

BK-674-M  
Rev. May 1980

## Intersite Transfer of Estimated Response Surfaces

CONSTANCE L. WOOD

Department of Statistics, University of Kentucky,  
Lexington, Kentucky 40506, U. S. A.

FOSTER B. CADY

Biometrics Unit, Cornell University,  
Ithaca, New York 14853, U. S. A.

### Summary

Transferability of agrotechnology assumes the feasibility of extrapolating a response-input relationship, estimated from experimental sites, to other sites with similar conditions. One specific conjecture is that crop production technology is transferable across sites within a soil family classification. The general approach to evaluating the transfer conjecture incorporates into the data analysis the prediction of yields not used in the estimation of the transfer function. A transfer model, using a second order response surface and measured site variable information, is formulated and yields for each experimental site are predicted from a transfer function estimated from the other sites. The resulting transfer residuals are compared with the ordinary within-site residuals. Based on a sum of squares criterion, a prediction test statistic is developed and shown to have a distribution of a ratio of independent quadratic forms. The methodology of transfer residuals is applied to data from the Benchmark Soil Project, where the major objective is to assess the feasibility of transferring agrotechnology among sites having soil of the same taxonomic classification.

---

Key Words: Regression; Prediction; Extrapolation; Controlled and uncontrolled variables; Agrotechnology transfer.

## 1. Introduction

Agrotechnology transfer is the extrapolation of a response-input relationship, estimated from a series of experiments, to new but similar sites. A major objective of the Benchmark Soils Project, established by U.S.A.I.D. (Agency for International Development) in cooperation with the Universities of Hawaii and Puerto Rico, is to assess the feasibility of crop production technology transfer from one tropical site to another on the basis of similarity of soils as indicated by the soil family in the Soil Taxonomy Classification System (Soil Survey Staff 1975). The conjecture is that experimental results, specifically the response surface relating maize yield to applications of phosphorus and nitrogen, obtained from a set of sites can be applied to new sites on the same soil family.

The soil family was selected because the family classification integrates soil factors with long-term environmental factors that influence crop yield. However, because of natural and past management variability, soil properties are not constant within discrete soil families. Consequently, homogeneous response to applied fertilizer treatments usually will not be found in practice. Interpreting this to mean that agrotechnology is not transferable can be faulty. In particular, the individual site response surface may be affected by the specific biotic environment of the site. Only by measuring uncontrolled site variables which reflect differences in environments, and including them in the response surface, will the response to the applied variables be clearly focused.

Statistical analysis of the transfer conjecture involves

(i) identification and estimation of a response surface model which adequately relates maize yields from several experimental sites to both the applied fertilizer levels and measured site variables,

(ii) evaluation of the predictive ability of the resulting estimated response surface for sites within the same soil family, but not included in the estimation process.

Least squares estimation procedures for incorporating site variable information in the analysis of orthogonal response surface models have been discussed generally by Cochran and Cox (1957, Chapter 14), and specifically by Colwell (1967), and are used throughout. Here the focus is on quantitative evaluation of the predictive ability of the resulting response surface.

In order to assess the predictive ability, the actual transfer of agrotechnology to sites where experimentation has not been carried out needs to be simulated. Our approach is to predict the yields for one of the  $k$  experimental sites, say the  $i$ th site, using the response surface estimated from the other  $(k-1)$  sites. Since the  $i$ th site was excluded in the estimation procedure, the resulting  $(n \times 1)$  column vector of predicted yields for the  $i$ th site will be denoted by  $\hat{Y}_{[-i]}$ . This prediction procedure is then repeated for each of the  $k$  sites, i.e., the predicted yields at a particular site are based on a response surface estimated from the other  $(k-1)$  sites. The resulting prediction error is reflected in the  $k$  vectors of transfer residuals  $Y_i - \hat{Y}_{[-i]}$ , where  $Y_i$  is the  $(n \times 1)$  vector of observed yields for the  $i$ th site.

A quantitative evaluation of the predictive ability of the estimated response surface can then be based on a comparison of the transfer residuals  $Y_i - \hat{Y}_{[-i]}$ , with the least squares within-site residuals,  $Y_i - \hat{Y}_i$ , calculated by fitting the  $p$  variate response surface individually to each of the  $k$  sites. The specific objectives here are (i) to develop the transfer residual methodology for evaluating prediction and (ii) to demonstrate the methodology with yield and site variable information from maize transfer experiments on

the thixotropic, isothermic soil family of Hydric Dystrandeps. The first step is the development of a statistic for evaluating the transfer residuals.

## 2. Prediction Using Site Variables

Our approach utilizes a sum of squares criterion to compare the magnitude of the transfer residuals,  $Y_{\sim i} - \hat{Y}_{\sim i[-i]}$ , to the ordinary within-site residuals,  $Y_{\sim i} - \hat{Y}_{\sim i}$ . In particular, Cady (1974, Experimental strategy for transferring crop production information, Technical Report 502 in the Biometrics Unit Mimeo Series, Cornell University, Ithaca, New York) proposed the ratio of the pooled sum of squared transfer residuals to the pooled sum of squared within-site residuals; i.e.,

$$P = \frac{\sum_{i=1}^k (Y_{\sim i} - \hat{Y}_{\sim i[-i]})' (Y_{\sim i} - \hat{Y}_{\sim i[-i]})}{\sum_{i=1}^k (Y_{\sim i} - \hat{Y}_{\sim i})' (Y_{\sim i} - \hat{Y}_{\sim i})} .$$

For two sites (P-1) is a symmetrized version of Gardner's (1972) ratio bias statistic used for assessing the predictive ability of one sample for a second sample. In the more general case of k sites, P is the natural extension comparing the predictive ability of the ith site for itself with the predictive ability of the remaining (k-1) sites.

