
Estimating Uncertainty in Map Intersections

Ronald E. McRoberts¹, Mark A. Hatfield², and Susan J. Crocker³

Abstract.—Traditionally, natural resource managers have asked the question “How much?” and have received sample-based estimates of resource totals or means. Increasingly, however, the same managers are now asking the additional question “Where?” and are expecting spatially explicit answers in the form of maps. Recent development of natural resource databases, access to satellite imagery, development of image classification techniques, and availability of geographic information systems has facilitated construction and analysis of the required maps. Unfortunately, methods for estimating the uncertainty associated with map-based analyses are generally not known, particularly when the analyses require maps to be combined. Motivated by the threat of the emerald ash borer in southeastern Michigan, the number of ash trees at risk was estimated by intersecting a forest/non-forest map and an ash tree distribution map. The primary objectives of the study were to quantify the uncertainty of the estimate and to partition the uncertainty by source. An important conclusion of the study is that spatial correlation—an often ignored component of uncertainty analyses—made the greatest contribution to the uncertainty in the estimate of the total number of ash trees.

Introduction

The emerald ash borer (*Agrilus planipennis* Fairmaire, Coleoptera: Buprestidae) (EAB) is a wood-boring beetle native to Asia that was initially discovered in the United States in

June 2002. It most likely entered the country in solid-wood packing material such as crates and pallets and was transported to Detroit, Michigan, at least 10 years before it was discovered there in 2002 (Cappaert *et al.* 2005, Herms *et al.* 2004). Ash trees are the only known host, and damage is the result of larval activity. Once eggs hatch, larvae bore into the cambium and begin to feed on and produce galleries in the phloem and outer sapwood. Larval feeding disrupts the translocation of water and nutrients and eventually girdles the tree. Tree mortality occurs within 1 to 3 years, depending on the severity of the infestation (McCullough and Katovich 2004, Haack *et al.* 2002). All of Michigan’s native ash species (*Fraxinus* spp.) and planted cultivars are susceptible (Cappaert *et al.* 2005). Since 2002, southeastern Michigan has lost an estimated 15 million ash trees due to the EAB (Cappaert *et al.* 2005). The natural rate of EAB dispersal is estimated to be less than 1 km per year in low-density sites. Natural dispersal has been enhanced by human transportation of infested firewood, ash logs, and nursery stock. This artificial spread of EAB has initiated the majority of outlier infestations (Cappaert *et al.* 2005). Continued spread outside of the core zone increases the threat to ash trees across the United States.

The objective of the study was twofold: (1) to estimate the uncertainty in forest/nonforest maps, ash tree distribution maps, and intersections of the two maps, and (2) to partition the total uncertainty in areal estimates of the total number of ash trees by source.

Methods

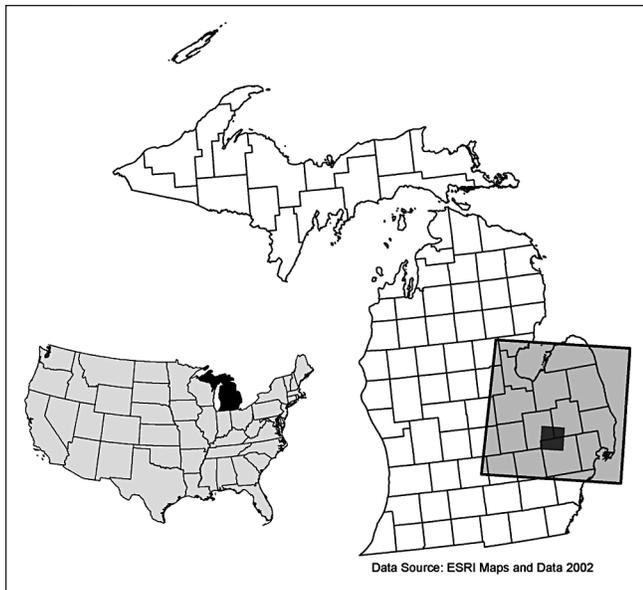
The motivating problem for the study was to calculate an estimate, A_{total} , of the total number of ash trees in a region of southeastern Michigan (fig. 1) that is susceptible to infestation

¹ Mathematical Statistician, U.S. Department of Agriculture (USDA) Forest Service, Northern Research Station, Forest Inventory and Analysis Unit, 1992 Folwell Avenue, St. Paul, MN 55108. E-mail: rmcroberts@fs.fed.us.

