**ELSEVIER**

# Stratified estimation of forest area using satellite imagery, inventory data, and the *k*-Nearest Neighbors technique

Ronald E. McRoberts *, Mark D. Nelson, Daniel G. Wendt

*Forest Inventory and Analysis, North Central Research Station, USDA Forest Service, 1992 Folwell Avenue, St. Paul, MN 55108, USA*

## Abstract

For two large study areas in Minnesota, USA, stratified estimation using classified Landsat Thematic Mapper satellite imagery as the basis for stratification was used to estimate forest area. Measurements of forest inventory plots obtained for a 12-month period in 1998 and 1999 were used as the source of data for within-stratum estimates. These measurements further served as calibration data for a *k*-Nearest Neighbors technique that was used to predict forest land proportion for image pixels. The continuum of forest land proportion predictions was separated into strata to facilitate stratified estimation. The *k*-Nearest Neighbors technique is carefully explained, five precautions are noted, and a plea is made for an objective approach to calibrating the technique. The variances of the stratified forest area estimates were smaller by factors as great as 5 than variances of the arithmetic mean calculated under the assumption of simple random sampling. In addition, when including all plots over a 5-year plot measurement cycle, the forest area precision estimates may be expected to satisfy national standards.
© 2002 Elsevier Science Inc. All rights reserved.

## 1. Introduction

The five regional, multi-state Forest Inventory and Analysis (FIA) programs of the Forest Service, US Department of Agriculture, are required to report estimates of forest land area for their respective regions every 5 years. Each estimate is obtained as the product of total area inventoried and the mean, over a systematic array of field plots, of the proportion of each plot in FIA-defined forest land. The FIA definition of forest land includes commercial timberland, some pastured land with trees, forest plantations, unproductive forested land, and reserved, noncommercial forested land. In addition, forest land must satisfy minimum stocking levels, a 0.405-ha (1-ac) minimum area, and a minimum continuous bole-to-bole canopy width of 36.58 m (120 ft), therefore excluding lands such as wooded strips, idle farmland with trees, and narrow windbreaks. A combination of budgetary constraints and natural variability among plots prohibits

sample sizes sufficient to satisfy national FIA precision standards for forest area estimates unless the estimation process is enhanced using ancillary data.

Traditionally, FIA has enhanced the estimation process by using stratified estimation with aerial photography as the basis for stratification (Bickford, 1960; Hansen, 1990; Loetsch & Haller, 1964). First, an extensive array of photo plots on aerial photographs is interpreted and stratified using ocular methods, and the proportions of photo plots assigned to strata are used as stratum weights. Then, field crews visit a subset of the photo plots and observe plot attributes. Finally, estimates of forest land area are obtained with these data using stratified estimation techniques (Cochran, 1977).

A second approach to enhancing the estimation process is to use stratified estimation with classified satellite imagery as the basis for the stratification. With this approach, image pixels for the area of interest are classified with respect to predictions of land cover attributes into homogeneous classes, and the classes are then used as strata in the stratified analyses. Strata weights are the proportions of pixels in strata, and plots are assigned to strata on the basis of the strata assign-

---

* Corresponding author. Tel.: +1-651-649-5174; fax: +1-651-649-5285.
*E-mail address:* rmcroberts@fs.fed.us (R.E. McRoberts).

ments of their associated pixels. If the stratification is done before sampling and the within-stratum variances of the inventory variables are well estimated, then maximum precision may be achieved by selecting within-strata sampling intensities to be proportional to within-strata variances. However, even when the within-strata sampling intensities are independent of the stratification, stratified estimation may still yield increases in precision.

Satellite imagery has been used as a basis for stratification for variance reduction in forestry applications. Poso, Hame, and Paananen (1984) and Poso, Paananen, and Simila (1987) used Landsat Thematic Mapper (TM) imagery to obtain stratified estimates of volume and age in Finland, and Deppe (1998) used satellite imagery to stratify for estimating forest area in Brazil. For estimating forest land area, Hansen and Wendt (2000) used the GAP classification (Scott et al., 1993) to increase the precision of inventory estimates for Indiana and Illinois, USA; Hoppus, Arner, and Lister (2001) investigated the utility of both GAP and the National Land Cover Dataset (NLCD) (Vogelmann et al., 2001) for the same purpose for Connecticut, USA; and McRoberts, Wendt, Nelson, and Hansen (2002) investigated methods for optimizing the stratification utility of the NLCD for estimating forest area for four states in the North Central region of the USA. Both GAP and the NLCD are land cover classifications based on nominal 1992 TM imagery and ancillary data. Neither classification was designed for forest inventory estimation; neither has been demonstrated to produce stratifications that satisfy FIA precision standards; and neither has been replaced yet with a more current version. The Multi-Resolution Land Characterization Consortium (MRLC) (Loveland & Shaw, 1996), which produced the NLCD, plans to produce such classifications at 10-year intervals, but because release is not expected until 5 years after the date of the imagery, the classifications will always be 5–15 years out of date.

The timeliness of FIA estimates is enhanced when cycles for obtaining and classifying imagery are comparable to the 5-, 7-, or 10-year plot measurement cycles, depending on region. Because externally produced classifications such as GAP and NLCD apparently will not be produced on such cycles, the regional FIA programs may find it necessary to produce their own classifications. On average, a regional FIA program on a 5-year plot measurement cycle will need to classify images for approximately 125 TM scenes over the cycle. In addition, sufficient training data to guide the classifications must be obtained in close temporal proximity to the imagery dates. These are important tasks that merit FIA investigation of efficient means of obtaining training data and processing images. The objective of this study is to investigate the utility of the k-Nearest Neighbors technique in processing TM imagery for use as the basis for enhancing forest area estimates through stratification.

## 2. Data

### 2.1. Study areas

The study was conducted in two areas in Minnesota, USA, designated St. Louis and St. Cloud (Fig. 1). The St. Louis study area encompasses most of St. Louis County, includes approximately 2.1 million hectares of which approximately 75% is forest land and is dominated by Aspen–Birch and Spruce–Fir associations. The St. Cloud study area contains the St. Cloud urban area, includes approximately 3.3 million hectares of which approximately 20% is forest land and is characterized by prairie agriculture and a diverse mixture of forest lands including both coniferous and deciduous species.

