model design matrix was used to estimate the parameters  $\mu_1, \beta_{11}, \beta_{12}, \beta_{13}, \mu_2, \beta_{21}, \beta_{22}, \beta_{23}, \mu_3, \beta_{31}, \beta_{32},$  and  $\beta_{33}$ . These OLS estimates are given in Table 3. A test of the sum of squares model corrected for the mean ( $SS_{mcfm}$ ) gives  $SS_{mcfm} = 171.4980, df = 11; MS_e = 0.303821, df = 18,$  and  $F = 51.32$  with  $p \leq .0001$ . A

TABLE 2  
Means for Huitema data

Group	Y	Cov 1	Cov 2	Cov 3
1	10.75	4.10	8.10	3.80
2	9.05	5.90	7.40	4.60
3	11.80	7.00	6.80	4.40

TABLE 3  
OLS parameter estimates of Huitema data

Parameter	Group 1	Group 2	Group 3
$\mu$	8.47578	1.12989	2.82128
$\beta_1$	0.26751	1.02299	1.00184
$\beta_2$	0.11584	0.49388	-0.17383
$\beta_3$	0.06294	-0.38483	0.71542

test of the homogeneity of the regression slopes using  $df = 2$  provides the following: for Covariate 1,  $F = 11.844$ ,  $p \leq .0005$ ; for Covariate 2,  $F = 1.149$ ,  $p \leq .3390$ ; and for Covariate 3,  $F = 4.900$ ,  $p \leq .020$ . Figure 9 shows the three hyperplanes plotted for contour values equal to the means of the dependent variables for the three groups, that is, 10.75, 9.05, and 11.80.

The symbolic contrast matrix to define the polynomial relating  $SS_{mcfm}$  to the unknown values  $p$ ,  $q$ , and  $s$  of the covariates to which an adjustment is sought is

$$\mathbf{K}'_1 = \begin{bmatrix} 1 & p & q & s & -1 & -p & -q & -s & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & p & q & s & -1 & -p & -q & -s \end{bmatrix}.$$

Calculating the polynomial given by the left-hand part of (2), Mathematica provides the following symbolic result:

$$\begin{aligned} & (0.698468 (19461.2 - 5843.39 p + 637.897 p^2 - 32.161 p^3 \\ & + 0.816059 p^4 - 3478.45 q + 925.765 p q - 83.8804 p^2 q \\ & + 2.04099 p^3 q + 165.35 q^2 - 22.1883 p q^2 + 2.33551 p^2 \\ & q^2 - 8.51403 q^3 - 1.20645 p q^3 + 1. q^4 - 2609.86 s + \\ & 481.937 p s - 2.27356 p^2 s - 1.3041 p^3 s + 257.735 q s \\ & - 84.5905 p q s + 1.92134 p^2 q s + 19.6921 q^2 s + \\ & 3.27911 p q^2 s - 3.52683 q^3 s + 336.326 s^2 - 40.3918 p \\ & s^2 + 2.32929 p^2 s^2 - 30.6719 q s^2 + 1.47897 p q s^2 + \\ & 4.3459 q^2 s^2 - 13.9853 s^3 - 0.652421 p s^3 - 3.47824 q \\ & s^3 + 2.95241 s^4)) / \\ & (8465.94 - 858.975 p + 77.2914 p^2 - 3.21349 p^3 + \\ & 0.0791991 p^4 - 3312.59 q + 243.235 p q - 14.8958 p^2 q \\ & + 0.336888 p^3 q + 506.148 q^2 - 23.7958 p q^2 + 0.764708 \\ & p^2 q^2 - 35.826 q^3 + 0.813569 p q^3 + 1. q^4 - 1051.35 s - \\ & 33.8838 p s + 3.22571 p^2 s - 0.268882 p^3 s + 255.995 q \\ & s + 9.97869 p q s - 0.592037 p^2 q s - 20.2637 q^2 s - \end{aligned}$$