² Forester, USDA Forest Service, Northern Research Station, Forest Inventory and Analysis Unit, 1992 Folwell Avenue, St. Paul, MN 55108. E-mail: mahatfield@fs.fed.us.

³ Forester, USDA Forest Service, Northern Research Station, Forest Inventory and Analysis Unit, 1992 Folwell Avenue, St. Paul, MN 55108. E-mail: scrocker@fs.fed.us.

Figure 1.—State of Michigan, USA with Landsat scene path 20, row 30 (large rectangle) and 30-km x 30-km study area (small rectangle).



by the EAB. The estimation approach entails intersecting two 30-m x 30-m resolution maps, one depicting the spatial distribution of forest land and the other depicting the spatial distribution of ash trees. The technical objective was to estimate the uncertainty of A_{total} for a selected region. The map depicting the distribution of forest land was derived from a forest probability layer constructed using forest inventory plot observations, Landsat Thematic Mapper (TM) satellite imagery, and a logistic regression model. The ash tree distribution layer was constructed using the same forest inventory plot observations and inverse distance weighted spatial interpolation. Uncertainty in A_{total} was estimated using an analytical technique for estimating the covariances of the logistic regression model parameters, a sample-based technique for estimating the uncertainty in interpolated ash tree counts per hectare, and Monte Carlo techniques for generating forest/nonforest maps and for combining the components of uncertainty.

Data

The study area is wholly contained in Landsat scene path 20, row 30 (fig. 1), for which three dates of Landsat TM/ETM+ imagery were obtained: May 2002, July 2003, and October 2000. Preliminary analyses indicated that Normalized Differ-

ence Vegetation Index (NDVI) and the tasseled cap (TC) transformations (brightness, greenness, and wetness) (Kauth and Thomas 1976, Crist and Cicone 1984) were superior to both the spectral band data and principal component transformations with respect to predicting forest attributes. Thus, the predictor variables were the 12 satellite image-based variables consisting of NDVI and the three TC transformations for each of the three image dates. Mapping units for all analyses consisted of the 30-m x 30-m TM pixels.

The Forest Inventory and Analysis (FIA) program of the Forest Service, U.S. Department of Agriculture has established field plot centers in permanent locations using a sampling design that is assumed to produce a random, equal probability sample (Bechtold and Patterson 2005, McRoberts *et al.* 2005). The plot array has been divided into five nonoverlapping, interpenetrating panels, and measurement of all plots in one panel is completed before measurement of plots in the next panel is initiated. Panels in the study area are selected on a 5-year rotating basis. Over a complete measurement cycle, the Federal base sampling intensity is approximately one plot per 2,400 ha. The State of Michigan provided additional funding to triple the sampling intensity to approximately one plot per 800 ha. In general, locations of forested or previously forested plots were determined using global positioning system receivers, while locations of nonforested plots were determined using aerial imagery and digitization methods. Each plot consists of four 7.31-m radius circular subplots. The subplots are configured as a central subplot and three peripheral subplots with centers located at a distance of 36.58 m and azimuths of 0, 120, and 240 degrees from the center of the central subplot. Data for 2,995 FIA plots or 11,980 subplots with centers in the selected TM scene that were measured between 2000 and 2004 were available for the study. For each subplot, the proportion of the subplot area that qualified as forest land was determined from field crew observations. The FIA program requirements for forest land are at least 0.4 ha in size, at least 10 percent stocking, at least 36.58 crown-to-crown width, and forest land use. For each subplot, the number of observed ash trees with diameter at breast height of at least 12.5 cm was scaled to a count/hectare basis. The ash tree count/hectare for the i^{th} subplot is denoted A_i^o , where the superscript denotes a subplot observation and

is distinguished from the ash tree count/hectare for the i^{th} pixel which is denoted A_i .