### 2.2. Satellite imagery

The St. Louis study area is covered by the Landsat TM Path 27, Row 27 scene and includes all of St. Louis County except the northern portion. For this scene, Landsat-7 ETM+ images were obtained for two seasons: autumn (5 November 1999) and spring (31 May 2000). The St. Cloud study area is covered by the Landsat TM Path 28, Row 28 scene. For this scene, Landsat-7 ETM+ images were obtained for three seasons: summer (23 July 1999), autumn (27 October 1999), and winter (3 March 2000). The following attributes pertain to all five images: (1) $30 \times 30$ m pixels from bands 1 to 5 and band 7, (2) absolute radiance units scaled to 8 bits, (3) processing to level 1G (processing level 08; radiometrically and geometrically corrected using satellite model and platform/ephemeris information), and (4) geo-referencing to Albers Equal Area projection, NAD83. In addition, for the St. Louis study area, the November image was rectified using 40 ground control points with resulting root mean square error of 12.1 m. The May image was registered to the November image using 26 ground
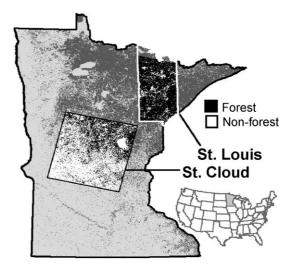


Fig. 1. Minnesota study areas.

control points and resampled using first-order polynomial and nearest neighbor techniques with resulting root mean square error of 31.9 m. For the St. Cloud study area, all three images were rectified using ground control points and digital elevation model terrain correction (processing level 10) and resampled using cubic convolution with resulting root mean square error of less than 8.5 m. Bands are distinguished using an alphanumeric character representing the first letter of the month of the image and a numeric character designating the band. The context of band references indicates whether they refer to St. Louis or St. Cloud images.

### 2.3. FIA plot data

Under the FIA program's annual inventory system (McRoberts, 1999), field plots are established in permanent locations using a systematic sampling design. In each state, a fixed proportion of plots are measured annually; plots measured in a single federal fiscal year (e.g., FY1999: 1 October 1998 to 30 September 1999) make up a single panel of plots, and panels are selected for annual measurement on a rotating basis. In aggregate, over a complete measurement cycle, a plot represents 2403 ha (slightly less than 6000 ac). In general, locations of forested or previously forested plots are determined using global positioning system receivers, while locations of nonforested plots are determined using digitization methods.

Each field plot consists of four 7.31-m (24-ft) radius circular subplots. The subplots are configured as a central subplot and three peripheral subplots with centers located at 36.58 m (120 ft) and azimuths of 0°, 120°, and 240° from the center of the central subplot. Among the observations field crews obtained are the proportions of subplot areas that satisfy specific ground land use conditions. Subplot-level estimates of forest land proportion are obtained by aggregating these ground land use conditions consistent with the FIA definition of forest land, and plot-level estimates are obtained as means over the four subplots.

For both study areas, measurements for the FY1999 panel of inventory plots were available. For the St. Louis study area, measurements for 133 plots or 532 subplots were used of which 387 subplots were completely forested, seven subplots were partially forested, and 138 subplots were nonforested. For the St. Cloud study area, measurements for 268 plots or 1072 subplots were used of which 226 subplots were completely forested, 13 subplots were partially forested, and 833 subplots were nonforested.

## 3. Methods

### 3.1. k-Nearest Neighbors technique

The k-Nearest Neighbors (k-NN) technique is a nonparametric approach to predicting values of point variables on the basis of similarity in a covariate space between the point and other points with observed values of the variables. Nearest neighbor techniques have not been used extensively in forestry estimation except in the Nordic countries (Fazakas & Nilsson, 1996; Katila & Tomppo, 2001; Tokola, 2000; Tokola, Pitkanen, Partinen, & Muinonen, 1996; Tomppo, 1991). Elsewhere, Moeur (1988) reported on the use of nearest neighbor techniques for multivariate forestry estimation, Trotter, Drymond, and Goulding (1997) used the k-NN technique to estimate volume, and McRoberts, Franco-Lopez, Ek, and Bauer (2000) and Franco-Lopez, Ek, and Bauer (2001) reported using the k-NN technique to classify satellite imagery.

For this application, consider a TM pixel to be a point, let $Y_i$ denote a ground attribute (e.g., forest land proportion, cumulative volume of individual trees, tree density) for the $i$th pixel, and let $\mathbf{X}_i$ denote the $i$th pixel's vector of TM spectral values. For a finite number, $N$, of image pixels of which n correspond to FIA subplots, the data points ($Y_i,\mathbf{X}_i$) may be reordered without loss of generality so that ($Y_i,\mathbf{X}_i)_{i=1,\ldots,n}$ denote the points corresponding to pixels associated with FIA subplots and ($Y_i,\mathbf{X}_i)_{i=n+1,\ldots,N}$ denote the points for the remaining pixels. With the k-NN technique, a prediction for any $Y_{j,j=1,\ldots,N}$ is obtained in two steps:

1. for each $Y_j$, reorder $Y_{i,\ i=1,\ldots,n}$ with respect to increasing distance, $d_{ji}$, between $\mathbf{X}_j$ and each $\mathbf{X}_{i,\ i=1,\ldots,n}$, excluding $Y_j$ from the ordering if $1 \leq j \leq n$, and denote the resulting ordering $\{Y_{ji}\}$;
2. for each $Y_j$,

$$\hat{Y}_j = \left(\frac{1}{k}\right)\left(\sum_{i=1}^{k} w_{ji}\right)^{-1}\left(\sum_{i=1}^{k} w_{ji} Y_{ji}\right) \quad (1)$$

where $k$ is a predetermined constant, $1 \leq k < n$, and $\{w_{ji}\}$ are point weights to be selected.

The quality of predictions may be assessed using $Y_{i,\ i=1,\ldots,n}$, an appropriate objective criterion, and the leaving-one-out method. With the leaving-one-out method, a k-NN prediction, $\hat{Y}_i$ is sequentially obtained for each $Y_{i,\ i=1,\ldots,n}$, but with the provision that $Y_i$ itself cannot be included in the mean forming its own k-NN prediction. In addition, to avoid issues related to the high correlation expected among attributes for subplots of the same plot, for this study the prediction for a subplot was constrained against including an observation for any of the other three subplots of the same plot. By comparing the observations, $Y_{i,\ i=1,\ldots,n}$, and the corresponding predictions with respect to a selected objective criterion, the quality of predictions may be evaluated.