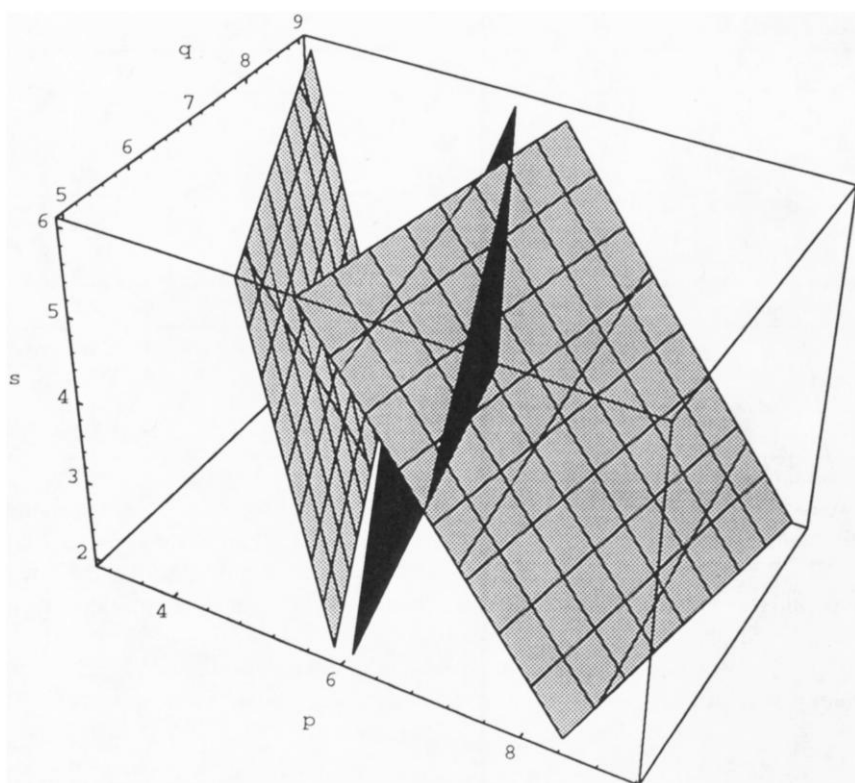


FIGURE 9. Contour plot of regression hyperplanes at  $Y = (10.75, 9.05, 11.80)$

$$0.777418 p q^2 s + 0.431562 q^3 s + 182.69 s^2 - 5.88228 p s^2 + 0.707141 p^2 s^2 - 34.3571 q s^2 + 1.09478 p q s^2 + 1.92591 q^2 s^2 - 6.04068 s^3 - 0.850662 p s^3 - 0.0481645 q s^3 + 0.692588 s^4)$$

This polynomial is a quartic having 35 terms in the numerator and 35 terms in the denominator. The right-hand part of (2) is 3.65369, based on  $df_k = 2$ ,  $df_e = 18$ ,  $\alpha = .01$ , and  $MS_e = .303821$ . The next step is to plot the polynomial for those values of  $p$ ,  $q$ , and  $s$  that provide a result equal to 3.65369, because then the expression

$$(\mathbf{K}'\mathbf{B})[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}\mathbf{K}'\mathbf{B} = 3.65369$$

is true and represents the point at which significance is just reached. Figure 10 shows the region of significance. In this plot, any covariate point  $(p, q, s)$  that falls on the surface is significant at  $\alpha = .01$ . Plotting for a value of  $\alpha = .001$  (Figure 11) shows that any covariate point falling on the surface or within the