The distribution of (P-1) is considered for a prediction model describing an experimental situation for which

(a) the same equally replicated treatment design was used at each site generating independent, normally distributed yields with common unknown experimental error variance,  $\sigma^2$ ,

(b) the same p variate response surface model; e.g., a quadratic polynomial in two treatment factors (p=5), adequately fits each site,

(c) differences among the site means, which would not affect the economically optimal rates of the treatment factors, have been eliminated by subtracting the site mean from yield data within each site, and

(d) differences in the estimated model parameters can be explained by interactions between the  $p$   $x$ -variables of the prediction equation and the measured site variables. All  $x$ -variables used in interactions are required to be orthogonal and centered at zero.

For this prediction model, the statistic  $(P - 1)$  can be written as the ratio of two quadratic forms; i.e.,

$$(P - 1) = \frac{k^2}{(k-1)^2} \cdot \frac{\underline{\underline{Y}}' \underline{\underline{B}}_1 \underline{\underline{Y}}}{\underline{\underline{Y}}' \underline{\underline{B}} \underline{\underline{Y}}},$$

where  $\underline{\underline{B}}_1$  and  $\underline{\underline{B}}$  are  $(kn \times kn)$  symmetric matrices with  $\underline{\underline{B}}_1 \underline{\underline{B}} = 0$ ,  $\underline{\underline{B}} \underline{\underline{B}} = \underline{\underline{B}}$  and  $\underline{\underline{Y}}' = [\underline{\underline{Y}}'_1 : \underline{\underline{Y}}'_2 : \dots : \underline{\underline{Y}}'_k]$ . (See Appendix for details.) Since  $\underline{\underline{B}}_1 \underline{\underline{B}} = 0$ , the numerator and denominator are independent. Also,  $\underline{\underline{B}} \underline{\underline{B}} = \underline{\underline{B}}$  implies that the denominator is distributed as  $\sigma^2 \chi^2 [k(n-p-1) \text{d.f.}]$ . The distribution of the numerator, unfortunately, cannot be so easily identified. However, the quadratic form,  $\underline{\underline{Y}}' \underline{\underline{B}}_1 \underline{\underline{Y}}$ , can be represented as  $\sigma^2 \sum_{\ell=1}^{kn} \theta_{\ell} \chi^2_{\ell}$  (1 d.f.), where  $\chi^2_{\ell}$  (1 d.f.) are independent  $\chi^2$  random variables, each with 1 degree of freedom (d.f.), and  $\theta_{\ell}$  are the eigenvalues of  $\underline{\underline{B}}_1$ .

In general, the distribution of a linear combination of  $\chi^2$  variables cannot be simplified. However, in a particular problem, once  $P$  has been computed from the data and the eigenvalues of  $\underline{\underline{B}}_1$  have been determined, the attained significance level of  $P$  can be accurately estimated by Monte Carlo methods.

One complication arises because the dimension of  $\underline{\underline{B}}_1$  is extremely large for moderate  $k$  and  $n$ . This makes direct numerical computation of the eigenvalues not feasible. Due to the orthogonality and centering at zero of the  $x$ -variables used in interactions, this problem can be reduced to the computation of the eigenvalues of several matrices, each of dimension  $k$ , a problem

readily handled by available computer packages.

Specifically, if two orthogonal, centered, x-variables are used to form interactions with measured site variables, the eigenvalues of  $\underline{B}_1$  can be thought of as arising in three groups. The eigenvalues in the first are either zero or one, the number depending only on  $(p - 2)$  and the number of sites. Each of the remaining groups corresponds to one of the x-variables and the eigenvalues depend only on the site variables used as interactions with that x-variable. Two matrices,  $\underline{C}_1$  and  $\underline{C}_2$ , with the required eigenvalues, are computed as follows:

(i) Construct  $\underline{T}_{(-i)1}$  as a  $(k - 1) \times m_1$  matrix of  $m_1$  site variables used to form interactions with the first x-variable for all sites except the ith site. Each column of  $\underline{T}_{(-i)1}$  is centered at zero. The  $k$  excluded  $(1 \times m_1)$  row vectors are denoted as  $\underline{T}_{i1}$ . Similarly, construct  $\underline{T}_{(-i)2}$  and  $\underline{T}_{i2}$  for the second x-variable.

(ii) Form  $\underline{T}_1$  as a  $(k \times k)$  matrix with diagonal elements equal to zero and the remaining elements in the ith row are given in order by the elements of the alias matrix,  $\underline{T}_{i1} \left( \underline{T}'_{(-i)1} \underline{T}_{(-i)1} \right)^{-1} \underline{T}'_{(-i)1}$ . Similarly form  $\underline{T}_2$ .