Before implementation, the k-NN technique must be calibrated. First, the particular spectral bands used to calculate the distances, $d_{ji}$, between $\mathbf{X}_j$ and each element of the set, $\mathbf{X}_{i,\ i=1,\ldots,n}$, must be selected. Without loss of generality, the spectral band components, $X_{im}$, of $\mathbf{X}_i$ may be reordered

so that $m = 1, \ldots, M$ denote the selected bands. Second, a distance metric, $d$, must be selected; among the alternatives are weighted Euclidean distance,

$$d_{ji} = \left[ \sum_{m=1}^{M} v_m (X_{jm} - X_{im})^2 \right]^{1/2} \tag{2}$$

where $\{v_m\}$ are variable weights, and Mahalanobis distance,

$$d_{ji} = (\mathbf{X}_j' - \mathbf{X}_i')' \, \mathbf{V}^{-1} (\mathbf{X}_j' - \mathbf{X}_i') \tag{3}$$

where only the selected $M$ components of $\mathbf{X}$ are used and $\mathbf{V}$ is the covariance matrix for the $M$ components of $\mathbf{X}$. If weighted Euclidean distance is selected, then the variable weights $\{v_m\}$ for Eq. (2) must also be selected. Third, the value of $k$, the number of nearest neighbors to be included in the calculation of predictions (Eq. (1)), must be selected. Finally, the point weights, $\{w_{ji}\}$, for Eq. (1) must be selected; common alternatives include constant weighting for which $w_{ji} = 1$, inverse distance weighting for which $w_{ji} = d_{ji}^{-1}$, and inverse distance squared weighting for which $w_{ji} = d_{ji}^{-2}$.

The $k$-NN analyses were conducted at the subplot-pixel level, because a plot-level approach would require calibration using means of inventory observations over the four subplots and either means of TM spectral values over the four pixels corresponding to the four subplots or means over a block of pixels covering the plot. Predictions for image pixels must likewise then be based on the mean over four pixels in the same configuration as the four pixels corresponding to the four subplots or the mean over a block of pixels of the same size and configuration as the block covering the plot. For this study, subplot-pixel-level analyses entail a more simple approach without sacrificing statistical validity. Thus, each subplot was associated with the TM pixel with center closest to the subplot center.

### 3.1.1. Precautions

Calibration of the $k$-NN technique is guided by optimization of the objective criterion using the leaving-one-out method with the set $Y_{i, i = 1, \ldots, n}$, but with attention to five precautions and one plea for objectivity. The nature of and rationale for the precautions are illustrated using forest land proportion as the ground attribute, inventory subplot observations as surrogates for pixel-level observations, root mean square error,

$$\mathrm{RMS_e} = \left[ \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 \right]^{1/2}, \tag{4}$$

as the objective criterion, constant variable weighting, and constant point weighting. The denominator, $n$, in the expression for $\mathrm{RMS_e}$ (Eq. (4)) is the number of inventory subplot observations; no provision is made to adjust the denominator of Eq. (4) for the number of bands selected to

calculate the distance, $d$. The precautions are discussed and illustrated only as a means of creating awareness; the consequences of ignoring them will vary by application.

The first precaution is that small $k$-values may result in $\mathrm{RMS_e}$ values that are larger than the standard deviation of the observations. For the St. Louis study area and TM bands N3 and M5, the $\mathrm{RMS_e}$ value for $k = 1$ was greater than the $\mathrm{RMS_e}$ value that resulted when the overall mean was used for each prediction (Fig. 2). The general pattern of the curve was typical: large initial decreases in $\mathrm{RMS_e}$ with increasing $k$, a more gradual decrease in $\mathrm{RMS_e}$ as the optimal $k$ was approached, and then a very gradual increase in $\mathrm{RMS_e}$ approaching the $\mathrm{RMS_e}$ value that corresponded to the overall mean as $k$ approached the number of observations. Thus, users are advised that with $k$-NN analyses, unlike with simple linear regression, it is possible to obtain results that are worse than using the mean over all observations for every prediction.

The second precaution is that the $k$-NN technique produces biased estimates for pixels corresponding to the extremes of the distribution of observations. The precaution is illustrated for the St. Louis study area with simulated data generated in two steps: (1) a logistic model,

$$E(Y) = [1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)]^{-1} \tag{5}$$

was fit to the actual forest land proportion observations ($Y$) where $x_1$ and $x_2$ represent the spectral values of TM bands N3 and N4, respectively, for pixels associated with inventory subplots, and the $\beta$'s are parameters estimated from the observations; and (2) simulated data were calculated as the sum of residuals randomly generated from selected Gaussian distributions and expected values from Eq. (5) using the estimated parameters and actual spectral values. In this manner, the simulated data were generated in appropriate numbers and with appropriate characteristics to more clearly
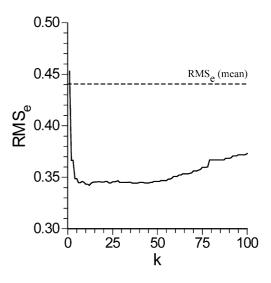


Fig. 2. $\mathrm{RMS_e}$ versus $k$ for predicting forest land proportion for the St. Louis study area with TM bands N3 and M5.
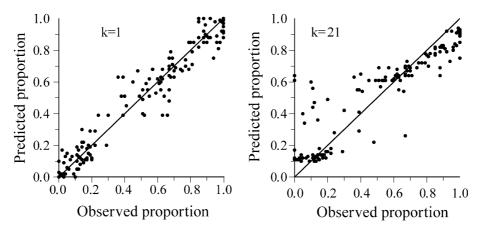
Fig. 3. *k*-NN prediction bias for forest land proportion for the St. Louis study area using TM bands N3 and N4 for *k* = 1 and *k* = 21.

illustrate the precaution. The *k*-NN technique was applied to the simulated data with *k* = 1 and *k* = 21, and predictions were graphed against their respective simulated observations. For *k* = 1, the *k*-NN predictions exhibited little systematic bias and fairly constant residual variability, while for *k* = 21, the predictions corresponding to extreme observations exhibited considerable bias and heterogeneous residual variability (Fig. 3). The bias occurred because predictions corresponding to extreme observations were calculated as means of *k* observations that were mostly larger, or smaller, respectively, than the observations themselves. Bias was worse and extended to more pixels for larger *k*-values and for smaller numbers of observations.

The third precaution is that unrelated variables included in the subset of covariates used to calculate distances, *d*, may not only fail to improve the objective criterion, but actually may have adverse effects. For the best combination of two autumn St. Louis bands, N3 and N4, and corresponding optimal *k*-value, including only one of the remaining four other autumn bands, N2, improved $RMS_e$ (Table 1a). Including the other three bands actually increased $RMS_e$ from 1% to 7%. For both the St. Louis and St. Cloud areas, the smallest $RMS_e$ value for each combination of $m \leq 6$ bands shows that overall the smallest $RMS_e$ occurred for *m* = 3 for St. Louis and *m* = 4 for St. Cloud (Table 1b). Although not shown in the table, the smallest $RMS_e$ values for *m* > 6 were all larger than the smallest $RMS_e$ value for *m* = 6. Finally, the five band combinations with the overall smallest $RMS_e$ values, regardless of the number of bands, included far fewer than the maximum number of 12 bands

for St. Louis and 18 bands for St. Cloud (Table 2). Franklin (2001) reported similar results for the maximum likelihood classifier for a forestry application, noting that accuracies may decrease as additional data layers are included, even when the additional layers include new information. Hughes (1968) investigated and reported the same phenomenon in a generic information theory context. In regression analyses, $RMS_e$ may increase as more variables are included but only when the denominator in Eq. (4) is expressed as degrees of freedom, not when it is simply the number of observations. Thus, *k*-NN analyses are also unlike regression analyses in that inclusion of additional predictor variables may actually increase residual uncertainty.