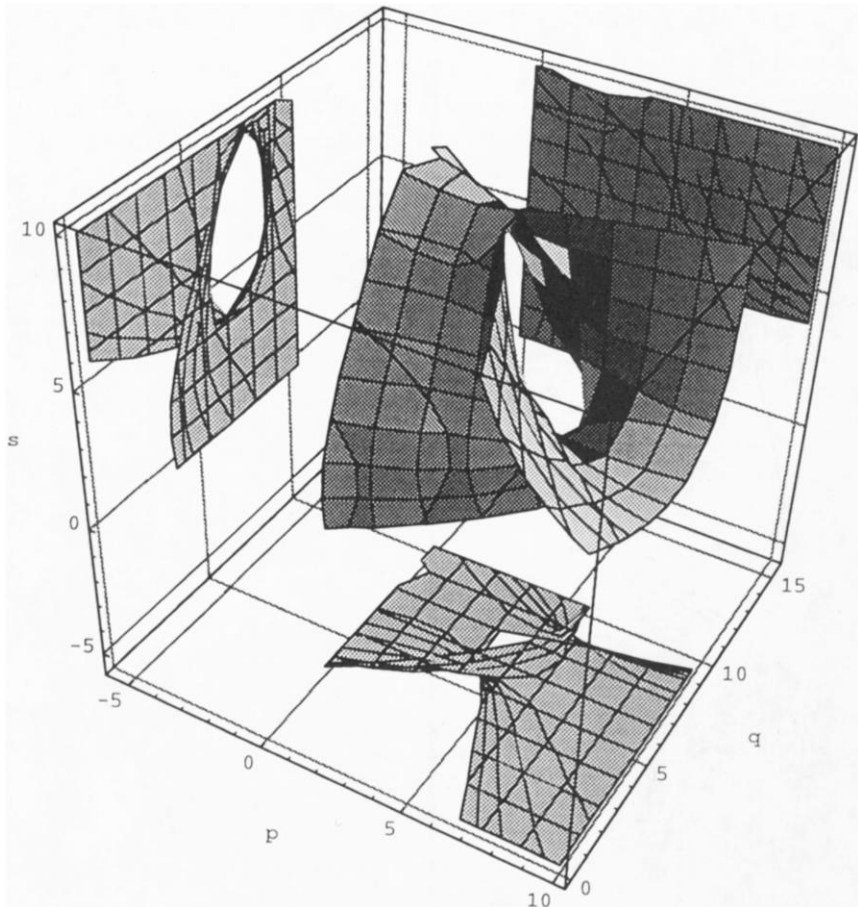


FIGURE 10. *Region of significance for  $SS_{mcfm}$  for  $\alpha = 0.01$*

enclosed volume represented in Figure 10, or within the outer boundaries of the shadow projections, will be significant at  $\alpha = .01$  or less.

By examining Figures 10 and 11 it can be seen that the basic form of the region of significance involves two paraboloids. Those in Figure 11 (for  $\alpha = .001$ ) are nested within the volume enclosed by those in Figure 10 (for  $\alpha = .01$ ).

#### *Pairwise Contrasts*

Defining regions of significance for contrasts can be handled in the same manner as is the region of significance for  $SS_{mcfm}$ . That is, a polynomial is defined for each row of the contrast matrix and then contour plotted for a specific value of  $\alpha$ . Since the Helmert contrast matrix given below, which is orthogonal by row, can be derived by elementary row operations from the contrast matrix  $\mathbf{K}_1$  used earlier, it will produce the same sum of squares.

$$\mathbf{K}'_2 = \begin{bmatrix} 1 & p & q & s & -1/2 & -p/2 & -q/2 & -s/2 & -1/2 & -p/2 & -q/2 & -s/2 \\ 0 & 0 & 0 & 0 & 1 & p & q & s & -1 & -p & -q & -s \end{bmatrix}$$

Contour plots of the two polynomials, produced by using Rows 1 and 2, respectively, of  $\mathbf{K}_2$  to define the regions of significance at  $\alpha = .01$ , are given in Figures 12 and 13.

Although the sums of squares provided by the range of values used for  $p$ ,  $q$ , and  $s$  are not expected to be orthogonal, if one envisages combining the contour plots of Figures 12 and 13 additively for the range of  $p$  from  $-5$  to  $10$ , the result approximates the contour plot for  $SS_{mcfm}$  given in Figure 10.

Searle (1971, p. 200) shows that contrasts  $\mathbf{K}_j$  and  $\mathbf{K}_k$  of one degree of freedom each are orthogonal contrasts if

$$\mathbf{K}'_j(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}_k = 0. \tag{3}$$

