

Partitioning Error Components for Accuracy-Assessment of Near-Neighbor Methods of Imputation

Albert R. Stage and Nicholas L. Crookston

Abstract: Imputation is applied for two quite different purposes: to supply missing data to complete a data set for subsequent modeling analyses or to estimate subpopulation totals. Error properties of the imputed values have different effects in these two contexts. We partition errors of imputation derived from similar observation units as arising from three sources: observation error, the distribution of observation units with respect to their similarity, and pure error given a particular choice of variables known for all observation units. Two new statistics based on this partitioning measure the accuracy of the imputations, facilitating comparison of imputation to alternative methods of estimation such as regression and comparison of alternative methods of imputation generally. Knowing the relative magnitude of the errors arising from these partitions can also guide efficient investment in obtaining additional data. We illustrate this partitioning using three extensive data sets from western North America. Application of this partitioning to compare near-neighbor imputation is illustrated for Mahalanobis- and two canonical correlation-based measures of similarity. *FOR. SCI.* 53(1):62–72.

Keywords: most similar neighbor, k -nn inference, missing data, landscape modeling

IMPUTATION METHODS are important tools for completing data sets in which some observation units lack observed values for a portion of their attributes. The objective is to impute a value as close to “truth” for each missing value in the observation unit as if it were examined in great detail for all attributes. Criteria for imputations to support this objective are essentially different from criteria for estimates of population totals. The difference is that pure error, rather than being a nuisance, is of real value for subsequent resource analyses and displays. These analyses are often nonlinear optimizations or simulations. For them to be realistic, the structure of the variances and covariances among attributes inherent in the population should be preserved in the data set. Even for display purposes, omission of pure error will cause the range of the displayed values to be contracted. Unfortunately, these inherently useful variances may be combined with variances attributable to the methodology used in the sampling and imputation processes. This mixture complicates choice among analytical methods for imputation. In this report we provide statistics based on a partitioning of the error components that facilitate finding a closer approximation of “truth.” We partition imputation errors independently for each variable in the data set, although the joint distribution of their error components would be of interest for some applications.

Imputation uses values of variables measured for all observation units (\mathbf{X}) to guide the imputation of values of \mathbf{Y} that are measured only for a sample subset of the observation units (the *Reference* set) to those units for which the \mathbf{Y} values are missing (the *Target* set). Both \mathbf{X}_i and \mathbf{Y}_i may be vectors of attributes for

the i th observation unit. Near-neighbor imputation selects units from the reference set to serve as surrogates for members of the target set using a measure of similarity based on the \mathbf{X} values. Choice of a particular measure of similarity, in turn, may depend on the relation of the \mathbf{Y} values to the \mathbf{X} values. Elements of \mathbf{Y}_i and \mathbf{X}_i , y_i and x_i , will be subscripted only to identify the i th observation unit. T and R will be used as additional subscripts when it is relevant to indicate that a Reference observation unit is being used as if it were a Target unit (hence a “pseudo-target”). Unit identifying subscripts (i or j) will be omitted when the variables are referred to collectively. $\text{Var}(\cdot)$ capitalized will be used for expected values, lower case $\text{var}(\cdot)$ for statistics calculated from the data.

Imputation from near-neighbor observations is often used for classification. However, when the “classes” are arbitrary intervals on scales of essentially continuous variables, we argue that the imputation should be based directly on the scales of the underlying continuous variables. If classes are needed for display purposes, the classification algorithm should use the imputed data. We will not consider in this article errors in classification in which the classes are inherently discrete, requiring the concept of “membership.” For discrete classes, other methods for classification such as using a discriminant function may be more appropriate than near-neighbor. For example, classification by a discriminant function may assign different classes to members of a target/reference pair of near neighbors because the discriminating boundary passes between them, whereas near-neighbor imputation would assign the target to the same class as the reference member of the pair. However, there is a parallel process of partitioning the error

Albert R. Stage (Retired), Moscow Forestry Sciences Laboratory, Rocky Mountain Research Station, 1221 S. Main St., Moscow, ID 83843—Fax: (208) 883-2318; astage@moscow.com. Nicholas L. Crookston, Moscow Forestry Sciences Laboratory, Rocky Mountain Research Station, 1221 S. Main St., Moscow, ID 83843—ncrookston@fs.fed.us.

Acknowledgments: Interested readers are greatly in debt to one of our anonymous reviewers for insisting in great detail on precise use of wording in this paper. The authors also thank the other reviewer, who said the results are so intuitively obvious as to wonder why it had not been done earlier for giving us the incentive to press on with revision. We are convinced that the result is much improved, but any remaining deficiencies are certainly our responsibility.

Manuscript received June 5, 2005, accepted October 11, 2006

This article was written by U.S. Government employees and is therefore in the public domain.

sources in imputation of discrete variables that is beyond the scope of this article.

Error properties of estimates derived from imputation differ from those of regression-based estimates because the two methods include a different mix of error components. For example, the reference-set data may not be beyond reproach because of measurement error. These error properties influence how we evaluate quality of the imputations, compare alternative methods for imputation, and invest in data collection. Commonly computed statistics that compare imputed values to those of a presumably similar observation unit mask methodological differences in this cloud of variation. We address this problem by partitioning the variation into components that can be estimated from the reference set. Then, new statistics based on this partitioning are presented for assessing the accuracy of imputation methods.

Several questions may be answered using these error components:

1. How does the accuracy of imputed \mathbf{Y} values compare to accuracy of estimates from regression, stratum means, or other model-based estimates?
2. How large is the error caused by imputing values to a target unit from reference units where there is substantial difference in their \mathbf{X} values? Is there room for improvement by obtaining additional reference observations to fill gaps in their distribution? How is this error component affected by the choice of a particular measure of similarity?
3. How is the accuracy of imputation affected by the choice of variables and their transformations?
4. What is the effect on imputed values of pooling k reference observations?
5. How do the measurement accuracies compare to components of variation from other sources?
6. Would investments in additional data be more efficient if used to obtain information on variables to be added to the target set (new \mathbf{X} variables), to refine the estimates of the \mathbf{X} values already included, or to obtain data on additional units for the reference set?

Resolution of these questions requires quantitative estimates of the sources of imputation error. These estimates can be obtained from the information in the n observation units in the reference data. In analysis of data in the reference-set data, although we will use some of the data as though they were targets, there is no difference in their approximation of “truth,” no intrinsic differences between “observed” and “predicted.” We are simply describing the properties of differences between members of pairs of observations. When the value to be imputed is a weighted average of k near-neighbors, then its error properties are derived from the error properties of the k separate pairs and the weights defined by the particular k -nn procedure.

Bootstrap and cross-validation methods for answering some of these questions have been developed for imputation methods other than near-neighbor (Shao and Sitter 1996) or for classification with k -nn (Mullin and Sukthankar 2000).

Neither of these papers has addressed the problem of partitioning the errors as to sources. Moeur and Stage (1995) used data-splitting and jackknife methods to evaluate capability of most-similar neighbor (MSN) to reproduce the variance and covariance structure of the reference data and to compare error rates to those obtained by stratified sampling and regression. Their analyses of errors also included variation in the coefficients in the measure of similarity caused by sequentially omitting 1.7% of their data as well as the difference between the observed and imputed \mathbf{Y} values for the pair selected by the calculated similarity measure.