The fourth precaution is that observations for pixels separated by large spectral distances may be negatively correlated. In the geographic space of latitude and longitude, observations for pixels in close proximity are expected to be positively correlated, while observations for pixels separated by sufficiently large distances are expected to be uncorrelated. Thus, variograms are used in kriging analyses to estimate the distances at which there are no longer relationships among observations and beyond which observations should not be used in predictions. Because observations are at worst uncorrelated in kriging analyses in geographic space, the consequences of violating the distance restriction may not be severe. However, when using *k*-NN techniques in spectral space, large *k*-values may cause negatively correlated observations to be included in the *k*-NN predic-

Table 1a

Effect on $RMS_e$ of including an additional variable for the St. Louis study area using the November TM image for predicting forest land proportion

| Bands | $RMS_e$ |
|---|---|
| N3 N4 | 0.300 |
| N3 N4 N1 | 0.303 |
| N3 N4 N2 | 0.297 |
| N3 N4 N5 | 0.319 |
| N3 N4 N7 | 0.321 |

Table 1b

Best combinations by number (*m*) of band for predicting forest land proportion

| *m* | St. Louis study area | | | St. Cloud study area | | |
|---|---|---|---|---|---|---|
| | $RMS_e$ | *k* | Bands | $RMS_e$ | *k* | Bands |
| 1 | 0.3243 | 1 | N4 | 0.3241 | 51 | J6 |
| 2 | 0.2989 | 13 | N2 N4 | 0.2726 | 41 | M3 M4 |
| 3 | 0.2652 | 9 | N3 N4 M4 | 0.2420 | 23 | J3 M3 M4 |
| 4 | 0.2687 | 7 | N1 N3 N4 M4 | 0.2392 | 21 | J2 J3 M3 M4 |
| 5 | 0.2693 | 13 | N2 N3 N4 M1 M4 | 0.2399 | 29 | J2 J3 M1 M3 M4 |
| 6 | 0.2706 | 11 | N1 N2 N3 N4 M2 M4 | 0.2446 | 37 | J1 J2 J3 N2 M3 M4 |

Table 2
Mean forest land proportion estimates

| Rank[a] | $RMS_e$ | Bands | $k$ | Optimal between strata boundaries[b] | | | Mean | SE[c] | RE[d] | PREC[e] | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | 1 panel | 5 panels |
| *St. Louis study area* | | | | | | | | | | | |
| 1 | 0.2652 | N3 N4 M4 | 9 | 0.20 | 0.75 | 0.95 | 0.7547 | 0.0175 | 4.4836 | 0.0455 | 0.0203 |
| 2 | 0.2687 | N1 N3 N4 M4 | 7 | 0.10 | 0.70 | 0.80 | 0.7493 | 0.0188 | 3.8864 | 0.0490 | 0.0219 |
| 3 | 0.2690 | N2 N3 N4 M4 | 9 | 0.20 | 0.50 | 0.75 | 0.7593 | 0.0177 | 4.3943 | 0.0459 | 0.0205 |
| 4 | 0.2690 | N3 N4 M1 M4 | 11 | 0.55 | 0.60 | 0.90 | 0.7845 | 0.0157 | 5.5912 | 0.0400 | 0.0179 |
| 5 | 0.2693 | N2 N3 N4 M1 M4 | 13 | 0.50 | 0.70 | 0.95 | 0.7681 | 0.0182 | 4.1655 | 0.0469 | 0.0210 |
| *St. Cloud study area* | | | | | | | | | | | |
| 1 | 0.2392 | J2 J3 M3 M4 | 21 | 0.10 | 0.40 | 0.75 | 0.2312 | 0.0107 | 4.9189 | 0.0635 | 0.0285 |
| 2 | 0.2399 | J2 J3 M1 M3 M4 | 29 | 0.15 | 0.45 | 0.70 | 0.2313 | 0.0107 | 4.9587 | 0.0635 | 0.0284 |
| 3 | 0.2406 | J2 J3 M2 M3 M4 | 33 | 0.20 | 0.40 | 0.75 | 0.2346 | 0.0109 | 4.7837 | 0.0643 | 0.0287 |
| 4 | 0.2420 | J3 M3 M4 | 23 | 0.25 | 0.45 | 0.60 | 0.2367 | 0.0103 | 5.3847 | 0.0605 | 0.0270 |
| 5 | 0.2423 | J3 M1 M3 M4 | 23 | 0.25 | 0.45 | 0.65 | 0.2398 | 0.0105 | 5.1238 | 0.0612 | 0.0274 |

[a] Rank based on $RMS_e$.
[b] The optimality criterion is maximization of the relative efficiency (RE) of the stratification.
[c] Standard error of mean, calculated as the square root of the variance of the mean (Eq. (7)).
[d] Relative efficiency of the stratification, calculated as the ratio of the variance of the mean assuming simple random sampling and the variance of the mean based on stratified analyses (Eq. (7)).
[e] FIA precision calculated from Eq. (9).

tions, thus adversely affecting the objective criterion. To illustrate this phenomenon, a correlogram was constructed for the St. Louis study area data using bands J2, J3, M3, and M4. Inventory subplot observations were grouped into categories based on their spectral separation distances calculated using Eq. (2), and the within category correlations were graphed against the mean spectral separation distance for the category (Fig. 4). Sufficient field data to span the ranges of both image and field variables lessens the risk of this phenomenon.

The fifth precaution, for which no illustration is provided, is that multiple pixels may be at the same spectral distance, $d$, from a pixel for which a $k$-NN prediction is desired, particularly when $M$, the number of selected bands
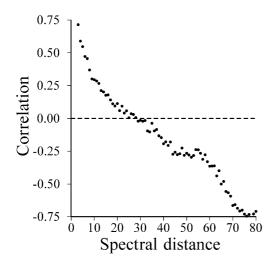
used to calculate $d$ is small. In this situation, either $k$-NN predictions should include all observations at the same distance as the $k$th ordered observation or a random selection procedure should be used to select observations at the $k$th ordered distance from the set of all such observations.