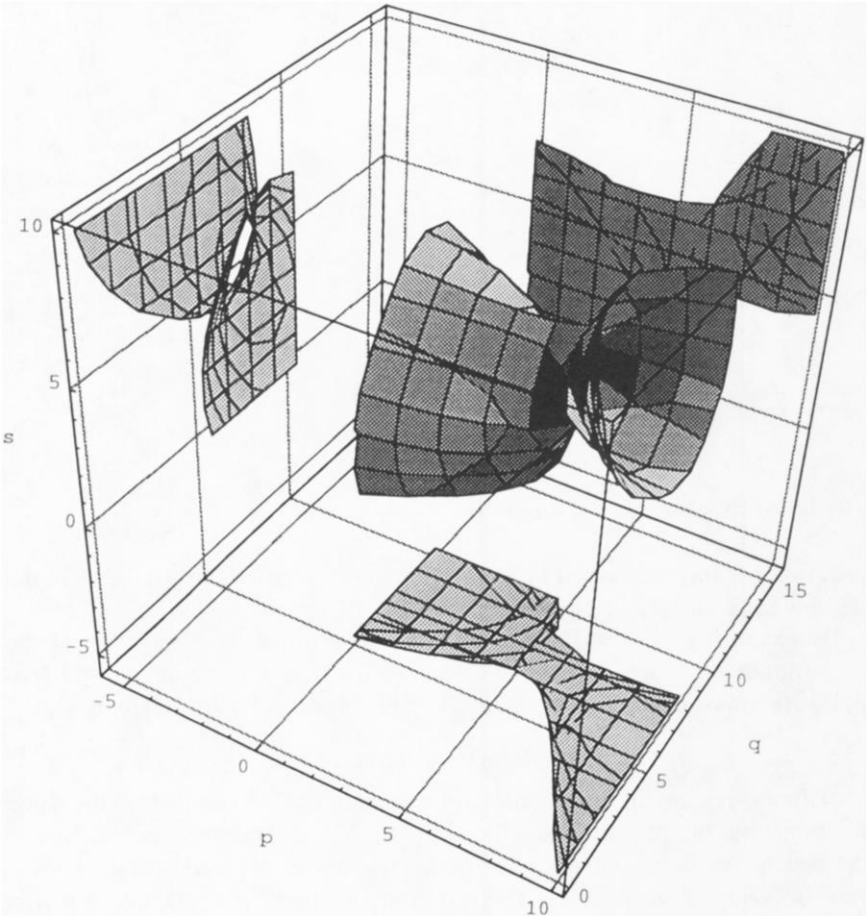


FIGURE 11. Region of significance for  $SS_{mcfm}$  for  $\alpha = 0.001$



Using this criterion, setting the  $j$ th contrast to be Row 1 of the Helmert contrast and the  $k$ th contrast to be Row 2, and using Mathematica to expand the result in rational form produces the following polynomial:

$$(-515638292 + 40451866 p - 152241 p^2 + 117622882 q - 5500481 p q - 6719617 q^2 - 20114016 s - 153500 p s + 3582710 q s - 126540 s^2) / 25233988$$

A numerical evaluation of this polynomial over the range of 0 to 10 for the covariate values  $p$ ,  $q$ , and  $s$  provides the plot given in Figure 14. The order of evaluating this polynomial is as follows:  $s$  varies the fastest,  $q$  next, and then  $p$ . The horizontal axis represents the number of evaluations made, that is,  $11 \times 11 \times 11 = 1,331$  iterations. An ANCOVA does not generally produce additive sums of squares, even when the rows of the contrast matrix are orthogonal; however, the plot in Figure 14 shows that the polynomial evaluates to zero at many points, and therefore indicates that there are specific values of the covariates which will produce orthogonal contrasts and sums of squares. A plot of this polynomial for a contour equal to 0 in order to satisfy (3) is given in Figure 15. Any set of points  $(p, q, s)$  that falls on the surface will produce orthogonal contrasts by Searle's definition. Using the same procedure as in Example 1, we substitute the values  $p = 5$  and  $q = 10$  into the polynomial in order to reduce it to a quadratic, and then solve for  $s$ . The result is  $s = 6.21073$  and  $111.899$ . Selecting the former

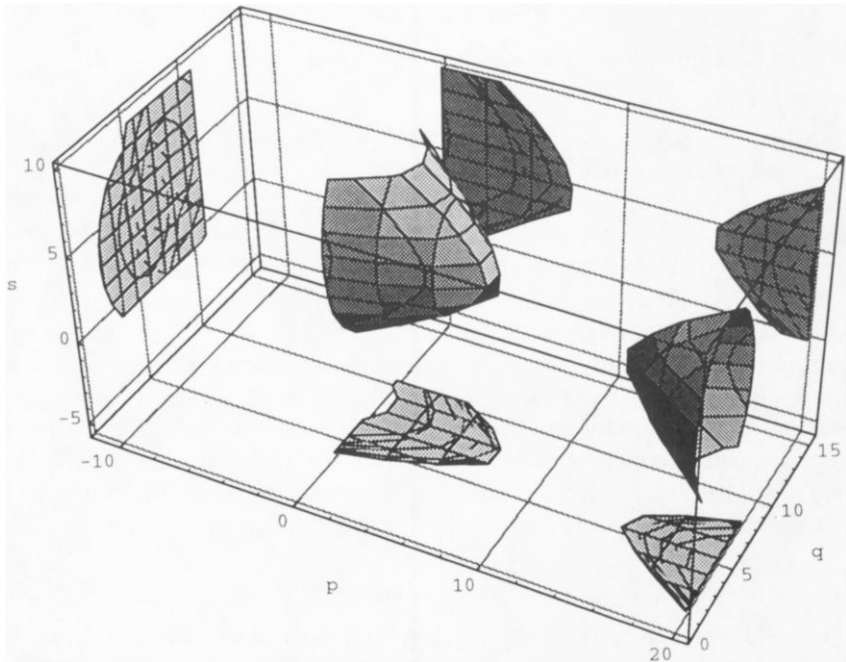


FIGURE 12. Region of significance for Row 1 of Helmert contrast for  $\alpha = 0.01$