Splitting data into “calibration” and “validation” subsets, which was intended to reduce bias in error estimates, introduces a different bias into estimates of imputation errors. The withheld reference observations in sparsely represented parts of \mathbf{X} -space could have supplied imputations for nearby target observations. In the analysis of imputation error, however, those targets will be paired with a more remote reference observation, thereby increasing the *estimated* error. A further disadvantage of the jackknife procedure is that it may increase the estimate of error by increasing the mean-square bias. Targets in the midst of a cloud of reference observations may be paired with an observation from any direction. Targets at the edge of a cloud, however, will likely be paired with a more central point. If there is a trend in the \mathbf{Y} values with distance from the center of the cloud, then the asymmetry of direction to the reference introduces bias in the imputed value. Withholding data unnecessarily increases this bias. The jackknife procedure using a single reference observation as if a target minimizes this bias by using the full range of data. Other methods to reduce this bias in k -nn imputation have been evaluated by Malinen (2003).

A statistic commonly used to evaluate imputation error estimates the root-mean-square differences between reference and target observations by withholding each observation unit in the reference set while searching for its similar neighbor in the remainder of the reference set. The term RMSE (root-mean-square error) used for this statistic is unfortunate (e.g., Moeur and Stage 1995, Crookston et al. 2002). The term as used in imputation includes different components of error than the same term used in a regression or sampling context. Therefore, we use the term mean-squared difference (MSD) for the statistic describing squared differences in a pair of similar observations. Thus, our partitioning is applicable for evaluating any of the near-neighbor methods of imputation that are judged on the basis of sums of squared errors.

We use the term “distance” for the value produced by the function measuring dissimilarity between the i th and j th pair of observation units. Although Podani (2000) cites more than 60 distance functions, those most widely used for imputation are of the quadratic form,

$$d_{ij}^2 = (\mathbf{X}_i - \mathbf{X}_j)\mathbf{W}(\mathbf{X}_i - \mathbf{X}_j)', \quad (1)$$

where \mathbf{X}_i is the $(1 \times p)$ vector of \mathbf{X} -variables for the i th target observation unit, \mathbf{X}_j is the $(1 \times p)$ vector of \mathbf{X} -variables for the j th reference observation unit, and \mathbf{W} is a $(p \times p)$ symmetric matrix of weights.

If the weight matrix, \mathbf{W} , is the diagonal identity matrix, then we have a simple Euclidean distance (squared). As a variation of Euclidean distance, some analysts empirically vary the diagonal elements to improve the imputation. If correlations among the variates are to be considered, then the inverse of their correlation matrix is used for \mathbf{W} to produce a Mahalanobis distance—a distance function that plays a key role in estimating the error components. MSN distances are of the same form with \mathbf{W} derived from analyses of canonical correlation (Moeur and Stage 1995), canonical regression (Stage and Crookston 2002), or of canonical correspondence (Ohmann and Gregory 2002). With a simple transformation of the \mathbf{X} values to $x_i/\sqrt{\sum_{i=1}^p x_{ii}^2}$, the quadratic form with identity matrix for \mathbf{W} also includes spectral analysis imputation as used by Sohn et al. (1999).

Our following presentation is in four sections: (1) defining error sources in the process of imputation, (2) partitioning MSD into components arising from these sources, (3) presenting some new statistics based on the partitioning relevant to the key questions stated above, and (4) applying these statistics to three extensive data sets.

Components of Error

Variation in imputed values arises from both natural variability of attributes of the ecosystem and from the measurement and analytical procedures used to describe the ecosystem. Although natural variability is useful in analyses requiring the completed data set, variation introduced by measurement and analytical procedures is a nuisance to be reduced.

Imputation error arises from four sources for a given set of \mathbf{X} and \mathbf{Y} variables:

1. Measurement Errors of the \mathbf{Y} Values in the Reference Set.—These errors are defined as

$$\varepsilon_{Yj} = y_j - y_j^* \quad (2)$$

in which the starred variables represent the true, but unknown, values. The ε_{Yj} are not properties of the ecosystem being described, but properties of the accidents of how we observed it. The measurement errors may arise from using a sample-based estimate as if it were a complete census within the j th unit, from changes during elapsed time since observation, from lack of standardization among different observers or their instruments, or any combination of such causes. These errors often are assumed to be zero (e.g., Moeur and Stage 1995). We now relax that assumption because in some applications, errors from this source have been quite large relative to total error. We assume that the measurement errors can be rendered unbiased and are independent of the true y_j^* and of the observed \mathbf{X} values.

2. Pure Error.—That there exists a relation between the \mathbf{Y} and the \mathbf{X} variables is a key premise of near-neighbor inference. For a given set of \mathbf{X} variables, the departure of an element of \mathbf{Y}_j^* from the underlying true but unknown model is termed pure error:

$$\varepsilon_{Pj} = y_j^* - g(\mathbf{X}_j). \quad (3)$$

Magnitude of the pure error (ε_{Pj}) depends on the particular choice of \mathbf{Y} and \mathbf{X} variables. By definition, pure error, which arises from effects not associated with the \mathbf{X} vari-

ables is independent of the \mathbf{X} variables and has zero expectation. Examples of omitted factors are myriad, but would include predicting species composition (the \mathbf{Y} variables) from Landsat spectra (the \mathbf{X} variables), but omitting elevation as an additional \mathbf{X} variable that might improve the imputation.

Not so obvious as a source of pure error would be the effect of lack of accurate registration between the \mathbf{Y} -variable observation units located on the ground and the paired \mathbf{X} -variable observation units from a remote sensing platform. In effect, the observed values of \mathbf{X}_j from a complete census from the erroneous position are just a differently defined variable for imputation than the \mathbf{X}_j values from a properly registered observation unit. Therefore, variation from lack of registration would contribute to pure error that might be reduced by improving registration.

From Equations 2 and 3,

$$y_j = g(\mathbf{X}_j) + \varepsilon_{Pj} + \varepsilon_{Yj}, \quad (4)$$

in which the error components include measurement error (ε_{Yj}), and pure error (ε_{Pj}). These two components of imputation error are shown graphically in Figure 1. Pure error and measurement error are inseparable in many data sets. To estimate pure error alone requires an external estimate of the measurement error. For example, if the observation unit is a spatial polygon represented by the mean of each of the attributes over a number of plots within the polygon, then the estimated variance-of-the-mean would provide the sampling portion of measurement variance to be subtracted to leave pure error.