The spatial configuration of the FIA subplots with centers separated by 36.58 m and the 30-m x 30-m spatial resolution of the TM/ETM+ imagery permits individual subplots to be associated with individual image pixels. The subplot area of 167.87 m² is an approximately 19-percent sample of the 900 m² pixel area, and subplot observations are assumed to adequately characterize entire pixels.

Areal Estimation

The estimate, A_{total} , for a region was calculated in three steps: (1) generate a 30-m x 30-m resolution forest/nonforest map; (2) construct a 30-m x 30-m ash tree count/hectare layer, and for each pixel, estimate the number of ash trees as the product of the ash tree count/hectare and the 0.09 ha pixel area; and (3) estimate A_{total} as the sum over forest pixels from step 1 of pixel-level estimates of the number of ash trees from step 2. Thus, two layers were required: a forest/nonforest layer and an ash tree count/hectare layer. Both layers were constructed specifically for this study to make known their pixel-level uncertainties. Two sets of analyses were conducted. First, forest/nonforest, ash tree count/hectare, and the two associated uncertainty maps were constructed for a 30-km x 30-km study area in the selected TM scene (fig. 1). Second, uncertainty analyses for A_{total} were conducted for a smaller 2-km x 2-km portion of the larger study area (figs. 2 and 3). The restriction of the latter analyses to the smaller area was due to technological constraints as is noted in a later section.

Forest/Nonforest Layer

Because satellite image pixels with different ground covers often have similar spectral signatures, assignment of classes to individual pixels is often probability based. A layer depicting the probability of forest was constructed using a logistic regression model (McRoberts 2006),

$$E(p_i) = f(X_i; \beta) = \frac{\exp \beta_0 + \sum_{j=1}^{12} \beta_j x_{ij}}{1 + \exp \beta_0 + \sum_{j=1}^{12} \beta_j x_{ij}}, \quad (1)$$

where p_i is the probability of forest for the i^{th} pixel, X_i is the vector of the 12 spectral transformations for the i^{th} pixel with x_{ij} being the j^{th} element, the β s are parameters to be estimated, $\exp(\cdot)$ is the exponential function, and $E(\cdot)$ denotes statistical expectation. When estimating the parameters of (1), only data for the 7,920 completely forested or completely nonforested subplots were used. The covariance matrix for the vector of parameter estimates was estimated analytically as,

$$\text{Var}(\hat{\beta}) = (Z' V_e^{-1} Z)^{-1}, \quad (2)$$

where the elements of the matrix Z are defined as,

$$z_{ij} = \frac{\partial f(X_i; \hat{\beta})}{\partial \beta_j},$$

the elements of V_e are defined as,

$$v_{ij} = \sqrt{\hat{p}_i(1 - \hat{p}_i)} \sqrt{\hat{p}_j(1 - \hat{p}_j)} \hat{\rho}_{ij},$$

and $\hat{\rho}_{ij}$ is the spatial correlation among the standardized residuals estimated using a variogram (McRoberts 2006).