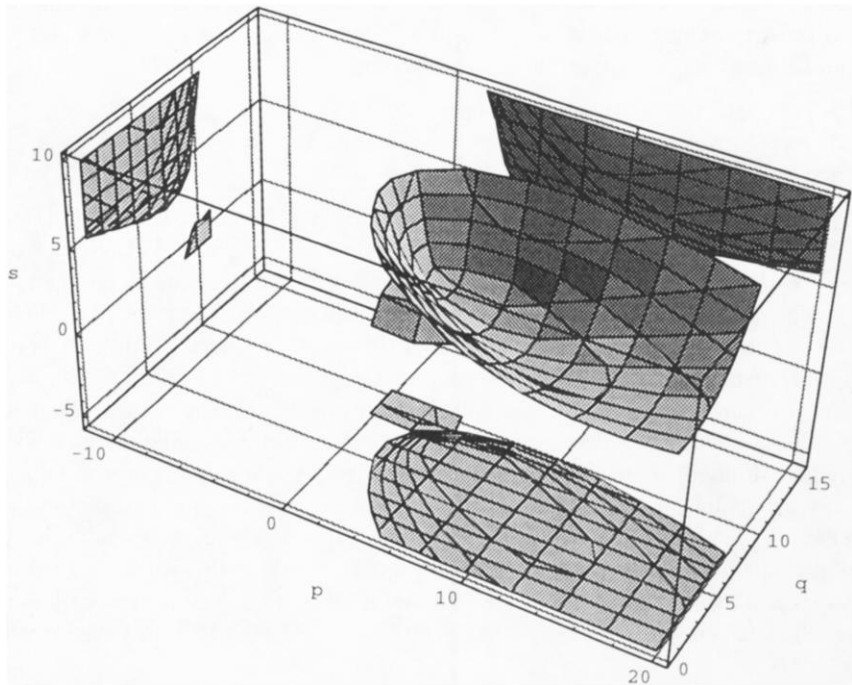


FIGURE 13. *Region of significance for Row 2 of Helmert contrast for  $\alpha = 0.01$*

value for  $s$  (because it is within the range of the data), using these values of  $p$ ,  $q$ , and  $s$  in the Helmert contrast matrix, and evaluating (3) yields  $8.45975 \times 10^{-7}$ , which indicates that orthogonality exists. Checking further, we find that the full Helmert contrast matrix produces a sum of squares of 1.18614, which is equal to the sum of squares produced by Row 1 (0.775682) and Row 2 (0.410463).

Although Llabre and Ware (1980) examined nonorthogonality in the ANCOVA, their study essentially examined the fitting order of the parameters and not the orthogonality of contrasts. Differences in sums of squares due to fitting order in the ANCOVA are best explained using Searle's (1971)  $R(.l.)$  notation. The criterion for orthogonality of contrasts, as well as the specificity of covariate values forming orthogonal contrasts found here, suggests that examining the orthogonality of contrasts by generating data sets randomly, as was done by Llabre and Ware, would not be very fruitful.

### Discussion

Using the general linear model together with the symbolic contrast matrix to generate a polynomial expression relating sums of squares to values of the covariates to which an adjustment is made can be handled in a simple and