(iii) Calculate  $\underline{C}_1$  as a  $(k \times k)$  matrix by

$$\underline{C}_1 = \left( \underline{I}_k - k^{-1} \underline{J}_k \right) + k^{-2} (k-1)^2 \left[ \underline{T}'_1 \underline{T}_1 - k(k-1)^{-1} (\underline{T}_1 + \underline{T}'_1) + (k-1)^{-1} (\underline{J}_k \underline{T}_1 + \underline{T}'_1 \underline{J}_k) \right]$$

where  $\underline{I}_k$  is a  $(k \times k)$  identity matrix and  $\underline{J}_k$  is a  $(k \times k)$  matrix of ones.

Similarly, form  $\underline{C}_2$  from  $\underline{T}_2$ .

If  $\theta_1, \dots, \theta_k$  denote the eigenvalues of  $\underline{C}_1$  and  $\theta_{k+1}, \dots, \theta_{2k}$  denote the eigenvalues of  $\underline{C}_2$ , then the remaining eigenvalues of  $\underline{B}_1$  will either be zero or one. In particular, only  $(k - 1)(p - 2)$  of the remaining eigenvalues will be equal to one; the rest will be zero. (See Appendix, equation (2).) The

term  $(p - 2)$  is the number of x-variables which are not used to form interactions with site variables.

Note that if site variables are not included in the response surface model and the above approach is followed using only x-variables, then  $(P - 1)$  is proportional to the usual F-statistic for testing equality of the k response surfaces.

### 3. Example

The Benchmark Soils Project is described by Silva and Beinroth (1978, Research on agrotechnology transfer in the tropics based on the soil family (Progress Report 1, Benchmark Soils Project), Technical Report, Department of Agronomy and Soil Science, University of Hawaii, Honolulu, Hawaii). As indicated in Section 1, a major objective of the project is to assess the feasibility of transferring agrotechnology in the tropics on the basis of soil taxonomic units, thereby reducing the amount of site specific experimentation. Specifically, the conjecture that an estimated response-input relationship can be transferred within the same soil family needs to be evaluated. This example uses grain yield (kg/ha) data from five maize experiments on the Hydric Dystrandept soil family; two sites (PUC-K and BUR-B) are in the Philippines, two in Hawaii (KUC-C and KUK-D) and one in Indonesia (LPH-E). The same 13-point treatment design with three replications was used at each site, a partial  $5 \times 5$  factorial with applied phosphorus and applied nitrogen as the controlled variables. An estimated second order response surface model in the two factors adequately fits the treatment means.

Given here are the numerical details of calculating the P statistic under a prediction model which introduces site variable information in the transfer function as interactions between the site variables and the linear effects of applied phosphorus (P) and applied nitrogen (N). Table 1 gives



TABLE 1  
Site Variable Data, Residual Sums of Squares  
and Transfer Sums of Squares

Site	EXTN (ppm)	MINT (°C)	EXTP (ppm)	Residual SS (x10 <sup>3</sup> )	Transfer SS (x10 <sup>3</sup> )
PUC-K	79	23.00	10	5,869	14,700
BUR-B	29	21.50	5	25,055	36,584
KUK-C	46	18.83	74	13,602	18,695
KUK-D	29	17.90	62	25,599	32,792
LPH-E	119	16.76	23	<u>17,880</u>	<u>23,660</u>
				88,005	126,431

Next, since interactions with the first two columns of  $X$ , say  $x_{\sim 1}$  and  $x_{\sim 2}$ , respectively, are used in the prediction equation, the data matrices for the first site, PUK-K, for calculating  $\hat{Y}_{\sim}[-1]$  are

$$X_{\sim 1} = [X : \begin{matrix} P \times \text{EXTP} \\ (10)_{x_{\sim 1}} \end{matrix} : \begin{matrix} P \times \text{EXTN} \\ (79)_{x_{\sim 1}} \end{matrix} : \begin{matrix} N \times \text{MINT} \\ (23)_{x_{\sim 2}} \end{matrix} : \begin{matrix} N \times \text{EXTN} \\ (79)_{x_{\sim 2}} \end{matrix}]$$

and

$$X_{\sim}(-1) = \begin{bmatrix} \begin{matrix} P \times \text{EXTP} \\ X \\ \sim \end{matrix} & \begin{matrix} P \times \text{EXTN} \\ (5)_{x_{\sim 1}} \\ (29)_{x_{\sim 1}} \end{matrix} & \begin{matrix} N \times \text{MINT} \\ (21.5)_{x_{\sim 2}} \\ (18.83)_{x_{\sim 2}} \end{matrix} & \begin{matrix} N \times \text{EXTN} \\ (29)_{x_{\sim 2}} \\ (46)_{x_{\sim 2}} \end{matrix} \\ \begin{matrix} P \times \text{EXTP} \\ X \\ \sim \end{matrix} & \begin{matrix} P \times \text{EXTN} \\ (74)_{x_{\sim 1}} \\ (46)_{x_{\sim 1}} \end{matrix} & \begin{matrix} N \times \text{MINT} \\ (21.5)_{x_{\sim 2}} \\ (18.83)_{x_{\sim 2}} \end{matrix} & \begin{matrix} N \times \text{EXTN} \\ (29)_{x_{\sim 2}} \\ (46)_{x_{\sim 2}} \end{matrix} \\ \begin{matrix} P \times \text{EXTP} \\ X \\ \sim \end{matrix} & \begin{matrix} P \times \text{EXTN} \\ (62)_{x_{\sim 1}} \\ (29)_{x_{\sim 1}} \end{matrix} & \begin{matrix} N \times \text{MINT} \\ (17.90)_{x_{\sim 2}} \\ (17.90)_{x_{\sim 2}} \end{matrix} & \begin{matrix} N \times \text{EXTN} \\ (29)_{x_{\sim 2}} \\ (29)_{x_{\sim 2}} \end{matrix} \\ \begin{matrix} P \times \text{EXTP} \\ X \\ \sim \end{matrix} & \begin{matrix} P \times \text{EXTN} \\ (23)_{x_{\sim 1}} \\ (119)_{x_{\sim 1}} \end{matrix} & \begin{matrix} N \times \text{MINT} \\ (16.76)_{x_{\sim 2}} \\ (16.76)_{x_{\sim 2}} \end{matrix} & \begin{matrix} N \times \text{EXTN} \\ (119)_{x_{\sim 2}} \\ (119)_{x_{\sim 2}} \end{matrix} \end{bmatrix}$$