The most probable forest/nonforest classification of the imagery is constructed by comparing the probability, \hat{p} , from (1) for each pixel to 0.5: if $\hat{p} \geq 0.5$, the pixel is classified as forest and assigned a numerical value of 1; if $\hat{p} < 0.5$, the pixel is classified as nonforest and assigned a numerical value of 0; however, because the assignment of forest or nonforest to pixels is based on probabilities, it is uncertain whether this procedure correctly assigns forest or nonforest to individual pixels. Forest/nonforest realizations that reflect the uncertainty in the classification were obtained using a four-step procedure designated Procedure A:

- A1. Using the procedure of Kennedy and Gentle (1980: 228-213), generate a vector of random numbers from a multivariate Gaussian distribution with mean 0 and covariance $\text{Var}(\hat{\beta}_i)$ from (2); add these random numbers to the logistic regression model parameter estimates to obtain simulated parameter estimates.
- A2. Using the simulated parameter estimates from step 1 with (1), calculate a probability, \hat{p} , of forest for each pixel.
- A3. For each pixel, generate a random number, r , from a uniform [0, 1] distribution; if $r \leq \hat{p}$, the pixel is designated

forest with a numerical value of 1; if $r > \hat{p}$, the pixel is designated nonforest with a numerical value of 0.

- A4. Repeat steps A1 through A3 many times and calculate the mean and variance of the numerical values assigned to each pixel.

The variance of the forest/nonforest classifications for each pixel is a measure of the uncertainty of the pixel's classification.

Ash Tree Distribution Layer

Because the ash tree distribution layer was for forest land, only data for the 1,953 FIA forest subplots were used in its construction. Using these plot data, an empirical variogram was constructed and an exponential variogram model was fit to the data. An interpolated surface was constructed for which the ash tree count/hectare, A_i , for the i^{th} pixel was estimated as:

$$\hat{A}_i = \sum_{j=1}^{1953} \frac{w_j}{W_i} A_j^o, \quad (3a)$$

where

$$w_{ij} = \begin{cases} \frac{1}{v_{ij}} - \frac{1}{\hat{\lambda}} & \text{if } \frac{v_{ij}}{\hat{\lambda}} \leq 0.95 \\ 0 & \text{if } \frac{v_{ij}}{\hat{\lambda}} > 0.95 \end{cases}; \quad (3b)$$

v_{ij} is the predicted covariance from the variogram model corresponding to the distance, d_{ij} , between the i^{th} pixel and the j^{th} plot; $\hat{\lambda}$ is the estimate of the variogram sill; and

$$W_i = \sum_{j=1}^{1953} w_j. \quad (3c)$$

The variance of \hat{A}_i was estimated as:

$$\text{Var}(\hat{A}_i) = \sum_{j=1}^{1953} \frac{w_j^2}{W_i} (A_j^o - \hat{A}_i)^2, \quad (4)$$

and is a measure of the uncertainty of the estimate, \hat{A}_i . Other equally valid approaches such as kriging could have been used to construct the ash tree county/hectare layer. Realizations of the ash tree count/hectare distribution were obtained by selecting for each pixel a random number from a normal distribution

with mean 0 and variance, $\text{Var}(\hat{A}_i)$ from (4), and adding the number to \hat{A}_i from (3a).

Uncertainty Estimation

Uncertainty in the areal estimate, A_{total} , is due to contributions from four sources: (1) uncertainty in the logistic regression model parameter estimates; (2) uncertainty in the pixel-level forest/nonforest classifications, given a set of parameter estimates; (3) uncertainty in the interpolated pixel-level ash tree count/hectare estimates, \hat{A}_i ; and (4) spatial correlation in forest/nonforest and ash tree observations not accommodated in the logistic regression model predictions and the ash tree count/hectare interpolated estimates.

The spatial correlation contribution to uncertainty in A_{total} results from two phenomena. First, forest areas tend to be clustered rather than independently and randomly distributed throughout the landscape. Thus, to mimic natural conditions, forest/nonforest realizations generated from the logistic regression model predictions of forest probability should exhibit clustering comparable to that observed among the FIA subplot observations of forest and nonforest. This feature requires that the random numbers used to generate the forest/nonforest realization in step A3 be drawn from a correlated uniform [0, 1] distribution. Second, the errors obtained as the differences between A^o and \hat{A} were expected to be spatially correlated; i.e., if an interpolation, \hat{A} , overestimates its true value, other interpolations in close spatial proximity would be expected to overestimate their true values also. For this investigation, however, the range of spatial correlation for the interpolation errors was only slightly more than the 30 m pixel width, regarded as negligible, and ignored.