3. Factors Affecting the Availability and Similarity of Reference Observation Units to Serve as Surrogates for the Target Units.—This component depends on both the choice of a distance function and on the distribution of observation units in the space spanned by the \mathbf{X} variables. Ideally, all the target data should be within the span of the reference data. The denser the data, the shorter will be the

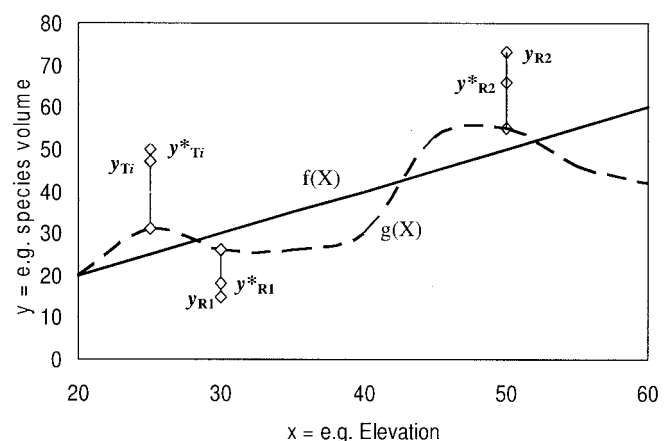


Figure 1. Error components for imputing y_{Rj} (e.g., species volume) to a target observation at x_{Ti} from one of two reference observations in a one-dimensional space of \mathbf{X} (e.g., elevation). Pure error (ε_{Pj}) is the vertical distance from y_i^* to the dashed line $g(\mathbf{X})$. Measurement error (ε_{Yi}) is the vertical distance between y_i^* and y_i . Model lack-of-fit ($\varepsilon_{L(Xj)}$) is the vertical separation between the dashed $g(\mathbf{X})$ and solid $f(\mathbf{X})$ lines.

average distance between a target unit and its nearest surrogate in the reference set, and shorter distances usually imply greater similarity. The magnitude of this effect can be appreciated by comparing the distribution of distances to nearest neighbors among the reference data to the distribution of distances from each target observation to its nearest neighbor in the reference set. The distances between the real targets and their near neighbors in the reference set usually would be, on average, shorter than the distances among members of the reference set. Thus, estimated errors based only on the reference set will be biased upward. Effects of the density and range of the data apply to all methods of imputation and are determined by the inventory design.

4. The Choice of k , the Number of Reference Observations, and their Relative Weights in k -nn Methods of Estimating Y Values as a Weighted Average of k Near Neighbors.

Error analyses we propose are based on the data in the reference set. Inferences about the error properties of the estimates for the entire population based on these analyses depend on the extent to which the reference set represents the target set. As with inferences about any population parameter, appropriate randomization is a prerequisite to the assumption that the partitioning of error based on the reference set will apply to imputations for the real target set.

Imputation Error Statistics Based on the Reference Set

In the imputation context $\sum_i (y_{Ti} - y_{Rj})^2/n$ is the statistic commonly reported as “squared error” based on the n observation units in the reference set. We use the term mean square difference (MSD) for it to emphasize that it is not an “error”—rather, it is simply a function of the difference between two co-equal observations, neither of which is any more “true” than the other. In this and the expressions to follow, the subscript j identifies the reference observation to be imputed to the i th pseudo-target observation unit. For each observation unit i , the value of j is determined by the minimum of d_{ij}^2 in Equation 1. In k -nn imputation y_{Rj} is replaced by an average of k values of y_m using a weighting rule for the particular flavor of k -nn inference, where m is from the set of indices of the k observations selected as near-neighbors. We will develop the partitioning of error components for $k = 1$ because the notation is much more compact. However, the extension to $k > 1$ introduces no new concepts and will be treated when we discuss the choice of k as an error source.

Each member of the pairs being averaged in MSD includes stochastic components that do not change whether the observation unit is playing the role of target or reference. Each pair also includes a component determined by the distribution of the \mathbf{X} values within the reference set. Thus, the statistics we compute are conditional on distribution of \mathbf{X} values in the reference set, and may be used to guide decisions on how or whether to augment that reference set. The stochastic components, pure error and measurement error, are assumed to be drawn from distributions having zero mean and zero covariance. Therefore, for both stochastic error sources:

$$E[\varepsilon_{pj}] = E[\varepsilon_{yj}] = 0; \quad \text{Var}(\varepsilon_p) = E[\sum_j \varepsilon_{pj}^2/n];$$

$$\text{Var}(\varepsilon_Y) = E[\sum_j \varepsilon_{Yj}^2/n]; \quad E[\varepsilon_{Yj}\varepsilon_{pj}] = 0. \quad (5)$$

We will define the estimated variances of the stochastic error terms $\text{var}(\varepsilon_p)$ and $\text{var}(\varepsilon_Y)$ as the average over the reference set, dividing by n rather than $(n - p)$ because the error terms are defined relative to true values rather than from a computed mean.

We first introduce the measurement error from Equation 2 into an addend of MSD:

$$(y_{Ti} - y_{Rj})^2 = (y_{Ti}^* + \varepsilon_{Yi} - y_{Rj}^* - \varepsilon_{Yj})^2. \quad (6)$$

Expanding Equation 6 on the starred terms from Equation 3, we have

$$(y_{Ti} - y_{Rj})^2 = [g(\mathbf{X}_{Ti}) + \varepsilon_{pi} + \varepsilon_{Yi} - g(\mathbf{X}_{Rj}) - \varepsilon_{pj} - \varepsilon_{Yj}]^2. \quad (7)$$

Averaging over the n pseudo-target units (y_{Ti}) in Equation 7 assuming ε_{pj} and ε_{Yj} are independent of each other and using Equation 6, the expectation of MSD becomes:

$$\begin{aligned} E[\text{MSD}] &= E[\sum_i (y_{Ti} - y_{Rj})^2/n] \\ &= \sum_i [g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n + 2\text{Var}(\varepsilon_Y) + 2\text{Var}(\varepsilon_p). \end{aligned} \quad (8)$$

The term $\sum_i [g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n$ in Equation 8 is, therefore, the error component arising from the distance between a pseudo-target point and its selected surrogate reference point. Note that in addition to the distance error component, the other error variances are included twice in MSD.

Estimating Pure Error and Measurement Error

In a regression context, sums of squares for pure error plus measurement error can be estimated from differences between the y values for observations having the same \mathbf{X} values. The corresponding concept in imputation is for observations separated by zero Mahalanobis distance. Mahalanobis distances are calculated in the space spanned by the normalized, but uncorrelated, \mathbf{X} variables. The Mahalanobis distance was selected because other distance functions may transform the \mathbf{X} variables such that the dimension of the space spanned by the transformed \mathbf{X} variables is of lower dimension than the original \mathbf{X} space. Zero distances in the space of reduced dimension would not necessarily indicate that \mathbf{X}_{Ti} for a target unit is identical to the \mathbf{X}_{Rj} for the selected reference unit. We argue that an estimate of twice the sum of variances of pure error and measurement error can be obtained by averaging the squared differences for some fraction of the units with short Mahalanobis distances. We call this estimate MMSD(0), adding an initial M and the (0) to suggest it is derived from pairs of units with Mahalanobis distances of close to zero. Using Equation 8,

$$E[\text{MMSD}(0)] = 2 \text{Var}(\varepsilon_p) + 2 \text{Var}(\varepsilon_Y) + \text{bias}, \quad (9)$$

where the bias equals the amount by which the mean of the squared distance component (as in Equation 8 but averaged over only the observation units with close-to-zero distances)

differs from zero. Note that whereas MSD may be derived from any of the many distance functions, MMSD(0) always uses Mahalanobis distance.

The estimate is biased by the average of $[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2$ in MMSD(0). The bias might be reduced by regressing the values of $(y_{Ti} - y_{Rj})^2$ on their distances where the near-neighbor pairings are determined using a Mahalanobis distance function. The intercept of this regression may provide an improved estimate of MMSD(0) by extrapolation to zero distance. However, for some obstreperous \mathbf{Y} variables, the squared deviations decline with increasing distance so that the intercept is above the mean. This circumstance indicates that the \mathbf{X} variables do not measure similarity for those elements of \mathbf{Y} or that their stochastic components are heteroskedastic.