*3.1.2. Objective calibration criteria*

Finally, the plea relates to selecting, using, and reporting objective criteria for selecting $k$-values. A variety of criteria have been used: some quite objective, others more subjective. Trotter et al. (1997) selected $k = 15$ for predicting wood volume on the basis of maximizing $r^2$. Tokola et al. (1996) noted that for predicting volume, standard errors decreased rapidly as $k$ increased from 1 to 10 and finally selected $k = 15$ using the rationale that with additional plots there was little decrease in standard errors. Tokola (2000) selected $k = 15$ for another volume prediction application because it was adequate for the previous study (Tokola et al., 1996). Katila and Tomppo (2001) selected $k$ according to three criteria related to minimizing $RMS_e$: $1 \leq k \leq 30$; the significance of bias was controlled; and the proportional change in $RMS_e$ between $k$ and $k + 1$ was less than 0.005. Franco-Lopez et al. (2001) provide an excellent discussion of issues related to selecting $k$, noting that for predicting volume or basal area, $RMS_e$ decreased by 14% as $k$ increased from 1 to 5 and that after $k = 9$, the marginal increase in precision was less than 0.005. They recommended that if prediction variance similar to observation variance is desired, then $k = 1$ is the appropriate selection, while if minimization of $RMS_e$ is desired, then $k = 9$ is the appropriate selection.

These criteria for selecting $k$ represent varying levels of objectivity where objectivity is assessed in terms of whether an independent user applying the same criteria would be expected to make the same selection. Maximizing a quantity



Fig. 4. Correlogram for forest land proportion for the St. Louis study area using TM bands J2, J3, M3, and M4.

such as $r^2$ is certainly an objective criterion. However, selecting $k$ on the basis of its adequacy for a similar application is not objective and may entail substantial risk. For predicting forest land proportion for the St. Louis study area, optimal $k$-values were in the range $7 \leq k \leq 13$, while for the same application for the St. Cloud study area, optimal $k$-values were in the range $21 \leq k \leq 33$ (Table 2). These substantial differences in optimal $k$-values for the two study areas illustrate that $k$-NN calibrations for very similar applications may be quite different. In addition, a criterion based on an arbitrary specification of a maximum $k$-value is not objective. However, specification of a maximum $k$-value on the basis of correlation or similar analyses may be considered part of a set of objective criteria.

The criterion of selecting the first $k$-value that corresponds to a proportional decrease in $RMS_e$ below a specified value, such as 0.005, may lead to sub-optimal selections for at least two situations. First, a lengthy series of proportional decreases in $RMS_e$ below the specified value may still produce a large cumulative decrease in $RMS_e$. Second, consecutive proportional changes in $RMS_e$ are not guaranteed to decrease or even to be positive. For the St. Louis study area with bands N1, N3, N4, and M4, the proportional change in $RMS_e$ fell below 0.005 for $k = 3$, but subsequently rose above 0.005 before the optimal value of $k = 7$ was reached (Fig. 5). For $k = 3$, the corresponding $RMS_e = 0.2975$ is an increase of 11% over $RMS_e = 0.2687$ for the optimal $k = 7$. Results with this proportional change criterion would be expected to be variable and inconsistent.

Thus, the plea has three parts: first, $k$-values should be selected according to objective criteria; second, users should report the criteria they use; and third, users should justify their selection of $k$-values according to those criteria. The first part of the plea requires only that an objective criterion be used to select $k$; it does not require that the same $k$-value be selected for different applications or for the same application with different data, even when the criterion is
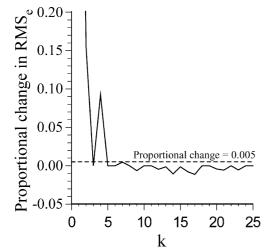
the same. Two intuitive objective choices are $k = 1$ and $k = k_{opt}$, the value of $k$ that optimizes the objective criterion. The rationale for $k = 1$ is that it incorporates into the predictions all the variability that exists in the observations, whereas $k > 1$ incorporates less variability because the predictions are based on means of multiple observations. However, when $RMS_e$ obtained for a small $k$-value is greater than $RMS_e$ corresponding to the overall mean, then the standard deviation of the variability incorporated in the $k$-NN predictions is greater than the standard deviation of the observations, an undesirable condition. The rationale for $k = k_{opt}$ is that the objective criterion is optimized, but if $k_{opt}$ is large, then predictions may not retain appropriate variability for mapping applications and prediction bias may be unacceptable. Thus, the selection of $k$ may require a compromise between a small $k$ that avoids prediction bias and retains appropriate variability in predictions and a large $k$ that optimizes the objective criterion. Nevertheless, the compromise selection must be made using objective criteria. Although the proportional change approach of Katila and Tomppo (2001) and Franco-Lopez et al. (2001) represents such a compromise, the compromise is not based on an objective criterion because the $RMS_e$ corresponding to $k_{opt}$ is not considered in the selection. An alternative compromise that is more objective is to select the smallest $k$-value such that the corresponding $RMS_e$ is not more than a specified percentage (e.g., 0.5%, 1%, 5%) greater than the $RMS_e$ corresponding to $k_{opt}$.

Although criteria for selecting $k$ must be allowed to vary to accommodate the particular objectives of each application, the criteria should be objective. The same objective criteria may lead to selections of different $k$-values for different applications or for similar applications with different data. However, when multiple users independently select $k$-values for the same application with the same data, objective criteria produce repeatable results. Repeatable results would not be expected for arbitrary selections of small $k$-values, constraints on maximum $k$-values, an arbitrary compromise between small $k$-values and $k$-values that optimize calibration criteria, previous experience, or arbitrarily small decreases in $RMS_e$. Repeatability would be expected when $k$ is selected as the value that maximizes $r^2$, the $k$-value that minimizes $RMS_e$, or the smallest $k$-value with $RMS_e$ not more than 1% larger than the minimum $RMS_e$ over all $k$-values.

### 3.2. Analyses

#### 3.2.1. Calibrating the k-NN technique

The compromises associated with selecting $k$-values may be partially avoided when $k$-NN predictions are used to assign pixels to categories or classes. Although prediction bias may be unacceptable for producing a continuum of predictions, the assignment of pixels to classes may be relatively unbiased. In addition, when the classes are obtained from pixel predictions and used as the basis for

Fig. 5. Proportional change in $RMS_e$ for forest land proportion predictions for the St. Louis study area with TM bands N1, N3, N4, and M4.

creating strata for stratified analyses, prediction or classi-fication bias adversely affects only the variance reduction obtained with the stratification; it does not produce bias in the stratified mean estimates. When using $k$-NN predictions to create strata, greater weight is attributed to assigning pixels to correct classes than to minimizing bias on the extremes of distributions or to preserving variability in the predictions. Thus, $k$-values were selected using the following steps:

1. initially, select $k = k_{opt}$, the $k$-value that optimizes the objective criterion, $RMS_e$;
2. on a prediction-by-prediction basis, reduce $k$ if necessary to ensure exclusion from the $k$-NN predictions of observations at spectral distances greater than the distance at which the correlogram indicates negative correlation;
3. on a prediction-by-prediction basis, increase $k$ if necessary to ensure that all observations at the same or smaller spectral distances (Eq. (2)) than the $k$th ordered observation are used in $k$-NN predictions.