Then  $\hat{Y}_{\sim}[-1] = X_{\sim 1} (X_{\sim}^i(-1) X_{\sim}^i(-1))^{-1} X_{\sim}^i(-1) Y_{\sim}(-1)$ , where  $Y_{\sim}(-1)$  is the  $(k-1)n \times 1$  vector of yields for all sites except the first, gives the transfer residuals for the first site,  $Y_{\sim 1} - \hat{Y}_{\sim}[-1]$ . The transfer residuals for the other sites are computed similarly.

From Table 1, we see that the prediction statistic is

$$P = \frac{\text{Transfer SS}}{\text{Residual SS}} = \frac{126,431,000}{88,005,000} = 1.44$$

In other words, a 44% increase in unexplained variability when predicting the ith site from the remaining sites is observed using the model with five design variables (quadratic polynomial) and four interactions with the site variables.

The next step is to assess whether this 44% increase is to be expected, or is so large as to contradict the ability to transfer results from one experiment to another. From Section 2, we have that

$$\frac{(k-1)^2}{k^2} (P-1) = 0.64 (P-1) \sim \frac{\sigma^2 \sum_{\ell=1}^{22} \theta_{\ell} \chi^2_{\ell} (1 \text{ d.f.})}{\sigma^2 \chi^2 (165 \text{ d.f.})}$$

where  $22 = k + k + (k-1)(p-2)$  and  $165 = k(n-p-1)$  and  $\theta_{\ell}$  are the eigenvalues of  $\tilde{B}_1$ .

Following the construction method outlined in Section 2,  $\tilde{C}_1$  is computed from the site variables EXTP and EXTIN. While the original values of the site variables are to be used in  $\tilde{X}_{(-i)}$ , the values in  $\tilde{T}_{(-i)j}$  are the deviations from the mean of the  $(k-1)$  sites involved. (See Appendix, equation (1).) In particular,  $\tilde{C}_1$  and its eigenvalues are:

$$\tilde{C}_1 = \begin{bmatrix} 3.870 & -2.037 & 0.067 & 0.262 & -2.163 \\ -2.037 & 1.240 & 0.459 & -0.707 & 1.045 \\ 0.067 & 0.459 & 1.457 & -1.670 & -0.313 \\ 0.262 & -0.707 & -1.670 & 1.945 & 0.170 \\ -2.163 & 1.045 & -0.313 & 0.170 & 1.261 \end{bmatrix} \quad \text{and} \quad \begin{array}{l} \theta_1 = 6.208 \\ \theta_2 = 3.564 \\ \theta_3 = 0.000 \\ \theta_4 = 0.000 \\ \theta_5 = 0.000 \end{array} .$$

Similarly,  $\tilde{C}_2$  is computed from minimum temperature (MINT) and extractable nitrogen (EXTN). This yields the eigenvalues

$$\theta_6 = 11.705, \quad \theta_7 = 2.666, \quad \theta_8 = 0.000, \quad \theta_9 = 0.000, \quad \text{and} \quad \theta_{10} = 0.000 .$$

Combining these facts we see that

$$0.64 (P-1) \sim \sigma^2 \sum_{\ell=1}^{22} \theta_{\ell} \chi^2_{\ell} (1 \text{ d.f.}) / \sigma^2 \chi^2 (165 \text{ d.f.}) .$$

We need to compare the observed value of  $[(k-1)^2/k^2](P-1) = 0.280$  with the quantiles of the distribution of  $\sum_{\ell} \theta_{\ell} \chi^2_{\ell} (1 \text{ d.f.}) / \chi^2 (165 \text{ d.f.})$ . As stated earlier, the distribution of such a linear combination of  $\chi^2_{\ell} (1 \text{ d.f.})$  as found in the numerator cannot be simplified, while the denominator is an independent  $\chi^2 (165 \text{ d.f.})$  variable.