Estimating Distance Component

The distance component depends only on the range and density of the \mathbf{X} values and on the measure of similarity used to select the near neighbor(s). Equation 8 showed that $E[\text{MSD}]$ is comprised of the distance component, $\sum_i [g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n$ plus two times the sum of variances of pure error and measurement error. Therefore, the distance component of MSD can be estimated by subtracting twice the components of pure error and measurement error estimated by Equation 9 in the previous section:

$$\sum_i [g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n \approx \text{MSD} - \text{MMSD}(0). \quad (10)$$

This error component does not depend on the specific functional form of the relations of the \mathbf{Y} variables to the \mathbf{X} variables, so any model lack-of-fit is not involved. Therefore, it applies equally to near-neighbor pairing of units without regard for the distance function. Unfortunately, $\text{MSD} - \text{MMSD}(0)$ is not constrained to be positive if $[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2$ decreases with increasing distance.

Using the Partitioning to Illuminate Key Questions

We now revisit key questions posed in the introduction, developing some new statistics based on the partitioning to provide answers.

Accuracy of Imputed Values

The fundamental variance statistic in sampling inference compares an estimate with its true value. In our notation that comparison is $y_{Rj} - g(\mathbf{X}_j)$ for k equal to one. Therefore, we propose that the efficacy of the imputation process should be based on a statistic we term the standard error of imputation (SEI).

$$\text{SEI}^2 = \sum_i [y_{Rj} - g(\mathbf{X}_{Ti})]^2/n$$

$$i = 1, \dots, n \text{ and } j \text{ minimizes } d_{ij}^2 \quad (11)$$

Unfortunately, the addends in the bracket of SEI cannot be computed directly from the data in the reference set because the true value, $g(\mathbf{X}_{Ti})$, is not directly observable. The proposed aggregate statistic (Equation 11), however, can be

obtained by replacing the “estimate” y_{Rj} in Equation 11 with Equation 4 evaluated for the j th reference unit.

$$\text{SEI}^2 = \sum_i [g(\mathbf{X}_{Rj}) + \varepsilon_{pj} + \varepsilon_{Yj} - g(\mathbf{X}_{Ti})]^2/n. \quad (12)$$

Then averaging with the same assumptions of error independence used in deriving Equation 8:

$$E[\text{SEI}^2] = E[\sum_i [y_{Rj} - g(\mathbf{X}_{Ti})]^2/n]$$

$$= \sum_i [g(\mathbf{X}_{Rj}) - g(\mathbf{X}_{Ti})]^2/n + \text{Var}(\varepsilon_p) + \text{Var}(\varepsilon_Y), \quad (13)$$

which differs from MSD (Equation 8) by omitting the terms for the variances of pure error and sampling error arising from the target members of Equation 11. If the distance component of MMSD(0) can be assumed to be trivially small when Equation 8 is averaged over only the shorter distances, then

$$E[\text{SEI}^2] = E[y_{Rj} - g(\mathbf{X}_{Ti})]^2 \approx \text{MSD} - \text{MMSD}(0)/2. \quad (14)$$

Imputation Compared to Estimates Using $f(\mathbf{X})$

The regression model is $y_j^* = f(\mathbf{X}_j) + \varepsilon_j$, where the ε_j includes pure error and the lack of fit of the assumed model. The regression model could be, but is not limited to, the familiar linear parameterization $f(\mathbf{X}_j) = \mathbf{B}\mathbf{X}'$. Error lack-of-fit for the linear model is indicated in Figure 1. Alternatively, it could be a nonlinear or nonparametric regression model or a collection of means for strata defined by the \mathbf{X} variables. The true model $y_j^* = g(\mathbf{X}_j) + \varepsilon_{pj}$ differs from the regression model by the lack-of-fit of the regression model:

$$\varepsilon_{L(Xj)} = g(\mathbf{X}_j) - f(\mathbf{X}_j). \quad (15)$$

The error statistic commonly calculated for a regression is the standard error of estimate (SEE) (ignoring the reduction of the divisor by the number of estimated parameters):

$$\text{SEE}^2 = \sum_j (y_j - f(\mathbf{X}_j))^2/n. \quad (16)$$

We assume that the lack-of-fit will sum to zero for the particular \mathbf{X} values (certain if $f(\mathbf{X})$ is fit by least-squares and includes an intercept) in the reference set.

Then, from Equations 2, 3, and 15,

$$(y_j - f(\mathbf{X}_j))^2 = (\varepsilon_{pj} + \varepsilon_{Yj} + \varepsilon_{L(Xj)})^2. \quad (17)$$

The terms for the model lack-of-fit were assumed to be independent of the \mathbf{X} values and of ε_{pj} and ε_{Yj} , so $E(\text{SEE}^2)$ is the sum of these three sources:

$$E[\text{SEE}^2] = E[\sum_j (y_j - f(\mathbf{X}_j))^2/n]$$

$$= \text{Var}(\varepsilon_p) + \text{Var}(\varepsilon_Y) + \sum_j [\varepsilon_{L(Xj)}^2]/n. \quad (18)$$

Comparison of $E[\text{SEI}^2]$ in Equation 13 with $E[\text{SEE}^2]$ in Equation 18 shows that they differ only by the substitution of the distance component, $\sum_i [g(\mathbf{X}_{Rj}) - g(\mathbf{X}_{Ti})]^2/n$, in imputation error variance for lack of fit, $\sum_j (\varepsilon_{L(Xj)})^2/n$, in regression estimation error variance.

Rearranging Equation 18 and substituting Equation 9,

$$\sum_j [\varepsilon_{L(Xj)}^2]/n = E[\text{SEE}^2] - E[\text{MMSD}(0)/2] \quad (19)$$

The ideal contents for a data set for subsequent analysis would be \mathbf{Y}_j^* , which would have variance about $g(\mathbf{X}_j)$ of

$\text{Var}(\varepsilon_p)$. Unfortunately, the best imputation can do for a given data set is \mathbf{Y}_{Tj} , which differs from the ideal by inclusion of measurement error variance plus the distance component. Alternatively, regression estimation could supply as estimates, $f(\mathbf{X}_j)$ plus a random element drawn from a distribution with variance $\text{Var}(\varepsilon_p)$. Using Equations 4 and 15 and the independence of pure error relative to the model lack-of-fit, these estimates would have variance about $g(\mathbf{X}_j)$ given by

$$\begin{aligned} E[\Sigma_j[f(\mathbf{X}_j) + \varepsilon_{pi} - g(\mathbf{X}_j)]^2/n] \\ = \Sigma_j(\varepsilon_{L(\mathbf{X}_j)})^2/n + \text{Var}(\varepsilon_p) \\ = E[\text{SEE}^2] - \text{Var}(\varepsilon_Y), \end{aligned} \quad (20)$$

which can be estimated by

$$\begin{aligned} \Sigma_j[f(\mathbf{X}_j) + \varepsilon_{pj} - g(\mathbf{X}_j)]^2/n \\ = \text{SEE}^2 - \text{MMSD}(0)/2 + \text{var}(\varepsilon_p) \\ = \text{SEE}^2 - \text{var}(\varepsilon_Y). \end{aligned} \quad (21)$$