For each study area, the optimal $k$-value was determined using these steps for each combination of spectral bands by comparing values of $RMS_e$ obtained with constant variable and constant point weighting. Constant variable weighting was selected to facilitate implementation of the $k$-NN algorithm. Constant point weighting was selected, because for the best band combinations for each $M \leq 6$, the proportional reduction in $RMS_e$ with inverse distance weighting compared to constant weighting never exceeded 0.006 for the St. Louis study area or 0.020 for the St. Cloud study area. In addition, for the St. Cloud study area, the overall smallest $RMS_e$ with inverse distance weighting was actually greater than the overall smallest $RMS_e$ with constant weighting. For each study area, the five spectral band combinations with smallest $RMS_e$ obtained using $k$-values determined using the above steps, but without regard to the number of bands, were selected for further evaluation.

### 3.2.2. Creating strata

For each of the five best spectral band combinations for each study area, forest land proportion was predicted for each pixel using the $k$-NN technique with the $k$-value determined using the above steps. For each study area and band combination, the continuum of predictions was divided into four strata by selecting strata separation boundaries subject to three constraints: first, the lower bound of the first stratum was always 0.00, and the upper bound of the fourth stratum was always 1.00; second, the minimum stratum width was 0.05; and third, at least five plots had to be assigned to each stratum. Strata were limited to four because the preponderance of observed forest land proportions were either 0.00 or 1.00. All possible stratifications, subject to the constraints, were evaluated with respect to the relative efficiency that the stratification produced. Relative efficiency, RE, is the ratio of

the variance of the simple arithmetic mean of inventory subplot observations calculated under an assumption of simple random sampling and the variance obtained using stratified analyses. The four optimal strata corresponded to the stratification with largest RE.

Stratified estimation was accomplished by assigning each pixel to a stratum based on its forest land proportion prediction, and strata weights were calculated as the proportions of pixels assigned to strata. To avoid the mathematical complexity necessary to accommodate the spatial correlation among the four subplot observations, FIA assigns plots rather than subplots to strata for stratified analyses. Plots were assigned to strata on the basis of the stratum assignment of the pixel corresponding to the center of the center subplot. Plots were stratified using $k$-NN predictions of forest land proportion for their corresponding pixels rather than observations so that the assignment of plots to strata would be consistent with the calculation of strata weights.

The optimal strata with respect to RE were expected to differ between the study areas, thereby raising the question of whether the sub-optimal nature of a common stratification across areas covered by multiple TM scenes would adversely affect the precision of stratified estimates for the combined area. For each of the 25 combinations of one of the five best St. Louis band combinations and one of the five best St. Cloud band combinations, the optimal strata boundaries were determined. Although the predictions of forest land proportion for each pixel and plot from the separate study area analyses were retained, the division of the continuum of these predictions into optimal strata was common for the two study areas.

### 3.2.3. Stratified estimation

Stratified estimates of mean forest land proportion, $\bar{Y}$, and estimated variance, $\widehat{Var}(\bar{Y})$, were calculated using standard methods (Cochran, 1977):

$$\bar{Y} = \sum_{j=1}^{J} w_j \bar{Y}_j \tag{6}$$

and

$$\widehat{Var}(\bar{Y}) = \sum_{j=1}^{J} w_j^2 \hat{\sigma}_j^2 / n_j, \tag{7}$$

where $j = 1, \ldots, J$ denotes stratum; $w_j$ is the weight for the $j$th stratum; $\bar{Y}_j$ denotes the mean forest land proportion for plots assigned to the $j$th stratum; $n_j$ is the number of plots assigned to the $j$th stratum; and $\sigma_j^2$ is the within-stratum variance for the $j$th stratum calculated as,

$$\hat{\sigma}_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_j)^2; \tag{8}$$

where $Y_{ij}$ is the forest land proportion observed by the field crew for the $i$th plot in the $j$th stratum. Variance estimates

obtained using Eq. (7) ignore the slight effects due to finite population correction factors and to variable rather than fixed numbers of plots per stratum.

The FIA program reports precision estimates as coefficients of variation scaled to compensate for varying sample sizes using as a reference standard the sample size corresponding to 404,694 ha (1 million acres) (USDA-FS, 1970). For forest area estimate, $\widehat{FA} = A\bar{Y}$ the scaled precision estimate, denoted PREC, is defined for this study as

$$
\begin{aligned}
\text{PREC} &= \frac{[\widehat{\text{Var}}(\widehat{FA})]^{1/2}}{\widehat{FA}} \left[\frac{\widehat{FA}}{404,694}\right]^{1/2} \\
&= \frac{[\widehat{\text{Var}}(\bar{Y})]^{1/2}}{\bar{Y}} \left[\frac{A\bar{Y}}{404,694}\right]^{1/2},
\end{aligned} \tag{9}
$$

where $\bar{Y}$ is again mean forest area proportion per plot, and $A$ is total area inventoried in hectares. Two values of PREC are reported: the value obtained from Eq. (9) that corresponds to the sample size resulting from a single panel of plot measurements and the value obtained from Eq. (9) divided by the square root of 5, which corresponds approximately to the value expected with the sample size resulting from all five panels of plots. The national FIA precision standard is PREC ≤ 0.03.

## 4. Results

In general, the $k$-NN algorithm was very simple to implement, straight forward to calibrate, and required no user intervention after initiation. The $k$-NN predictions captured much of the forest/nonforest detail and provided an excellent basis for stratifications. When compared to variances of forest area estimates obtained using simple random estimation, the variances obtained using stratified estimation were smaller by factors as great as 5. Specific results follow.

### 4.1. Precautions

In the Methods section, five precautions were noted:

1. Small $k$-values may result in $\text{RMS}_e$ values that are larger than the standard deviation of the observations.
2. The $k$-NN technique produces biased estimates for pixels corresponding to the extremes of the distributions of observations.
3. Unrelated variables used to calculate distances, $d$, may cause an increase in $\text{RMS}_e$.
4. Observations for pixels separated by large spectral distances may be negatively correlated.
5. Multiple pixels may be at the same spectral distance from a pixel for which a $k$-NN prediction is desired.

All five precautions were either observed or found not to apply for these analyses. The first precaution was observed by beginning the procedure for the selection of $k$ using the objective criterion of minimizing $\text{RMS}_e$. The second precaution was not relevant when using the classified images for stratification, because bias in the pixel predictions does not produce bias in the estimates of stratified means. The third precaution was observed by selecting only the spectral band combinations that were among the five best for optimizing the objective criterion. The fourth precaution was observed by adjusting $k$ for individual predictions to ensure that negatively correlated observations were not used. The correlogram analyses indicated that correlations among observations were positive for spectral distances less than about 15–25 for the St. Louis study area and about 20–25 for the St. Cloud study area. Decreasing $k$ to values less than $k_{\text{opt}}$ to avoid including negatively correlated observations in the $k$-NN predictions was seldom necessary. The fifth precaution was observed by adjusting the $k$-value to include all observations at the same spectral distance but was necessary only for 1- and 2-band combinations. In general, the base criterion of selecting $k = k_{\text{opt}}$ was usually sufficient for individual $k$-NN predictions. However, these results should not be generalized to other applications or other data sets.