To generate random numbers from an appropriately correlated uniform [0, 1] distribution as required to accommodate spatial correlation, an eight-step procedure designated Procedure B was used:

- B1. Construct an empirical variogram,

$$\hat{\gamma} = \frac{2}{2n} \sum_{i=1}^{n-1} \sum_{j=i+1}^n (F_j - F_i)^2$$

where F is the numerical designation for forest or nonforest subplot observations, and $n(d)$ denotes a collection of pairs, (F_i, F_j) , whose Euclidean distances in geographic space are within a given neighborhood of d .

B2. Fit an exponential variogram,

$$\hat{\gamma} = 1 - \exp(-\hat{\alpha}_2 d) \quad (5)$$

to the empirical variogram from step B1, where the estimate of the range of spatial correlation is

$$\frac{\ln(0.05)}{\hat{\alpha}_2}$$

B3. Construct a spatial correlation matrix by assigning to each pixel pair, (i, j) , a correlation, r_{ij} , calculated as,

$$\rho_{ij} = \exp(-\gamma d_{ij}),$$

where d_{ij} is the distance between the i^{th} and j^{th} pixel centers and, initially, $\gamma = \hat{\alpha}_2$ from step B2.

B4. Generate a vector of random numbers, one for each pixel, from a multivariate Gaussian distribution with the correlation structure constructed in step B3 using the technique described by Kennedy and Gentle (1980: 228-231).

B5. Convert the Gaussian random numbers from step B4 to Gaussian cumulative frequencies, resulting in a correlated uniform [0, 1] distribution.

B6. Generate a forest/nonforest realization using Procedure A with the correlated uniform [0, 1] distribution from step B5.

B7. Construct an empirical variogram of the forest/nonforest realization from step B5; fit an exponential variogram model; and estimate the range of spatial correlation as in step B2.

B8. Repeat steps B3–B7, adjusting the γ parameter in step B3 each iteration until the range of spatial correlation from step B7 is close to that obtained in step B2.

The exponential variogram model was used in step B2 because of its simplicity and the adequacy of the fit to the data. Construction of the multivariate Gaussian distribution in step B4 requires the Cholesky decomposition of a covariance matrix corresponding to the correlation matrix constructed in step B3. For a square region, n pixels on a side, the correlation matrix

will be $n^2 \times n^2$. Thus, the 30 km x 30 km study area, which has 1,000 TM pixels on a side, would require decomposition of a $10^6 \times 10^6$ matrix. To accommodate personal computer space and processing limitations, analyses involving spatial correlations were constrained to a 2-km x 2-km region, which is approximately 67 pixels on a side and requires decomposition of a smaller 4,489 x 4,489 matrix (figs. 2 and 3).

Uncertainty in A_{total} for the 2-km x 2-km region was estimated using a six-step Monte Carlo simulation procedure designated Procedure C:

- C1. Generate random numbers from a multivariate Gaussian distribution with mean 0 and variance matrix, $\text{Var}(\hat{\beta}_i)$ from (2); add these random numbers to the logistic regression model parameter estimates to obtain simulated parameter estimates; calculate the probability, \hat{p} , of forest for each pixel using the simulated parameter estimates with (1).
- C2. For each pixel, generate a random number, r , from a correlated uniform [0, 1] distribution using Procedure B; if $r \leq \hat{p}$, designate the pixel as forest; if $r > \hat{p}$, designate the pixel as nonforest.
- C3. Calculate the total forest area, F_{total} , as the product of the number of forest pixels from step C2 and the 0.09 ha pixel area.
- C4. For each pixel, generate a random number from a normal distribution with mean 0 and variance, $\text{Var}(\hat{A}_i)$ from (4); add the random number to the interpolated estimate of ash tree count/hectare, \hat{A} , to obtain a simulated ash tree count/hectare; multiply the simulated ash tree count/hectare and the 0.09 hectare pixel area to obtain a simulated ash tree count for the pixel.
- C5. Estimate A_{total} as the sum of the simulated ash tree counts from step C4 for forest pixels from step C2;
- C6. Repeat steps C1 through C5 many times; calculate the mean and variance of F_{total} and A_{total} over all repetitions; estimate the uncertainties in F_{total} and A_{total} as $\hat{\text{Var}}(F_{\text{total}})$ and $\hat{\text{Var}}(A_{\text{total}})$, respectively.