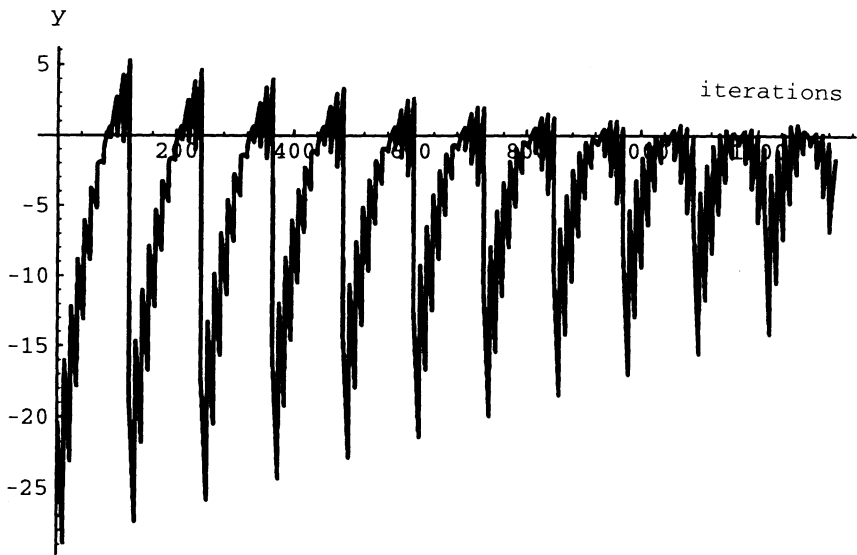


FIGURE 14. *Evaluation of the polynomial defining orthogonal contrasts for covariate values  $p$ ,  $q$ , and  $s$  ranging from 0 to 10*

straightforward manner using the numerical and symbolic processing capabilities of Mathematica. Point sets defining the Johnson-Neyman regions of significance for three covariates can be represented by plotting the polynomial expression in a three-dimensional space for a specific value at which significance occurs, or by reducing the plot to a two-dimensional space by fixing a covariate value. If a subset of orthogonal contrast exists among a set of contrasts, these can be identified by solving the polynomial produced by (3) either graphically or by setting the coefficients for a couple of the covariates and explicitly solving for the third. The latter approach is suggested, because reading points of three variables from a three-dimensional surface plot is difficult. Although common statistical programs may have adequate numerical and plotting capabilities, they do not have the capability to manipulate symbolic expressions. Therefore, it is unlikely that packages such as SPSS, BMDP, and SAS can be expanded to accommodate analysis with the same degree of generality as illustrated here.

Researchers contemplating the use of Mathematica need to be aware that symbolic computations and those required to produce high-resolution color plots can be time consuming and use a considerable amount of a computer's main memory. User options allow the allocation of main memory to the front-end processor of Mathematica and to the Mathematica kernel. The front-end processor looks after input to and output from the kernel, including the rendering of plots, while the kernel carries out and retains the results of all the computations.

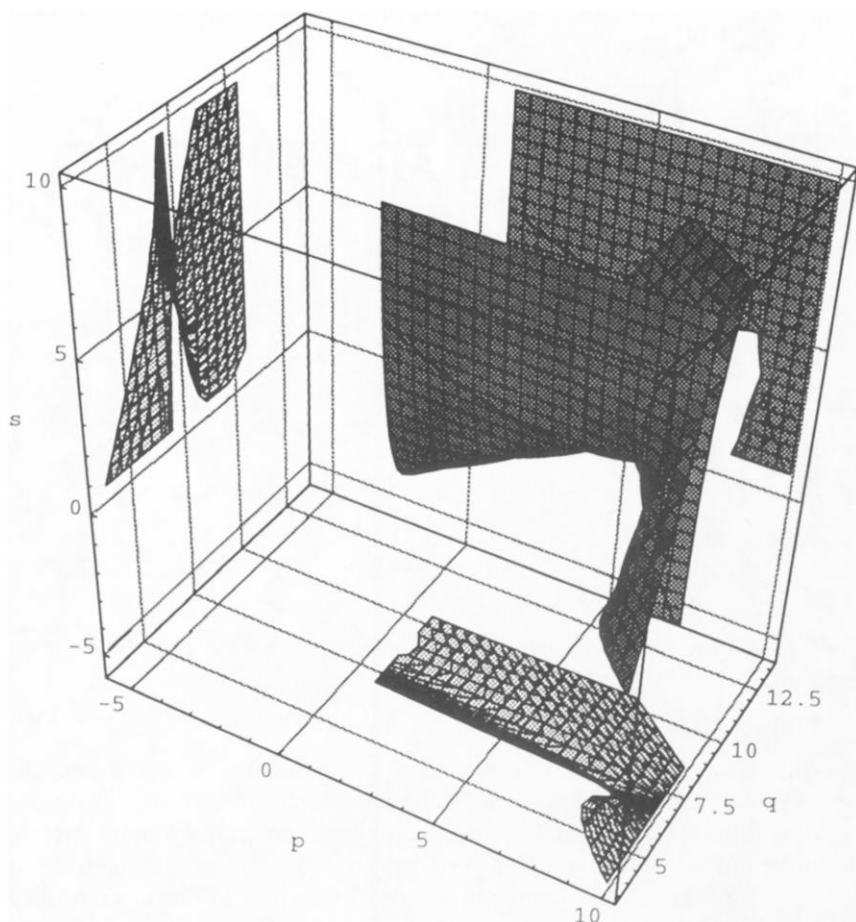


FIGURE 15. Contour plot of the polynomial defining orthogonal contrasts for values of  $(p, q, s)$  resulting in 0