Subtracting Equation 21 from Equation 14, the comparison of SEI^2 to Equation 21 becomes

$$\begin{aligned} E[y_{Rj} - g(\mathbf{X}_{Ti})]^2 - E[\Sigma_j[f(\mathbf{X}_j) + \varepsilon_{pi} - g(\mathbf{X}_j)]^2/n] \\ \approx \text{SEI}^2 - [\text{SEE}^2 - \text{var}(\varepsilon_Y)], \end{aligned} \quad (22)$$

which is the same as Equation 13 minus Equation 18 plus pure error variance:

$$\Sigma_i[g(\mathbf{X}_{Rj}) - g(\mathbf{X}_{Ti})]^2/n - [\Sigma_j(\varepsilon_{L(\mathbf{X}_j)})^2/n + \text{var}(\varepsilon_p)]. \quad (23)$$

Thus, the variance of the imputed values would be greater than regression estimated values for each y if Equation 22 or equivalently if Equation 23 is greater than zero. However, the regression alternative would not guarantee that the true correlation among the estimated y values within each observation unit would be retained.

Effects of Distribution of \mathbf{X} Values

The second key question concerning distributions of the \mathbf{X} values and alternative measures of similarity is addressed by considering the distance component of MSD, $\Sigma_i[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n$. This error component should be made as small as possible, either by adding new members to the reference set to reduce average distance between target units and their similar reference unit(s) or by adopting a better measure of similarity, or both.

An important consideration in accuracy assessment based only on the reference observation units is the relation between the distribution of the \mathbf{X} values in the target set in relation to that distribution in the reference set. Ideally, the reference set would completely cover the ranges of \mathbf{X} variables of the target set and have an approximately uniform distribution over the range of the combined sets. The distance function being invoked may weight variation of some of the \mathbf{X} variables heavier than others, thereby stretching and rotating the space spanned by the \mathbf{X} variables. Therefore, the distributions to be compared are those of the distances between the reference units and their near-neighbors from the pseudo-target set versus the distances between

the reference units and their near-neighbors from the real target set.

A statistic sensitive to the merits of alternative distance functions would reduce the influence of pure error and sampling error to focus on $\Sigma_i[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2$. At short distances, the values of $(y_{Ti} - y_{Rj})^2$ are dominated by the pure error plus sampling error. Therefore, a better alternative to MSD calculated as the average overall references is to average by only using pairs separated by the longer distances.

Choice of \mathbf{X} Values and their Transformations

How these decisions affect MSD for a particular variable y depends on the choice of the weight matrix \mathbf{W} in Equation 1. If \mathbf{W} gives little or no weight to a particular x , then that x is effectively omitted. Conversely, an x may be heavily weighted because of its contribution to $g(\mathbf{X})$ for other y values. Then, even though a subset of the x variables may effectively predict the y under consideration, their contribution will be diluted by differences in the extraneous x variables and MSD for that element, y of Y will be dominated by pure error and measurement error to such an extent that $[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2$ may decrease with distance. If it does decrease, then the distance component and model lack-of-fit will be underestimated.

Transformations in variables are typically invoked to simplify a model such as $y = f(\mathbf{X})$ and to render errors more homogeneous. Consideration of Equations 8, 10, and 17, as estimates of sources of imputation errors from the three sources shows that transformations of the \mathbf{X} variables, while modifying the fit of the regression model $y = f(\mathbf{X})$, affect MSD only through the distance component, $\Sigma_i[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj})]^2/n$, and homogeneity of the pure error component. Transformations affect the distance component through the selection of surrogates, which in turn depend on the choice of the weight matrix \mathbf{W} . In dense regions of the space spanned by the \mathbf{X} variables of the reference set, the distance component in MSD is small relative to pure error plus measurement error for any choice of near neighbor. However, where the \mathbf{X}_{Ti} are not closely spaced (sparse), their imputations to the \mathbf{X}_{Tj} will be few in number, so their effect on MSD will be small. This ambiguity explains a puzzling property of near-neighbor imputation: that it has not appeared to be very sensitive to monotonic transformations of the variables. However, for imputation methods that base \mathbf{W} on the relations of \mathbf{Y} to \mathbf{X} in distance calculations (e.g., MSN), the nonlinear components represented by lack-of-fit would change the selection of “near neighbors.” The extent of the change would be greatest in pairs of observation units in which model lack-of-fits were of opposite sign.

Choice of k

The partitioning of error provides useful insight concerning the choice of k for imputation using a weighted average of k near neighbors. The obvious effect is that larger k , by averaging over the errors of more reference observations, would seem to reduce the error of the imputed value.

However, it is not that simple. Following the same assumptions used in deriving Equation 8 MSD becomes

$$E[\Sigma_i[y_{Ti} - \Sigma_m w_{im} y_{Rm}]^2/n] \\ = \Sigma_i[g(\mathbf{X}_{Ti}) - \Sigma_m w_{im} g(\mathbf{X}_{Rm})]^2/n \\ + (1 + \Sigma_i \Sigma_m w_{im}^2/n)[\text{Var}(\varepsilon_Y) + \text{Var}(\varepsilon_P)]. \quad (24)$$

In k -nn imputation, y_{Rj} of Equation 8 is replaced by an average of k values of y_m using a weighting rule for the particular flavor of k -nn inference, where m is from the set of indices of the k observations selected as near neighbors to the i th target and $\Sigma_m w_{im} = 1$. When $w_{im} = 1/k$, the multiplier of the variances in Equation 24 becomes $(1 + 1/k)$. To the extent that it is pure error being reduced, increasing k is counterproductive for the subsequent analysis. Offsetting this effect, measurement error will also be reduced in the same proportion. Hence, there is a tradeoff; either lose valuable pure error or reduce undesirable measurement error. The net effect of changing k also depends on the change in $\Sigma_i[g(\mathbf{X}_{Ti}) - \Sigma_m w_{im} g(\mathbf{X}_{Rm})]^2/n$. Whether this component increases or decreases the total error depends on the change of $[g(\mathbf{X}_{Ti}) - \Sigma_m w_{im} g(\mathbf{X}_{Rm})]^2$ for the reference observation being added or omitted by changing k .

Application to Example Data Sets

Three data sets will be used to illustrate the estimation of error components and application of these estimates in evaluating alternative weight matrices. All three use suites of remotely sensed data and data from digital terrain models to impute data from ground-based observations. As examples of real imputation analyses, they illustrate the behavior of the statistics we propose. We do not purport to second-guess the analysis of these data sets, so the definitions of the 69 specific variables in these three data sets are mostly irrelevant to our purposes. Where we do discuss behavior of the partitioning as a consequence of the biological situation, we will define those variables explicitly in the text. Otherwise, readers desiring more detail are directed to the original sources.

The first example uses data used by Moisen and Frescino (2002) obtained by the USDA Forest Service, Rocky Mountain Experiment Station Forest Inventory and Analysis Unit (FIA). The ground-based data (\mathbf{Y} variables) are from routine FIA observations for Utah. The \mathbf{X} variables were obtained from LANDSAT and digital terrain data.