In Procedure C, the contribution of uncertainty due to the logistic regression model parameter estimates may be excluded by not adding uncertainty in step C1; the contribution of uncertainty in the classification, given the parameter estimates, may be

excluded by skipping step C2 and comparing the probabilities generated in step C1 to 0.5 using Procedure A; the contribution of uncertainty due to spatial correlation may be excluded by generating an uncorrelated uniform [0, 1] distribution in step C2; and the contribution of uncertainty due to the interpolated ash tree counts/hectare may be excluded by not adding uncertainty in step C4. The magnitude of the contributions of individual sources of uncertainty may be estimated by considering $\hat{V}\hat{a}r\left[\hat{F}_{total}\right]$ and $\hat{V}\hat{a}r\left[\hat{A}_{total}\right]$ obtained by including contributions from all sources individually and in combinations.

Results

The analysis of Monte Carlo results indicated that all estimates stabilized to within less than 1 percent by 25,000 simulations (table 1). Therefore, 25,000 simulations were used when applying Procedure C to estimate the contributions of the various sources of uncertainty.

The forest/nonforest maps constructed using logistic regression model predictions produced realistic spatial distributions, although no independent accuracy assessment was conducted (fig. 2a); however, considerable detail was revealed in the uncertainty map; e.g., the field structure and road networks (fig. 2b). Considerably less detail was revealed in the ash tree count/hectare map, but this result was expected because of the fewer FIA plots available and the more continuous nature of the layer

(fig. 3a). As expected with biological analyses, the greatest uncertainty in the latter map occurred in the same locations as the greatest estimated values (fig. 3b).

The estimates obtained using Procedure C dramatically revealed that the source of uncertainty making the greatest contribution to uncertainties in the estimates of both F_{total} and A_{total} was spatial correlation in the realizations of the forest/nonforest maps. The magnitude of this effect is highlighted by noting that when uncertainty from this source was included, 95-percent confidence intervals for both F_{total} and A_{total} included, or were close to including, 0. The contribution of the uncertainty in the underlying ash tree count/hectare layer to $\hat{V}\hat{a}r\left[\hat{A}_{total}\right]$ was much less than the contribution due to the uncertainty in the forest/nonforest layer.

Conclusions

Three conclusions may be drawn from this study. First, spatial correlation is a crucial contributor to uncertainty in map analyses that aggregate results over multiple mapping units. Ignoring this contribution inevitably leads to underestimates of variances and unwarranted statistical confidence in estimates. Unfortunately, the importance of this source of uncertainty is generally not known to those who conduct map-based analyses, and techniques for estimating its effects are generally unfamiliar. Second, researchers, authors, and university faculty should

Table 1.—Monte Carlo simulation estimates from Procedure C.

Model parameter estimates	Source of uncertainty			Estimates			
	Classification	Spatial correlation	Ash tree count/ha interpolation	F_{tot}		A_{tot}	
				Mean (ha)	SE* (ha)	Mean (count)	SE* (count)
No uncertainty				113.40		2860.11	
X				110.93	25.37	2801.15	643.23
	X			104.22	1.87	2637.11	47.61
	X	X		104.04	55.32	2541.56	1180.34
X	X	X		105.58	60.89	2633.00	1549.91
			X	113.40		2860.55	47.57
X	X		X	105.53	18.97	2670.79	484.46
X	X	X	X	106.25	61.92	2690.07	1576.47

* SE is standard error and is calculated as the square root of variance.