Even though no tables exist for the distribution of  $[(k-1)^2/k^2](P-1)$ , the attained significance level may be readily estimated by Monte Carlo simulation. Using the fact that a standard normal variable squared is  $\chi^2$  (1 d.f.) and that the numerator and denominator are independent, many variables with the above distribution may be computed and the proportion which falls above the computed value of 0.280 recorded. This will give an accurate estimate of the attained significance level. In this example, we may make a further simplification. Since the number of d.f. of the denominator is so large, the residual mean square is very close to  $\sigma^2$ , the unknown experimental error, with high probability. Rewriting

$$[(k-1)^2/k^2](P-1) \doteq \frac{\sigma^2 \sum_{\ell=1}^{22} \theta_{\ell} \chi_{\ell}^2 (1 \text{ d.f.})}{165\sigma^2} \quad \text{or} \quad 165[(k-1)^2/k^2](P-1) \doteq \sum_{\ell=1}^{22} \theta_{\ell} \chi_{\ell}^2 (1 \text{ d.f.}) .$$

This implies that we need only compare  $165(0.280) = 46.2$  with the quantiles of  $\sum_{\ell} \theta_{\ell} \chi_{\ell}^2$  (1 d.f.) .

Ten thousand random variables with the distribution given above were generated. In particular, at each iteration, 22 standard normal random variables, say  $N_{\ell}$ ,  $\ell = 1, \dots, 22$ , were generated using GGUSN from the IMSL Library (Houston, Texas, U.S.A.). Then each variable was formed as the linear combination of  $\chi_{\ell}^2$  (1 d.f.) ( $N_{\ell}^2$ ) variables given above. The attained significance level is 0.236 .

Constant experimental error variances across sites have been assumed. The residual sums of squares in Table 1 make this assumption dubious. However, if the error variances are heterogeneous,  $\sigma_1^2, \dots, \sigma_k^2$ , but known, the above analysis remains valid with minor modifications of the eigenvalues. In particular, if

$$\tilde{D} = \text{diag} \left[ \left( \sigma_1^2 / \sum_{i=1}^k \sigma_i^2 \right), \dots, \left( \sigma_k^2 / \sum_{i=1}^k \sigma_i^2 \right) \right],$$

then  $\theta_1, \dots, \theta_k$  are the eigenvalues of  $\tilde{D}^{\frac{1}{2}} \tilde{C}_1 \tilde{D}^{\frac{1}{2}}$ , and  $\theta_{k+1}, \dots, \theta_{2k}$  are the eigenvalues of  $\tilde{D}^{\frac{1}{2}} \tilde{C}_2 \tilde{D}^{\frac{1}{2}}$ . Unfortunately, the last  $(k-1)(p-2)$  eigenvalues are no longer ones but are the eigenvalues of  $\tilde{D}^{\frac{1}{2}} (\tilde{I}_k - k^{-1} \tilde{J}_k) \tilde{D}^{\frac{1}{2}}$ , each with multiplicity  $(p-2)$ .

For this example, we estimate  $\tilde{D} = \text{diag}(0.067, 0.285, 0.155, 0.291, 0.203)$ . More especially, the mean square residual for each site with 33 d.f. provides a sufficiently close estimate of the unknown error variance. Then the eigenvalues are:

$\theta_1 = 4.737$	$\theta_6 = 13.496$	$\theta_{11} = \theta_{12} = \theta_{13} = 1.439$
$\theta_2 = 3.554$	$\theta_7 = 2.699$	$\theta_{14} = \theta_{15} = \theta_{16} = 1.208$
$\theta_3 = 0.0$	$\theta_8 = 0.0$	$\theta_{17} = \theta_{18} = \theta_{19} = 0.881$
$\theta_4 = 0.0$	$\theta_9 = 0.0$	$\theta_{20} = \theta_{21} = \theta_{22} = 0.472$
$\theta_5 = 0.0$	$\theta_{10} = 0.0$	

A Monte Carlo simulation yielded a significance level of 0.240 for  $P-1$ .

#### Acknowledgments

The authors acknowledge their appreciation to the editors and referees for valuable suggestions, and to Dr. James A. Silva of the University of Hawaii and Dr. T. S. Gill of U.S.A.I.D. for their support and encouragement. We also acknowledge the use of the data from the Benchmark Soils Project, a U.S.A.I.D. Project with the Department of Agronomy and Soil Science, University of Hawaii.

References

- Cochran, W. G. and Cox, G. M. (1957). Experimental Designs, 2nd edition. John Wiley and Sons, New York.
- Colwell, J. D. (1967). Calibration and assessment of soil tests for estimating fertilizer requirements. Australian Journal of Soil Research 5, 275-293.
- Gardner, M. J. (1972). On using an estimated regression line in a second sample. Biometrika 59, 263-274.
- Searle, S. R. (1971). Linear Models. John Wiley and Sons, New York.
- Soil Survey Staff, Soil Conservation Service, U.S. Department of Agriculture. (1975). Soil Taxonomy: A Basic System of Soil Classification for Making and Interpreting Soil Surveys. Agriculture Handbook 436, U.S. Government Printing Office, Washington, D. C.

Appendix

Distribution of (P - 1)

Determining the distribution of P - 1 begins by expressing the combined vector of transfer residuals as a linear combination of kn independent normal errors  $\epsilon$ , each with common variance  $\sigma^2 > 0$ ; i.e.,  $\epsilon \sim N(0, \sigma^2 I)$ . Then

$$P - 1 = \frac{\epsilon' (R'R - B) \epsilon}{\epsilon' B \epsilon}$$

where  $\epsilon' B \epsilon$  is the pooled sum of squared residuals divided by  $\sigma^2$ . Next the numerators and denominators are shown to be independent. Lastly the eigenvalues of the quadratic form in the numerator, say  $\theta_1, \dots, \theta_q$  with  $q = q(n, k, p, m_1, m_2)$ , are determined.