The other two data sets use ground data from inventories of stands defined as polygons. One, from the Deschutes National Forest in Oregon, has been used in previously reported analyses by Moeur (2000) and is the example in the MSN User's Guide (Crookston et al. 2002). The third data set is from Tally Lake area in the Flathead National Forest in Montana. For these comparisons, the \mathbf{Y} variables will be limited to those measured on continuous scales. These analyses differ from those reported by Stage and Crookston (2002) in that all discrete and a few redundant y variables have been omitted to achieve approximately equal numbers of y variables in the three examples, and additional x variables (transformations of the original variables) have been

Table 1. Number of coefficients to be estimated in relation to number of samples for three data sets used as examples

	Tally Lake	Users Guide	Utah
Number of Y variables	8	6	10
Number of X variables (p)	21	12	12
Number of reference obs. (n)	847	197	1076
Significant canonical pairs (s)	7	5	4
$n/(s + p * s)$	5.50	3.03	16.55

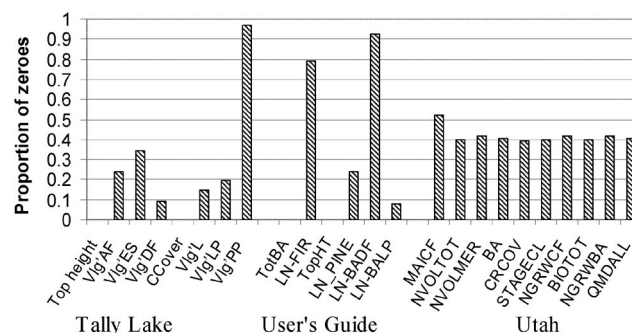


Figure 2. Proportion of zero values in example data sets.

added. Table 1 summarizes numbers of variables and sample sizes for the three data sets. Of the three data sets, the Users Guide has remarkably fewer observations in relation to the number of unique coefficients in the weight matrix being estimated (last line, Table 1).

The Utah data set differs from the other two in that it contains a notable portion of locations in nonforest, although the continuous \mathbf{Y} variables describe forest stand parameters. By contrast, y values of zero in the other two data sets indicate lack of stocking in otherwise forested polygons. The proportion of zeroes in the three data sets are indicated in Figure 2.

Table 2 summarizes the structure of the correlations between the canonical vectors for the three data sets. Multivariate regression R^2 values of y on \mathbf{X} are listed in col. B of Table 2. Correlations between \mathbf{Y} and \mathbf{X} variables were lowest in the Utah data because the measurement errors of the \mathbf{Y} variables from the FIA plot clusters were larger than in the two data sets based on inventories of stand polygons.

Components of Variance

Data for partitioning variance for the three example data sets are displayed in Table 3. Columns A–C contain statistics for each y variable considered independently of the remaining elements of \mathbf{Y} . Columns D–F contain statistics

Table 2. Comparison between three example data sets of first four squared canonical correlations between \mathbf{Y} and \mathbf{X}

Canonical Pair (m)	Tally Lake	User's Guide	Utah
1	0.697	0.686	0.450
2	0.477	0.456	0.153
3	0.325	0.376	0.109
4	0.292	0.244	0.034

Table 3. Components of variance for three example data sets

Y Variable	Total variance of Y variable in reference set (A)	Multivariate regression R^2 (B)	Squared error about regression of single Y SEE^2 1.-col. B (C)	Mean square between target and nearest reference for all pairs in MSD (D)	Mean square between target and nearest reference for 1/8 of shorter distances (MMSD(0)) (E)	Calculated distance component $D - E$ (F)
Tally Lake						
Top height	566.669	0.6713	0.3287	0.6990	0.2837	0.4153
Vlg'AF	8.69601	0.4716	0.5284	1.0380	0.8461	0.1919
Vlg'ES	9.01080	0.4322	0.5678	1.2098	1.4075	-0.1977
Vlg'DF	7.03682	0.3696	0.6304	1.0064	0.6292	0.3772
CCover	222.797	0.2999	0.7001	1.1628	1.0061	0.1567
Vlg'L	6.66466	0.2556	0.7444	1.3189	0.9271	0.3918
Vlg'LP	8.18933	0.1956	0.8044	1.4893	1.0475	0.4418
Vlg'PP	0.71893	0.1076	0.8924	1.1779	0.6486	0.5294
Users Guide						
TotBA	2822.19	0.5917	0.4083	0.7695	0.735	0.0345
LN-FIR	4.43945	0.5440	0.4560	0.8217	0.087	0.7347
TopHT	292.968	0.4839	0.5161	0.9453	0.287	0.6583
LN_PINE	7.0639	0.3858	0.6142	1.0768	0.087	0.9898
LN-BADF	1.85368	0.3548	0.6452	0.9557	0	0.9557
LN-BALP	3.55926	0.3225	0.6775	1.2927	0.6712	0.6215
Utah						
MAICF	684.346	0.3567	0.6433	1.1259	0.3334	0.7925
NVOLTOT	2064882.	0.3142	0.6858	1.3522	0.7868	0.5654
NVOLMER	1546287.	0.2976	0.7024	1.3472	0.8143	0.5329
BA	4211.15	0.2736	0.07264	1.3271	0.7525	0.5746
CRCOV	779.175	0.2621	0.7379	1.4426	0.8551	0.5876
STAGECL	3746.17	0.2528	0.7472	1.3367	1.0754	0.2613
NGRWCF	905.920	0.2434	0.7566	1.4868	2.1123	-0.6255
BIOTOT	636.018	0.2390	0.07610	1.4758	0.5886	0.8872
NGRWBA	0.75238	0.2280	0.7720	1.5302	1.0740	0.4562
QMDALL	19.5868	0.1711	0.8289	1.6427	0.6462	0.9965

Columns C–F are standardized by division by variance in column A. Columns B and C are for a linear model used as $y = f(\mathbf{X})$. Columns D–F are obtained with a Mahalanobis distance function.

for each y variable for pairs of near neighbors selected using a multivariate Mahalanobis distance measure.

Accuracy of Imputed Values

Standard error of imputation squared (SEI^2) (as a fraction of variance of each variable) of values imputed using a Mahalanobis distance function are shown in Figure 3. The error component arising from distance between target and reference, $\sum_i [(g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Rj}))^2]/n$, as estimated by Equation 10, is shown in Figure 3 by the shaded portions of the bars for each y variable. This figure also shows the combined components of pure error and measurement error as estimated by Equation 9.