Figure 2a.—Forest/nonforest classification for 30-km x 30-km study area.

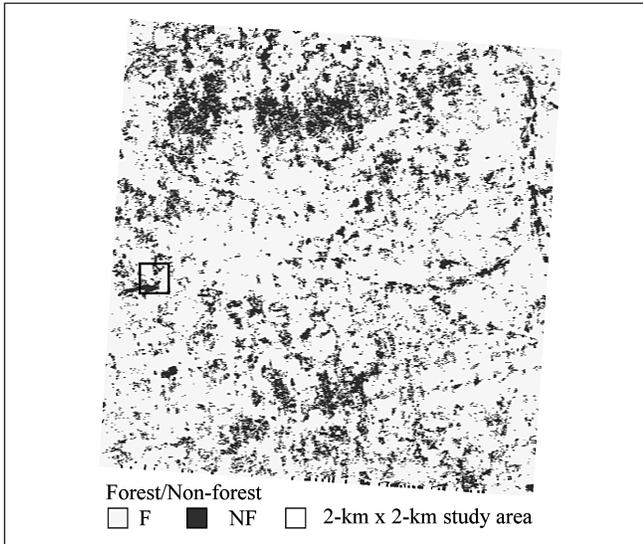


Figure 2b.—Standard errors of pixel forest/nonforest classifications for 30-km x 30-km study area.

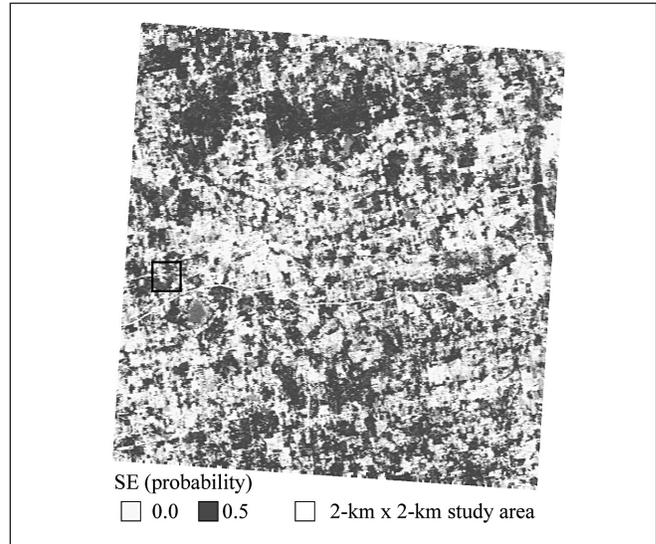


Figure 3a.—Ash tree count/hectare interpolations for 30-km x 30-km study area.