First we develop the data matrices used for estimating the response surface required at each step. Let  $\tilde{X}$  denote the  $(n \times p)$  matrix of x-variables with each column centered at zero. Further, let  $\tilde{x}_1$  and  $\tilde{x}_2$  be two orthogonal columns of  $\tilde{X}$  which will be used to form interactions with site variables. Then, using the notation developed in Section 2, the true response function at the  $i$ th site is given by  $X_i \beta$ , where

$$X_i = \left[ \tilde{X} : T_{i1} \otimes \tilde{x}_1 : T_{i2} \otimes \tilde{x}_2 \right] , \quad (1)$$

where  $\otimes$  denotes the right Kronecker product. Note that without loss of generality, centered site variables may be used to form interactions. This follows from the fact that  $\tilde{x}_1$  and  $\tilde{x}_2$  are centered at zero. Then the predicted values for the  $i$ th site based on the remaining  $(k-1)$  sites is  $\hat{Y}_{[-i]} = X_i b_{[-i]}$ , where  $b_{[-i]}$  is estimated only from the remaining  $(k-1)$  sites.

In order to compute  $b_{[-i]}$ , we require the data matrix for all sites except the  $i$ th, say  $X_{(-i)}$ . In particular,

$$X_{(-i)} = \left[ \mathbf{1}_{k-1} \otimes \tilde{X} : T_{(-i)1} \otimes \tilde{x}_1 : T_{(-i)2} \otimes \tilde{x}_2 \right] , \quad i = 1, \dots, k .$$

In Section 2, the yields were adjusted for the individual site means. In terms of the true response function, the adjusted yields are given by

$$Y_i = X_i \beta + (I_n - n^{-1} J_n) \epsilon_i , \quad i = 1, \dots, k$$

and, similarly, the combined  $(k-1)n \times 1$  vector of adjusted yields excluding the  $i$ th site is

$$Y_{(-i)} = X_{(-i)} \beta + \left[ I_{(k-1)n} \otimes (I_n - n^{-1} J_n) \right] \epsilon_{(-i)} , \quad i = 1, \dots, k ,$$

where  $\epsilon_{(-i)}$  is analogously defined.

Since the adjusted yields are not independent,  $\tilde{b}_{[-i]}$  is computed using the methods of generalized least squares (Searle 1971, Section 5.8). Noting that the covariance matrix of the adjusted yields is  $[\tilde{I}_{(k-1)} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n)]$ , which is idempotent and  $[\tilde{I}_{(k-1)} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n)]\tilde{X}_{(-i)} = \tilde{X}_{(-i)}$ ,

$$\begin{aligned} \tilde{b}_{[-i]} &= \left( \tilde{X}_{(-i)}' \tilde{X}_{(-i)} \right)^{-1} \tilde{X}_{(-i)}' \left[ \tilde{I}_{(k-1)} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n) \right] \left\{ \tilde{X}_{(-i)}\beta + \left[ \tilde{I}_{(k-1)} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n) \right] \tilde{\epsilon}_{(-i)} \right\} \\ &= \tilde{\beta} + \left( \tilde{X}_{(-i)}' \tilde{X}_{(-i)} \right)^{-1} \tilde{X}_{(-i)}' \tilde{\epsilon}_{(-i)}, \quad i = 1, \dots, k \end{aligned}$$

Therefore the  $(n \times 1)$  vector of transfer residuals is given by

$$\tilde{Y}_i - \hat{\tilde{Y}}_{[-i]} = (\tilde{I}_n - n^{-1}\tilde{J}_n)\tilde{\epsilon}_i - \tilde{X}_i \left( \tilde{X}_{(-i)}' \tilde{X}_{(-i)} \right)^{-1} \tilde{X}_{(-i)}' \left[ \tilde{I}_{(k-1)} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n) \right] \tilde{\epsilon}_{(-i)}$$

Next we must determine  $R$ . First consider  $\tilde{Y}_1 - \hat{\tilde{Y}}_{[-1]}$ . For simplicity, let  $\tilde{e}_i = (\tilde{I}_n - n^{-1}\tilde{J}_n)\tilde{\epsilon}_i$ ,  $\tilde{e}_{(-i)} = [\tilde{I}_{k-1} \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n)]\tilde{\epsilon}_{(-i)}$  and  $\tilde{e} = [\tilde{I}_k \otimes (\tilde{I}_n - n^{-1}\tilde{J}_n)]\tilde{\epsilon}$ . Then

$$\begin{aligned} \hat{\tilde{Y}}_{[-1]} &= \tilde{X}_1\beta + \left\{ \left[ (k-1)^{-1} \tilde{1}_{k-1}' \otimes P_X \right] + \left[ \left( \tilde{T}_{11}(\tilde{T}'_{(-1)1}\tilde{T}_{(-1)1}) \right)' \otimes P_1 \right] \right. \\ &\quad \left. + \left[ \left( \tilde{T}_{12}(\tilde{T}'_{(-1)2}\tilde{T}_{(-1)2}) \right)' \otimes P_2 \right] \right\} \tilde{e}_{(-1)}, \end{aligned}$$