Imputation Compared to Linear Regression

Figure 4 compares the distance component of imputations (plotted as its negative) with the model lack-of-fit computed as $SEE^2 - \text{Min}(\text{MMSD}(0)/2, SEE^2)$ for $f(\mathbf{X}) = \beta\mathbf{X}'$. As a corollary of the differences between imputation distance component and regression lack-of-fit, SEE^2 is almost always less than SEI^2 . The exceptions to the inequality

are Crown Cover (CCover) and logarithm of *Pinus ponderosa* volume (Vlg'PP) in the Tally Lake data set. We conjecture that the linear regression is just not a very effective model for crown cover, and that the large proportion of zero data for *Pinus ponderosa* preclude effective prediction of volume. Also, there would be two anomalies leading to negative estimates of lack-of-fit if the minimum of MMSD(0) and SEE^2 were not used: logarithm of Engelmann spruce volume (Vlg'ES) in the Tally Lake data and net growth in cubic feet (NGRWCF) in the Utah data. The larger values of MMSD(0) for these variables are the consequence of squared differences between y_{Ti} and y_{Rj} that decrease with increasing differences in the \mathbf{X} variables. As a result, MMSD(0) is larger than SEE^2 . We attribute this anomaly to unequal pure error in different regions of the \mathbf{X} space. Engelmann spruce in the Tally Lake area occurs bimodally with elevation—either very common at high elevations or as sparse stringers in valley bottoms. However, the density in the \mathbf{X} space of the observations representing valley bottoms and stands at similar elevations is higher than the density of data representing high elevations. Thus observation pairs with near-zero distances tend to come from low elevations, where the sporadic presence of spruce gives large squared differences, whereas at high

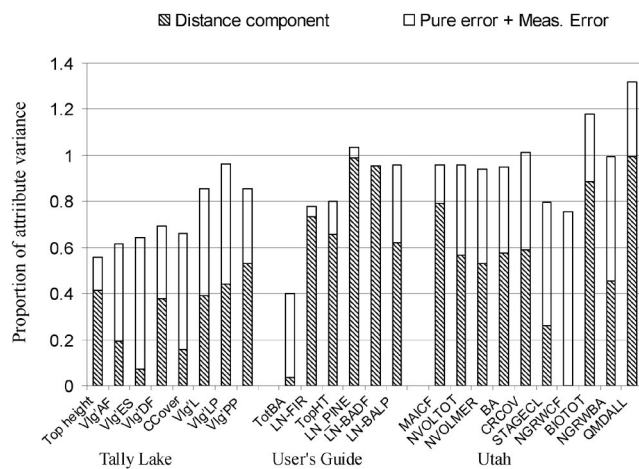


Figure 3. Partitioning of relative variance of imputed values (SEI Equation 13) for Mahalanobis distance function. Variables within a data set are ordered from left to right by increasing SEE. Values standardized by division by attribute variance.

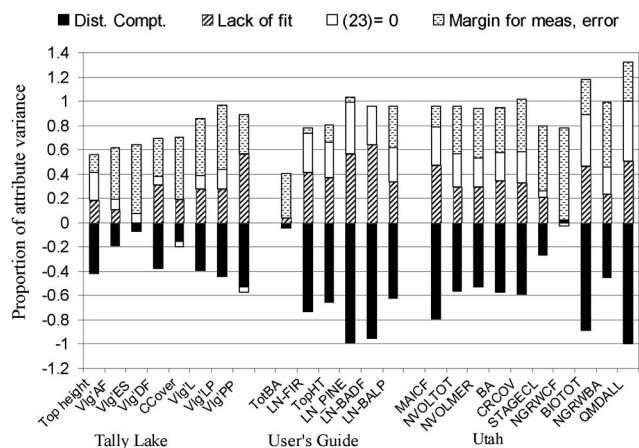


Figure 4. Distance error component of imputation (plotted as negative) compared to lack of fit of a linear regression, and pure error plus measurement error. The clear portion of the bar is the amount of error that would be added to lack of fit to make expression 23 equal zero. The stippled bar is the remaining portion of pure error plus measurement error. Variables within a data set are ordered from left to right by increasing SEE for a linear regression model. Values standardized by division by attribute variance.

elevations spruce is more ubiquitous, giving smaller differences in volume even at larger separations in \mathbf{X} space.

That SEE is almost always less than SEI is not surprising, because whereas SEE is a least-squares minimization of the model prediction, SEI is not the result of an explicit minimization and includes the pure error and measurement error components. When pure error should be included in estimates for subsequent analyses, the proportion of pure error that might be added to regression lack of fit that would just make Equation 23 equal to zero is indicated by the white bars in Figure 4. Unfortunately, we lack a direct estimate of measurement error that should be subtracted from SEI, so we can only show the margin from which it would be subtracted.

Effect of Distances between \mathbf{X} Values

The three data sets show differences in the proportions of variance attributable to the Mahalanobis distances between

target and reference (Figure 3, shaded bar). The low ratio of number of observations compared to number of coefficients to be estimated and large linear model lack-of-fit of the User's Guide data produces a relatively large distance component compared to the Tally Lake data. Utah data show an intermediate level because the effect of the larger number of data relative to the number of coefficients to be estimated is offset by the low correlations between \mathbf{Y} and \mathbf{X} values (Table 3) caused by the inclusion of nonforest observations (Figure 2).

In the Tally Lake application, average distances from reference observation units to actual target observation units is 2.04 times the average distance from each reference observation unit to its nearest neighbor also in the reference set. Nearly one-third of the targets are farther from their nearest reference than the ninth percentile of the distribution of distances among the references. The significance of this extrapolation might be determined by modeling squared differences for each element of \mathbf{Y} as a function of distance. Such analysis is beyond the scope of this report.

Comparison of Alternative Distance Functions

The difficulty of using MSD to compare alternative distance functions can be appreciated by considering that the influence of pure error plus sampling error would be double that shown in Figure 3. Although the absolute value of differences in MSD arising from different distance functions would not change, the relative importance of the differences among the alternative distance functions would be underestimated.

Figures 5, 6, and 7 compare three alternative distance functions, the Mahalanobis distance used heretofore in this

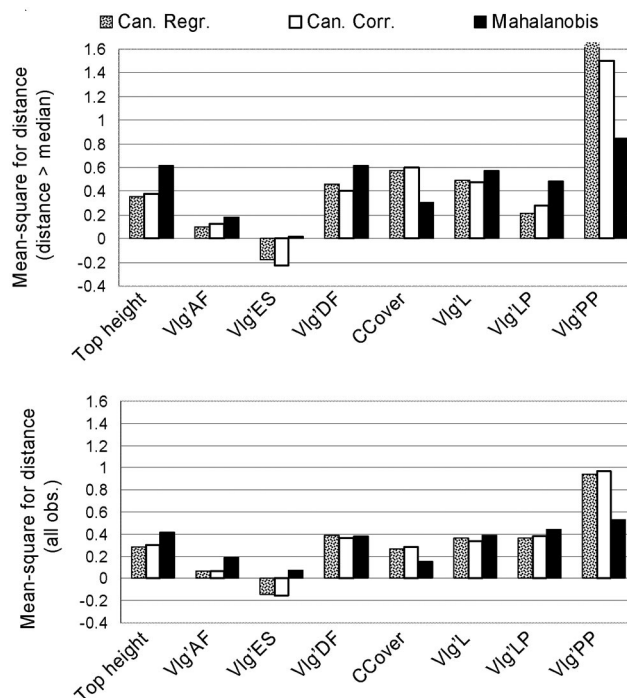


Figure 5. Tally Lake comparison of distance components (Equation 10) for two canonical-correlation-based distance functions with Mahalanobis distance function. Variables within a data set are ordered from left to right by increasing SEE.