### 4.2. k-NN calibrations

Both similarities and differences were noted among calibrations for the five best band combinations and the resulting stratified estimates (Table 2). The following similarities were noted.

(1) The means for the five best band combinations were comparable within study areas.

(2) Values of $\text{RMS}_e$, SE, RE, and PREC were generally of the same order of magnitude both within and between study areas.

(3) The bands selected for the five best band combinations were similar within study areas: N3, N4, and M4 were selected for all five combinations for the St. Louis study area; and J3, M3, and M4 were selected for all five combinations for the St. Cloud study area. Bands 3 and 4 were most commonly selected, and bands from the spring or summer 2000 images were selected for all five best band combinations for both study areas.

(4) For both study areas, the stratifications based on the $k$-NN analyses produced expected five-panel precision for forest land area estimates that satisfied the national FIA precision standards.

(5) For each best band combination, multiple sets of between strata boundaries produced similar values of RE.

The following differences were noted.

(1) The ordering of the band combinations with respect to RE, or equivalently PREC, was not the same as that with respect to $\text{RMS}_e$, suggesting that if the optimal band combination with respect to RE or PREC is desired, then multiple best band combinations selected with respect to $\text{RMS}_e$ should be evaluated as was done in this study.

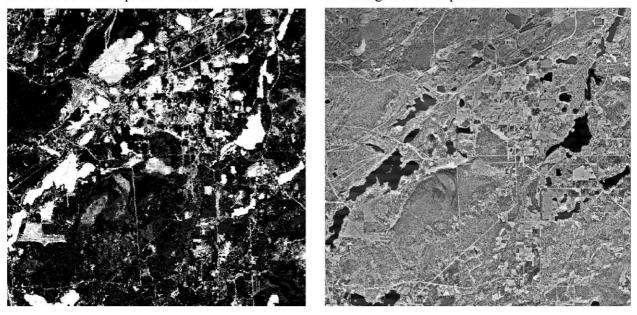## 1999-2000 k-NN predictions

## 1992 digital orthoquad



Fig. 6. Predicted forest land proportion and aerial photograph for the 15 × 15 km center of the St. Louis study area.

(2) Optimal between strata boundaries for the St. Louis study area differed considerably, although as noted previously, multiple sets of between strata boundary combinations for the same band combination produced similar values of RE.

(3) Optimal *k*-values for the St. Cloud study area were considerably larger than those for the St. Louis study area. However, even for a given criterion, optimal *k*-values would be expected to differ as a result of many factors including the attributes of interest; the number, size, and accuracy of locations of plots; the accuracy of image registration; and variability in the area of interest.

The 25 common stratifications across both study areas produced stratified estimates of the mean for the combined study area that were slightly less precise than the estimates for individual study areas. Nevertheless, 10 of the 25 five-panel expected PREC values for the combined estimates still satisfied the national precision requirement of PREC ≤ 0.03. The combined area estimates of the mean ranged from 0.6038 to 0.6298; the combined area estimated standard errors of the mean ranged from 0.0100 to 0.0126; and the combined area five-panel expected PREC values ranged from 0.0253 to 0.0328. On average, across the 25 combinations, the between strata boundaries were 0.42, 0.68, and 0.90. These results suggest that the precision resulting from a common set of optimal strata boundaries for an area covered by multiple TM images may also be expected to satisfy the national FIA precision standards.

### 4.3. Predicting forest land proportion

The *k*-NN predictions for the central portion of the St. Louis area obtained with the calibration corresponding to

the minimum overall RMS$_e$ (Table 2) portrayed a substantial portion of the forest/nonforest detail (Fig. 6). The water is clearly identified as is much of the road network, the airport in the lower left quadrant, nonforested areas in the upper left quadrant, and a large area in the upper left quadrant that was apparently cleared of vegetation between the date of the aerial photograph and the date of the TM imagery.

## 5. Conclusions and discussion

Four important conclusions may be drawn from this study. First, although the *k*-NN technique is conceptually easy to implement, careful attention must be paid to its calibration to achieve optimal results. Second, stratifications derived from classified TM imagery reduced variances of forest area estimates by factors as great as 5 for both a heavily forested area and a sparsely forested area. Third, the stratifications may be expected to produce forest land area estimates that satisfy national FIA precision standards for sample sizes corresponding to five panels of measurements. Fourth, the *k*-NN technique is a viable alternative for predicting forest land proportion from satellite imagery that is as fast and easy to implement as traditional image classification methods.

The implications of the latter three conclusions for the FIA program are considerable. First, in the absence of stratification, sample sizes would have to be increased by factors at least as great as 5 to achieve the same level of precision as obtained with the stratifications. The magnitude of the resulting cost saving is substantial. For the State of Minnesota with a sampling intensity of one plot for every 2403 ha, approximately 825 plots are field-measured annu-

ally at an FY1999 cost of approximately US$1000 per plot. Thus, the annual cost savings obtained with such stratifications is approximately US$3,300,000.

Second, the effectiveness of the *k*-NN algorithm frees the FIA program from more costly and less timely alternatives. The speed and automation of the *k*-NN technique make it vastly superior to FIA's time-consuming, labor-intensive, traditional approach based on interpreting aerial photographs. With a complete set of aerial photographs available, a crew of four photo-interpreters, working full-time, could be expected to complete the photo-interpretation and stratification task for the State of Minnesota in approximately 2 years. Alternatively, with the 19 rectified TM images for Minnesota available, prediction of forest land proportion and stratification could be expected to be accomplished using the *k*-NN technique in 2–3 weeks by a single remote sensing technician.

Third, the speed and ease of implementation of the *k*-NN technique relieves the FIA program from dependence on external agencies and programs. Currently, several regional FIA programs rely on or are investigating the NLCD as the basis for stratifications. This classification was based on nominal 1992 TM imagery and a suite of ancillary data, was not available until approximately 5 years after the date of the imagery, and will be replaced with a more current classification only every 10 years. With plot measurement cycles of 5 years for much of the US, it is inefficient to rely on stratifications based on classifications that are 5–15 years out of date. In addition, the NLCD has yet to be demonstrated to produce stratifications that satisfy the national precision standards. The *k*-NN technique permits FIA to implement effective stratified estimation using TM imagery with dates that are concurrent with plot measurement dates and to do so independently of other agencies and programs.