where  $\tilde{1}_k$  is a  $(k \times 1)$  vector of ones,  $P_X = \tilde{X}\tilde{X}'^{-1}\tilde{X}'$ , and  $P_j = \tilde{x}_j(\tilde{x}'_j\tilde{x}_j)^{-1}\tilde{x}'_j$ ,  $j = 1, 2$ . Note that  $\tilde{T}_{1j}(\tilde{T}'_{(-1)j}\tilde{T}_{(-1)j})^{-1}\tilde{T}'_{(-1)j}$ ,  $j = 1, 2$  are  $[1 \times (k-1)]$  row vectors. Augmenting each by a zero in the first position, we can define the  $(1 \times k)$  row vectors:

$$\tilde{t}_{11} = \left[ 0 : \tilde{T}_{11}(\tilde{T}'_{(-1)1}\tilde{T}_{(-1)1})^{-1}\tilde{T}'_{(-1)1} \right] \text{ and } \tilde{t}_{12} = \left[ 0 : \tilde{T}_{12}(\tilde{T}'_{(-1)2}\tilde{T}_{(-1)2})^{-1}\tilde{T}'_{(-1)2} \right]$$

Then

$$\hat{\tilde{Y}}_{[-1]} = \tilde{X}_1\beta + \left[ (k-1)^{-1} \tilde{1}_{k-1}' \otimes P_X \right] \tilde{e}_{(-1)} + \left[ \tilde{t}_{11} \otimes P_1 + \tilde{t}_{12} \otimes P_2 \right] \tilde{e}$$

and

$$\underline{Y}_1 - \hat{\underline{Y}}[-1] = \left\{ \left[ \underline{I}_{\underline{n}} : -(k-1)^{-1} \underline{1}_{k-1}' \otimes \underline{P}_X \right] - \left[ \underline{t}_{11} \otimes \underline{P}_1 + \underline{t}_{12} \otimes \underline{P}_2 \right] \right\} \underline{\epsilon} .$$

Next define  $\underline{t}_{i1}$  and  $\underline{t}_{i2}$ ,  $i=2, \dots, k$ , analogously to  $\underline{t}_{11}$  and  $\underline{t}_{12}$ , i.e.,  $\underline{t}_{21}$  is the  $(1 \times k)$  row vector formed from  $\underline{T}_{21} (\underline{T}'_{(-2)1} \underline{T}_{(-2)1})^{-1} \underline{T}'_{(-2)1}$  with a zero element inserted as the second element, etc. Then

$$\begin{aligned} & \left[ (\underline{Y}_1 - \hat{\underline{Y}}[-1])' : (\underline{Y}_2 - \hat{\underline{Y}}[-2])' : \dots : (\underline{Y}_k - \hat{\underline{Y}}[-k])' \right] \\ & = \left[ \underline{A} - (\underline{T}_1 \otimes \underline{P}_1) + (\underline{T}_2 \otimes \underline{P}_2) \right] \left[ \underline{I}_k \otimes (\underline{I}_n - n^{-1} \underline{J}_n) \right] \underline{\epsilon} , \end{aligned}$$

where  $\underline{A} = \left\{ \underline{I}_k \otimes [\underline{I}_n + (k-1)^{-1} \underline{P}_X] \right\} - [(k-1)^{-1} (\underline{J}_k \otimes \underline{P}_X)]$ ,

$$\underline{T}_1 = [\underline{t}_{11}' : \underline{t}_{21}' : \dots : \underline{t}_{k1}']' \quad \text{and} \quad \underline{T}_2 = [\underline{t}_{12}' : \underline{t}_{22}' : \dots : \underline{t}_{k2}']' .$$

With the above representation and  $\underline{P}_X \underline{1} = \underline{0}$ ,  $\underline{T}_j \underline{J}_k = \underline{0}$ ,  $\underline{P}_X \underline{P}_j = \underline{P}_j$ ,  $j=1,2$ , and

$$\underline{B} = \left[ \underline{I}_k \otimes (\underline{I}_n - n^{-1} \underline{J}_n) \right] \left[ \underline{I}_k \otimes (\underline{I}_k \underline{P}_X) \right] \left[ \underline{I}_k \otimes (\underline{I}_n - n^{-1} \underline{J}_n) \right] ,$$

it can be shown that

$$\underline{\epsilon}' (\underline{R}' \underline{R} - \underline{B}) \underline{\epsilon} = k^2 (k-1)^{-2} \underline{\epsilon}' \left\{ \left[ (\underline{I}_k - k^{-1} \underline{J}_k) \otimes \underline{P}_X \right] + (\underline{D}_1 \otimes \underline{P}_1) + (\underline{D}_2 \otimes \underline{P}_2) \right\} \underline{\epsilon} ,$$

where

$$\underline{D}_j = k^{-2} (k-1)^2 \left[ \underline{T}_i' \underline{T}_i - k(k-1)^{-1} (\underline{T}_i + \underline{T}_i') + (k-1)^{-1} (\underline{J}_k \underline{T}_i + \underline{T}_i' \underline{J}_k) \right] , \quad j=1,2 .$$

It now easily follows that  $\underline{B}_1 \underline{B} = \underline{0}$ ; i.e., the numerator and denominator of  $k^{-2} (k-1)^2 (\underline{P}-1)$  are independent.