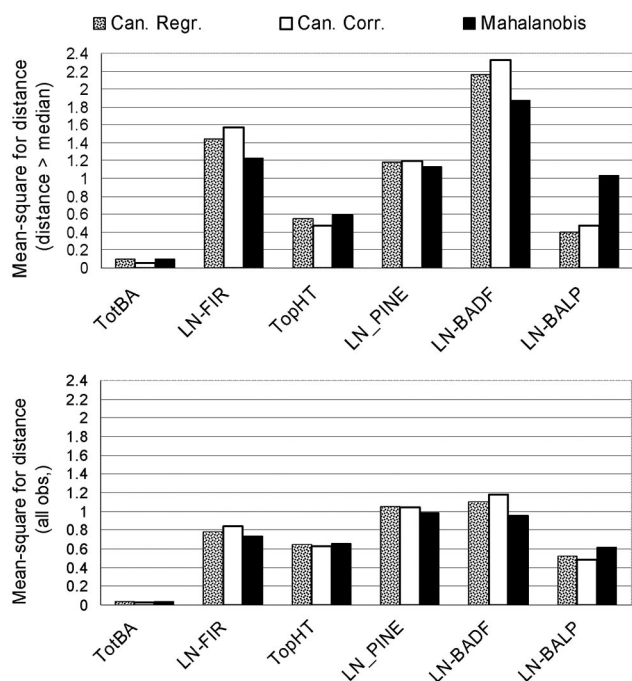


Figure 6. Users guide. Comparison of distance components (Equation 10) for two canonical-correlation-based distance functions with Mahalanobis distance function. Variables within a data set are ordered from left to right by increasing SEE.

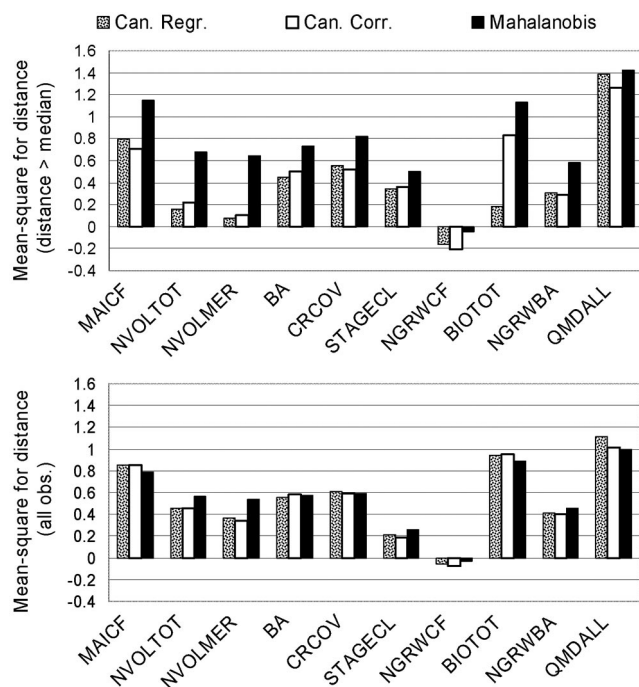


Figure 7. Utah. Comparison of distance components (Equation 10) for two canonical-correlation-based distance functions with Mahalanobis distance function. Variables within a data set are ordered from left to right by increasing SEE.

report, the original canonical-correlation-based distance (CC) of Moeur and Stage (1995), and the newer canonical-regression-based distance (CR) introduced by Stage and Crookston (2002). The panels present both estimated means of $[g(\mathbf{X}_{Ti}) - g(\mathbf{X}_{Ri})]^2$ based on all data for comparison to means for the 50% of the data separated by the longer

distances. Only the Utah data show the alternative similarity measures to rank differently in the full data set than in the reduced data set containing only the 50% longer distances. Also, the Utah data set was the only one to show a distinct advantage to using one or the other of the canonical-based distances over the Mahalanobis distances; the differences would be even greater if the nonforest data were masked because the Mahalanobis distances did slightly better at matching the zero data. The result seems anomalous because the Utah data had the lowest canonical correlations between \mathbf{Y} and \mathbf{X} . However, one of the merits of the canonical approach lies in its capability to ignore \mathbf{X} variables that are irrelevant. Moisen and Frescino (2002) found that several of the \mathbf{X} variables were superfluous. The Mahalanobis distance would have given these variables weights equal to the weights of the useful variables. The other two data sets were obtained after extensive analysis by others that probably had already screened the \mathbf{X} variables for utility.

Conclusions

This report concerns the error properties of imputation processes used to fill in a data set by imputing values from a sample of intensively measured observation units to interspersed, less completely measured units. The error statistics for the imputed, continuous-valued variables presented in this report are based on partitioning of the error components into measurement error, error inherent in the particular imputation method, and the pure error not associated with the variables measured on all observation units. These statistics can assist in the design of inventories and their analysis with near-neighbor imputation methods. It is now possible to consider the relative gains from reducing measurement error versus increasing the density of the sampled observation units. They also clarify comparisons to other inference methods such as regression or stratum-mean based estimators, and help to choose among alternative weight matrices in similarity measures.

Literature Cited

- CROOKSTON, N.L., M. MOEUR, AND D. RENNER. 2002. *Users guide to the most similar neighbor imputation program, version 2*. Gen. Tech. Rep. RMRS-GTR-96. US Department of Agriculture, Forest Service, Rocky Mountain Research Station, Ogden, UT. 35 p.
- MALINEN, J. 2003. Locally adaptable non-parametric methods for estimating stand characteristics for wood procurement planning. *Silva Fennica* 37(1):109–120.
- MOEUR, M. 2000. Extending stand exam data with most similar neighbor inference. P. 99–107 in *Proc. Soc. of Amer. Foresters National Convention*; Sept. 11–15, 1999. SAF Pub 00-1.
- MOEUR, M., AND A.R. STAGE. 1995. Most similar neighbor: An improved sampling inference procedure for natural resource planning. *For. Sci.* 41:337–359.
- MOISEN, G.G., AND T.S. FRESCINO. 2002. Comparing five modeling techniques for predicting forest characteristics. *Ecol. Modelling* 157:209–225.
- MULLIN, M., AND R. SUKTHANKAR. 2000. Complete cross-validation for nearest neighbor classifiers. P. 639–646 in *Proc. 17th International Conf. on Machine Learning* (June 29–July 2,

- 2000). P. Langley, (ed.). Morgan Kaufmann Publishers, San Francisco, CA.
- OHMANN, J.L., AND M.J. GREGORY. 2002. Predictive mapping of forest composition and structure with direct gradient analysis and nearest neighbor imputation in coastal Oregon, U.S.A. *Can. J. For. Res.* 32:725–741.
- PODANI, J. 2000. *Introduction to the exploration of multivariate biological data*. Backhuys, Leiden, The Netherlands.
- SHAO, J., AND R.R. SITTER. 1996. Bootstrap for imputed survey data. *J. Amer. Stat. Assoc.* 91(435):1278–1288.
- SOHN, Y., E. MORAN, AND F. GURRI. 1999. Deforestation in north-central Yucatan (1985–1995): Mapping secondary succession of forest and agricultural land use in Sotuta using the cosine of the angle concept. *Photogram. Eng. Remote Sensing* 65(8):947–958.
- STAGE, A.R., AND N.L. CROOKSTON. 2002. Measuring similarity in nearest neighbor imputation: Some new alternatives. P. 91–96. In *Symposium on statistics and information technology in forestry*, September 8–12, 2002. Virginia Polytechnic Institute and State University, Blacksburg, VA. Available online at www.forestry.ubc.ca/prognosis/documents/MSN_StageCrookston.pdf. Last accessed Jan. 4, 2007.