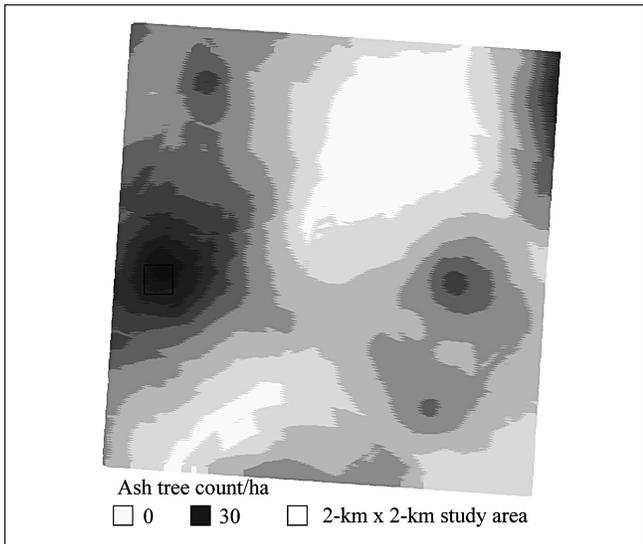
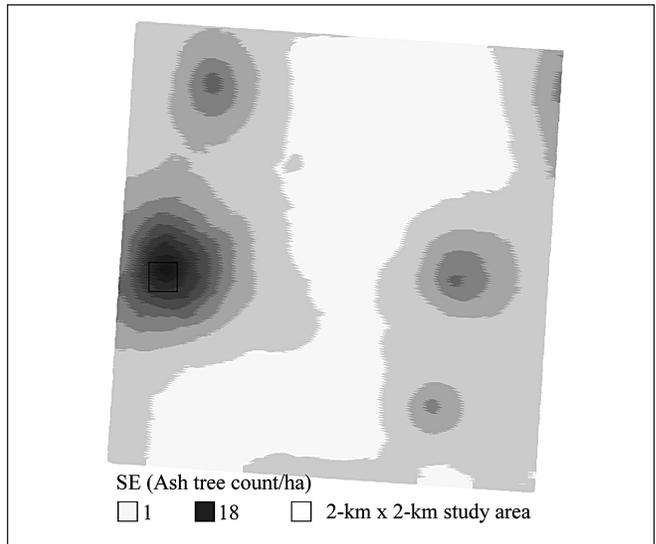


Figure 3b.—Standard error of pixel ash tree count/hectare interpolations for 30-km x 30-km study area.



give greater attention to uncertainty estimation. Third, estimation of uncertainty is not trivial, either conceptually or from a technical perspective. The necessity of decomposing very large matrices limits the size of regions that can be analyzed without high-speed computing facilities.

Literature Cited

Bechtold, W.A.; Patterson, P.L., eds. 2005. The enhanced Forest Inventory and Analysis program—national sampling design and estimation procedures. Gen. Tech. Rep. SRS-80. Asheville, NC: U.S. Department of Agriculture, Forest Service, Southern Research Station. 85 p.

-
- Cappaert, D.; McCullough, D.G.; Poland, T.M.; Siegert, N.W. 2005. Emerald ash borer in North America: a research and regulatory challenge. *American Entomologist*. 51: 152-165.
- Crist, E.P.; Cicone, R.C. 1984. Application of the tasseled cap concept to simulated Thematic Mapper data. *Photogrammetric Engineering and Remote Sensing*. 50: 343-352.
- Haack, R.A.; Jendek, E.; Liu, H.; Marchant, K.R.; Petrice, T.R.; Poland, T.M.; Ye, H. 2002. The emerald ash borer: a new exotic pest in North America. *Newsletter of the Michigan Entomological Society*. 47: 1-5.
- Herms, D.A.; McCullough, D.G.; Smitley, D.R. 2004. Under attack. *American Nurseryman*. <http://www.emeraldashborer.info/educational.cfm>.
- Kauth, R.J.; Thomas, G.S. 1976. The tasseled cap—a graphic description of the spectral-temporal development of agricultural crops as seen by Landsat. *Proceedings, symposium on machine processing of remotely sensed data*. West Lafayette, IN: Purdue University: 41-51.
- Kennedy, W.J., Jr.; Gentle, J.E. 1980. *Statistical computing*. New York: Marcel Dekker. 591 p.
- McCullough, D.G.; Katovich, S.A. 2004. Emerald ash borer. NA-PR-02-04. Newtown Square, PA: U.S. Department of Agriculture, Forest Service, State and Private Forestry, Northeastern Area. 2 p.
- McRoberts, R.E. 2006. A model-based approach to estimating forest area. *Remote Sensing of Environment*. 103: 56-66.
- McRoberts, R.E.; Bechtold, W.A.; Patterson, P.L.; Scott, C.T.; Reams, G.A. 2005. The enhanced Forest Inventory and Analysis program of the USDA Forest Service: historical perspective and announcement of statistical documentation. *Journal of Forestry*. 103(6): 304-308.