All computations illustrated here were carried out both on a Macintosh IIsi with 9MB of memory and on a Power Macintosh 6100/60 with 24MB of memory. The interactive computations of Example 1 were carried out cumulatively (i.e., plots were displayed and all computations saved) on the Macintosh IIsi with 2.5 MB allocated to the front end and 5.5MB allocated to the kernel. The program used in Example 2 ran to completion with 1.5MB allocated to the front end and 6MB to the kernel. Various features of Mathematica can be used to accommodate smaller memories—for example, initially obtaining lower-resolution plots and deleting intermediate results. Computational time is not as easily controlled by the user, except where the efficiency of computation can be optimized by

programming techniques and by modifying the parameters provided to standard Mathematica operations. For example, a higher-resolution plot of Figure 15, the lower-resolution version of which took 30 seconds to compute, would have taken approximately 8 minutes (i.e., about 16 times longer) if the recursion parameter had been 2 instead of 1. A small change in the granularity of the search for the zero condition of the polynomial would have raised the computation time to 2 hours. Likewise, if Example 2 had involved five groups rather than three, the polynomial produced would have contained 165 terms in the numerator and denominator (rather than 35 terms) and would have taken 1 hour of processing time rather than just 2 minutes. These constraints notwithstanding, Mathematica is a very powerful tool that a researcher or a student, using either a Macintosh or a PC, can use to explore methods of data analysis. The discussions of the Mathematica newsgroup on the Internet suggest that most researchers use PCs for problems in physics, engineering, and mathematics.

### Note

<sup>1</sup> This program, which carries out the analysis to define regions of significance as illustrated in the last example for use on a PC or a Macintosh platform, is available at <http://www.wolfram.com/cgi-bin/MathSource/Enhancements/Statistics/0208-066>.

### References

- Abell, M. L., & Braselton, J. P. (1994). *The Maples V handbook*. Boston: Academic Press.
- Borich, G. D., Godbout, R. C., & Wunderlich, K. W. (1976). *The analysis of aptitude-treatment interaction: Computer programs and calculations*. Chicago: International Educational Services.
- Butsch, R. L. C. (1944). A worksheet for the Johnson-Neyman technique. *Journal of Experimental Education*, 12, 226–241.
- Carroll, J. B., & Wilson, G. F. (1969). *An interactive-computer program for the Johnson-Neyman technique in the case of two groups, two predictor variables, and one criterion variable* (Research Bulletin RB69–68). Princeton, NJ: Educational Testing Service.
- Carroll, J. B., & Wilson, G. F. (1970). An interactive-computer program for the Johnson-Neyman technique in the case of two groups, two predictor variables, and one criterion variable. *Educational and Psychological Measurement*, 30, 121–132.
- Ceurvorst, R. W. (1979). Computer program for the Johnson-Neyman technique. *Educational and Psychological Measurement*, 39, 205–207.
- Fisher, R. (1932). *Statistical methods for research workers*. Edinburgh, Scotland: Oliver and Boyd.
- Fisher, R. (1935). *The design of experiments*. Edinburgh, Scotland: Oliver and Boyd.
- Huitema, B. E. (1980). *The analysis of covariance and alternatives*. New York: Wiley.
- Hunka, S. (1994). Using Mathematica to solve Johnson-Neyman problems. *Mathematica in Education*, 3(3), 32–36.
- Hunka, S. (1995). Identifying regions of significance in ANCOVA problems having non-homogeneous regressions. *British Journal of Mathematical and Statistical Psychology*, 48, 161–188.