Finally, better future results may be expected with the *k*-NN technique. Fine-tuning the calibration of the *k*-NN technique by including variable-weighting will increase the accuracy of classifications. Also, five panels of plot measurements will increase the density of observations in spectral space, allow each *k*-NN prediction to be based on subplot-pixel observations in closer spectral proximity, and, therefore, increase the accuracy of individual pixel predictions.

In conclusion, the *k*-NN technique is a viable and efficient method for processing TM images to obtain predictions of forest area proportion, and stratifications derived from these predictions produce forest area estimates that may be expected to satisfy national FIA precision standards.

## References

Bickford, C. A. (1960). A test of continuous inventory for National Forest management based upon aerial photographs, double sampling, and re-measured plots. In Anonymous (Eds.), *Proceedings of the Society of American Foresters Meeting* ( pp. 143–148).

Cochran, W. G. (1977). Sampling techniques. (3rd ed.). New York: Wiley.

Deppe, F. (1998). Forest area estimation using sample surveys and Landsat MSS and TM data. *Photogrammetric Engineering and Remote Sensing*, *64*, 285–292.

Fazakas, Z., & Nilsson, M. (1996). Volume and forest cover estimation over southern Sweden using AVHRR data calibrated with TM data. *International Journal of Remote Sensing*, *17*, 1701–1709.

Franco-Lopez, H., Ek, A. R., & Bauer, M. E. (2001). Estimation and mapping of forest stand density, volume, and cover type using the *k*-nearest neighbors method. *Remote Sensing of Environment*, *77*, 251–274.

Franklin, S. E. (2001). Remote sensing for sustainable forest management. Boca Raton, FL: CRC Press.

Hansen, M. H. (1990). *A comprehensive sampling system for forest inventory based on an individual tree growth model*. PhD dissertation. St. Paul, MN: University of Minnesota, College of Natural Resources.

Hansen, M. H., & Wendt, D. G. (2000). Using classified Landsat Thematic Mapper data for stratification in a statewide forest inventory. In R. E. McRoberts, G. A. Reams, & P. C. Van Deusen (Eds.), *Proceedings of the first annual forest inventory and analysis symposium* (General Technical Report, NC-213, pp. 20–27). St. Paul, MN: US Department of Agriculture, Forest Service, North Central Research Station.

Hughes, G. F. (1968). On the mean accuracy of statistical pattern recognisers. *IEEE Transactions on Informational Theory*, *IT-14*(1), 55–63.

Hoppus, M., Arner, S., & Lister, A. (2001). Stratifying FIA ground plots using a 3-year old MRLC forest cover map and current TM-derived variables selected by "decision tree" classification. In G. A. Reams, R. E. McRoberts, & P. C. Van Deusen (Eds.), *Proceedings of the second annual forest inventory and analysis symposium* (General Technical Report, SRS-47, pp. 19–24). Asheville, NC: US Department of Agriculture, Forest Service, Southern Research Station.

Katila, M., & Tomppo, E. (2001). Selecting estimation parameters for the Finnish multisource National Forest Inventory. *Remote Sensing of Environment*, *76*, 16–32.

Loetsch, F., & Haller, K. E. (1964). *Forest inventory, volume I: statistics of forest inventory and information from aerial photographs*. Munchen: BLV Verlagsgesselschaft.

Loveland, T. L., & Shaw, D. M. (1996). Multiresolution land characterization: building collaborative partnerships. In J. M. Scott, T. Tear, & F. Davis (Eds.), *Gap analysis: a landscape approach to biodiversity planning, proceedings of the ASPRS/GAP symposium, Charlotte, NC* ( pp. 83–89).

McRoberts, R. E. (1999). Joint annual forest inventory and monitoring system: the North Central perspective. *Journal of Forestry*, *97*(2), 21–26.

McRoberts, R. E., Franco-Lopez, H., Ek, A. R., & Bauer, M. E. (July 2000). Using remotely sensed imagery to extend forest inventory plot data to large areas. In G. B. M. Heuvelink, M. J. P. M. Lemmens (Eds.), *Accuracy 2000. Proceedings of the 4th International Symposium on spatial Accuracy Assessment in Natural Resources and Environmental Sciences, Amsterdam* ( pp. 465–472).

McRoberts, R. E., Wendt, D. G., Nelson, M. D., & Hansen, M. H. (2002). Using a land cover classification based on satellite imagery to improve the precision of forest inventory area estimates. *Remote Sensing of Environment*, *81*, 36–44.

Moeur, M. (1988). Nearest neighbor inference for correlated multivariate attributes. In A. R. Ek, S. R. Shifley, & T. E. Burk (Eds.), *Forest growth modelling and prediction. Proceedings of the IUFRO conference* (General Technical Report 120, pp. 716–723). St. Paul, MN: US Department of Agriculture, Forest Service, North Central Research Station.

Poso, S., Hame, T., & Paananen, R. (1984). A method of estimating the stand characteristics of a forest compartment using satellite imagery. *Silva Fennica*, *18*, 261–292.

Poso, S., Paananen, R., & Simila, M. (1987). Forest inventory by compartments using satellite imagery. *Silva Fennica*, *21*, 69–94.

Scott, J. M., Davis, F., Csuti, B., Noss, R., Buterfield, B., Groves, C., Anderson, H., Caicco, S., D'Erchia, F., Edwards Jr., T. C., Uliman, J., & Wright, R. G. (1993). Gap analysis: a geographic approach to protection of biological diversity. *Wildlife Monograph*, *123*, 1–41.

Tokola, T. (2000). The influence of field sample data location on growing stock volume estimation in Landsat TM-based forest inventory in eastern Finland. *Remote Sensing of Environment, 74*, 422–431.

Tokola, T., Pitkanen, J., Partinen, S., & Muinonen, E. (1996). Point accuracy of a non-parametric method in estimation of forest characteristics with different satellite materials. *International Journal of Remote Sensing, 17*, 2333–2351.

Tomppo, E. (1991). Satellite imagery-based national forest inventory of Finland. *International Archives of Photogrammetry and Remote Sensing, 28*, 419–424.

Trotter, C. M., Drymond, J. R., & Goulding, C. J. (1997). Estimation of timber volume in a coniferous plantation forest using Landsat TM. *International Journal of Remote Sensing, 18*, 2209–2223.

USDA Forest Service (USDA-FS) (1970). *Operational procedures, Forest Service Handbook 4809.11*. Washington, DC: US Department of Agriculture, Forest Service (Chapter 10:11-1).

Vogelmann, J. E., Howard, S. M., Yang, L., Larson, C. R., Wylie, B. K., & Van Driel, N. (2001). Completion of the 1990s National Land Cover Data Set for the conterminous United States from Landsat Thematic Mapper data and ancillary data sources. *Photogrammetric Engineering and Remote Sensing, 67*, 650–662.