Finally we need to find the eigenvalues of

$$\underline{B}_1 = \left[ (\underline{I}_k - k^{-1} \underline{J}_k) \otimes \underline{P}_X \right] + (\underline{D}_1 \otimes \underline{P}_1) + (\underline{D}_2 \otimes \underline{P}_2) .$$

Define

$$U_0 = [ \|x_1\|^{-2} x_1 : \|x_2\|^{-2} x_2 ] ,$$

where  $\|x_j\|^2 = \sum_{\ell=1}^n x_{j\ell}^2$ . Let  $U_1$  be any orthonormal matrix of rank  $p$  such that  $P_X[U_0 : U_1] = [U_0 : U_1]$ . Finally  $U$  will be any orthonormal matrix of rank  $n$  with  $U = [U_0 : U_1 : U_2]$  for some  $U_2$ . Then the eigenvalues of  $B_1$  are the same as the eigenvalues of  $(I_k \otimes U)' B_1 (I_k \otimes U)$ . By construction,  $U' P_X U = \text{diag}(I_p, 0_{n-p})$ ,  $U' P_1 U = \text{diag}(1, 0_{n-1})$ , and  $U' P_2 U = \text{diag}(0, 1, 0_{n-2})$ . Therefore,

$$\begin{aligned} (I_k \otimes U)' B_1 (I_k \otimes U) &= (I_k - k^{-1} J_k) \otimes \text{diag}(0, 0, I_{p-2}, 0_{n-p}) \quad (2) \\ &+ \left[ (I_k - k^{-1} J_k) + D_1 \right] \otimes \text{diag}(1, 0_{n-1}) \\ &+ \left[ (I_k - k^{-1} J_k) + D_2 \right] \otimes \text{diag}(0, 1, 0_{n-2}) . \end{aligned}$$

Since the above three matrices are orthogonal, the required eigenvalues of  $B_1$ ,  $\theta_1, \dots, \theta_q$ , are the  $(k-1)(p-2)$  eigenvalues of the first matrix which are identically equal to one, the  $k$  eigenvalues of  $C_1 = (I_k - k^{-1} J_k) + D_1$  and the  $k$  eigenvalues of  $C_2 = (I_k - k^{-1} J_k) + D_2$ . This is the result stated at the end of Section 2.

Finally, the effect of unequal experimental error variances across sites must be investigated. Let  $\sigma_i^2$ ,  $i = 1, \dots, k$ , denote the variances of the yields of the  $i$ th site; i.e.,  $E[\epsilon\epsilon'] = \text{diag}(\sigma_i^2 I_n) = D \otimes I_n$ , where  $D = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ . It follows that

$$\frac{(k-1)^2}{k^2} \cdot (P-1) = \frac{\left[ (D^{-\frac{1}{2}} \otimes I_n) \epsilon \right]' \left[ (D^{\frac{1}{2}} \otimes I_n) B_1 (D^{\frac{1}{2}} \otimes I_n) \right] \left[ (D^{-\frac{1}{2}} \otimes I_n) \epsilon \right]}{\left[ (D^{-\frac{1}{2}} \otimes I_n) \epsilon \right]' \left[ (D^{\frac{1}{2}} \otimes I_n) B (D^{\frac{1}{2}} \otimes I_n) \right] \left[ (D^{-\frac{1}{2}} \otimes I_n) \epsilon \right]}$$

where  $(D^{-\frac{1}{2}} \otimes I_n) \epsilon$  is a vector of independent standard normal random variables. Note that  $B(D \otimes I_n) = B(D \otimes I_n)B$ ,  $(D^{\frac{1}{2}} \otimes I_n)BDB_1(D^{\frac{1}{2}} \otimes I_n) = 0$ , and the numerator and denominator remain independent. Similarly,  $(I_k \otimes U)$  and  $(D^{\frac{1}{2}} \otimes I_n)$  also commute and the eigenvalues of the numerator are the eigenvalues of

$$\begin{aligned} & (I_k \otimes U)(D^{\frac{1}{2}} \otimes I_n)B_1(D^{\frac{1}{2}} \otimes I_n)(I_k \otimes U) \\ &= D^{\frac{1}{2}}(I_k - k^{-1}J_k)D^{\frac{1}{2}} \otimes \text{diag}(0, 0, I_{(p-2)}, 0_{n-p}) \\ & \quad + D^{\frac{1}{2}}C_1D^{\frac{1}{2}} \otimes \text{diag}(1, 0_{n-1}) + D^{\frac{1}{2}}C_2D^{\frac{1}{2}} \otimes \text{diag}(0, 1, 0_{n-2}) \end{aligned}$$

Thus the eigenvalues of the numerator consist of the eigenvalues of  $D^{\frac{1}{2}}(I_k - k^{-1}J_k)D^{\frac{1}{2}}$  with multiplicity  $(p-2)$  and the eigenvalues of both  $D^{\frac{1}{2}}C_1D^{\frac{1}{2}}$  and  $D^{\frac{1}{2}}C_2D^{\frac{1}{2}}$ , which may be computed as in the equal variance case.

Also, for small  $n$  and  $k$  the eigenvalues of the denominator may be similarly calculated. In practice, however, the degrees of freedom of sum of squared residuals for each site will be large enough to assume that the resulting consistent within-site estimated error variance sufficiently approximates  $\sigma_i^2$  and these eigenvalues will not be required. In particular,

$$k(n-p) \left[ \frac{(k-1)^2}{k^2} (p-1) \right] \xrightarrow{p} \frac{\epsilon' B_1 \epsilon}{\sum_i \sigma_i^2}$$

and the problem reduces to computing the eigenvalues of the numerator.