- Johnson, P. O., & Fay, L. C. (1950). The Johnson-Neyman technique, its theory and application. *Psychometrika*, 15, 349–367.
- Johnson, P. O., & Hoyt, C. (1947). On determining three dimensional regions of significance. *Journal of Experimental Education*, 15, 342–353.
- Johnson, P. O., & Jackson, R. W. B. (1959). *Modern statistical methods: Descriptive and inductive*. Chicago: Rand McNally.
- Johnson, P. O., & Neyman, J. (1936). Tests of certain linear hypotheses and their application to some educational problems. *Statistical Research Memoirs*, 1, 57–93.
- Karpman, M. B. (1980). ANCOVA—A one covariate Johnson-Neyman algorithm. *Educational and Psychological Measurement*, 40, 791–793.
- Karpman, M. B. (1983). The Johnson-Neyman technique using SPSS or BMDP. *Educational and Psychological Measurement*, 43, 137–147.
- Karpman, M. B. (1986). Comparing two non-parallel regression lines with the parametric alternative to analysis of covariance using SPSS-X or SAS—The Johnson-Neyman technique. *Educational and Psychological Measurement*, 46, 639–644.
- Kirk, R. E. (1982). *Experimental design*. Pacific Grove, CA: Brooks/Cole.
- Koenker, R. H., & Hansen, C. W. (1942). Steps for the application of the Johnson-Neyman technique—A sample analysis. *Journal of Experimental Education*, 10, 164–173.
- Kush, J. C. (1986). A Fortran V IBM computer program for the Johnson-Neyman technique. *Educational and Psychological Measurement*, 46, 185–187.
- Lautenschlager, G. J. (1987). JOHN-NEY: An interactive program for computing the Johnson-Neyman confidence region for nonsignificant prediction differences. *Applied Psychological Measurement*, 11, 194–195.
- Leighton, J. (1995). *Circumventing ANCOVA's assumptions using Mathematica*. Unpublished master's thesis, University of Alberta, Edmonton, Canada.
- Llabre, M. M., & Ware, W. B. (1980). Equal cell size and nonorthogonality in ANCOVA. *Educational and Psychological Measurement*, 40, 91–94.
- Maxwell, S. E., & Delaney, H. D. (1990). *Designing experiments and analyzing data: A model comparison perspective*. Belmont, CA: Wadsworth.
- Maxwell, S. E., O'Callaghan, M. F., & Delaney, H. D. (1993). Analysis of covariance. In L. K. Edwards (Ed.), *Applied analysis of variance in behavioral science* (pp. 63–104). New York: Marcel Dekker.
- Mood, A. M. (1950). *Introduction to the theory of statistics*. New York: McGraw-Hill.
- Pedhazur, E. J. (1982). *Multiple regression in behavioral research*. New York: Holt, Rinehart and Winston.
- Rogosa, D. (1981). On the relationship between the Johnson-Neyman region of significance and statistical tests of parallel within-group regressions. *Educational and Psychological Measurement*, 41, 73–84.
- Rutherford, A. (1992). Alternatives to traditional analysis of covariance. *British Journal of Mathematical and Statistical Psychology*, 45, 197–223.
- Schafer, W. D., & Wang, Y. (1991, April). *Graphical description of Johnson-Neyman outcomes for linear and quadratic regression surfaces*. Paper presented at the Annual Meeting of the American Educational Research Association, Chicago.
- Searle, S. R. (1971). *Linear models*. New York: Wiley.
- Searle, S. R. (1987). *Linear models for unbalanced data*. New York: Wiley.

- Searle, S. R., Speed, F. M., & Henderson, H. V. (1981). Some computational and model equivalences in analysis of variance of unequal-subclass-numbers data. *The American Statistician*, 35(1), 16–33.
- Walker, H. M., & Lev, J. (1953). *Statistical inference*. New York: Henry Holt.
- Wolfram, S. (1991). *Mathematica, a system for doing mathematics by computer*. Menlo Park, CA: Addison-Wesley.

### Authors

STEVE HUNKA is Professor Emeritus in the Department of Educational Psychology and the Centre for Research in Applied Measurement and Evaluation, University of Alberta, Edmonton, AB, Canada T6J 2G5; [steve.hunka@ualberta.ca](mailto:steve.hunka@ualberta.ca). He specializes in applied statistics and related computational methods.

JACQUELINE LEIGHTON is a graduate student in human development and measurement, Department of Psychology, University of Alberta, Edmonton, AB, Canada T6G 2E9; [leighton@psych.ualberta.ca](mailto:leighton@psych.ualberta.ca). She specializes in probabilistic reasoning.

Received November 13, 1995

Revision received June 10, 1996

Accepted July 8, 1996