

¹

² Classification and mapping of low-statured shrubland cover ³ types in post-agricultural landscapes of the US Northeast

⁴ Michael J Mahoney¹ , Lucas K Johnson¹ , Abigail Guinan¹ , Colin M Beier² 

⁵ ARTICLE HISTORY

⁶ Compiled December 21, 2022

⁷ ¹ Graduate Program in Environmental Science, State University of New York College
⁸ of Environmental Science and Forestry, Syracuse, NY, USA

⁹ ² Department of Sustainable Resources Management, State University of New York
¹⁰ College of Environmental Science and Forestry, Syracuse, NY, USA

¹¹ ABSTRACT

¹² Novel plant communities reshape landscapes and pose challenges for land cover
¹³ classification and mapping that can constrain research and stewardship efforts. In
¹⁴ the US Northeast, emergence of low-statured woody vegetation, or shrublands, in-
¹⁵ stead of secondary forests in post-agricultural landscapes is well-documented by
¹⁶ field studies, but poorly understood from a landscape perspective, which limits the
¹⁷ ability to systematically study and manage these lands. To address gaps in classifi-
¹⁸ cation/mapping of low-statured cover types where they have been historically rare,
¹⁹ we developed models to predict shrubland distributions at 30m resolution across
²⁰ New York State (NYS), using a stacked ensemble combining a random forest, gra-
²¹ dient boosting machine, and artificial neural network to integrate remote sensing of
²² structural (airborne LIDAR) and optical (satellite imagery) properties of vegetation

CONTACT: Michael J Mahoney. Email: mjmahone@esf.edu. Lucas K Johnson. Email: ljohns11@esf.edu.

Abigail Guinan. Email: aguinan@esf.edu. Colin M Beier. Email: cbeier@esf.edu.

Please cite as: Mahoney, M. J., Johnson, J. K., Guinan, A. Z., and Beier, C. M. 2022. Classification and mapping of low-statured 'shrubland' cover types in post-agricultural landscapes of the US Northeast. The International Journal of Remote Sensing, 43(19-24), 7117-7138. <https://doi.org/10.1080/01431161.2022.2155086>

23 cover. We first classified a 1m canopy height model (CHM), derived from a patchwork
24 of available LiDAR coverages, to define shrubland presence/absence. Next, these
25 non-contiguous maps were used to train a model ensemble based on temporally-
26 segmented imagery to predict shrubland probability for the entire study landscape
27 (NYS). Approximately 2.5% of the CHM coverage area was classified as shrubland.
28 Models using Landsat predictors trained on the classified CHM were effective at
29 identifying shrubland (test set AUC=0.893, real-world AUC=0.904), in discriminating
30 between shrub/young forest and other cover classes, and produced qualitatively
31 sensible maps, even when extending beyond the original training data. After ground-
32 truthing, we expect these shrubland maps and models will have many research and
33 stewardship applications including wildlife conservation, invasive species mitigation
34 and natural climate solutions. Our results suggest that incorporation of airborne
35 LiDAR, even from a discontinuous patchwork of coverages, can improve land cover
36 classification of historically rare but increasingly prevalent shrubland habitats across
37 broader areas.

38 **KEYWORDS**

39 LiDAR; Landsat; shrubland; machine learning; neural networks; land cover

40 **1. Introduction**

41 Human land use has fundamentally altered vegetation-environment relationships and
42 created legacies that include the emergence of novel communities and ecosystem
43 types (Foster, Motzkin, and Slater 1998; Cramer, Hobbs, and Standish 2008). In post-
44 agricultural landscapes of eastern North America, these legacies include loss of plant
45 diversity (Flinn and Vellend 2005) and widespread homogenization of vegetation compo-
46 sition and structure, relative to historical reconstructions (Foster, Motzkin, and Slater
47 1998; Flinn, Vellend, and Marks 2005). Widespread abandonment of crop, pasture and
48 industrial lands from the late-19th to middle-20th centuries created an expanding land
49 base for invasion and emergence of novel communities (Williams and Jackson 2007;
50 Fridley 2012; Alexander, Diaz, and Levine 2015), with variable outcomes depending on
51 prior land use practices (Stover and Marks 1998; Benjamin, Domon, and Bouchard 2005;
52 Kulmatiski, Beard, and Stark 2006). Entirely novel communities have emerged in old-
53 fields due to colonization by non-native plants, including invasive woody shrubs, which
54 are much more likely to establish and become dominant in post-agricultural (Johnson

55 et al. 2006; Cramer, Hobbs, and Standish 2008; McCay and McCay 2009) and post-
56 industrial sites (Spiering 2019) compared to closed-canopy forests of any successional
57 age. Meanwhile, secondary forests across the US Northeast, including those established
58 in old fields, often lack sufficient advance regeneration to maintain productivity and
59 resilience to changing disturbance regimes (Dey et al. 2019).

60 Among the outcomes of these changes, the emergence of low-statured vegetation or
61 shrublands as a more common cover type in the US Northeast has been suggested by
62 numerous field studies, but is poorly understood from a landscape perspective. Here
63 the term shrubland reflects a physiognomic definition following King and Schlossberg
64 (2014), which encompasses several types of plant communities found in the US North-
65 east, including: 1) young, regenerating or otherwise low-statured closed-canopy forests;
66 2) wetlands dominated by native shrubs (e.g., *Alnus* spp) or small-statured trees (e.g.,
67 *Picea mariana* in boreal peatlands); 3) uplands dominated by native shrubs (e.g., *Cor-*
68 *nus racemosa*); and most recently, 4) upland shrub/scrub dominated by invasive woody
69 (e.g., *Rhamnus* and *Lonicera* spp) and herbaceous (e.g., *Solidago* spp) perennials. Among
70 the types above, regenerating forests and invasive shrub/scrub communities are of grow-
71 ing interest for research and management purposes, given their anthropogenic origins,
72 their potential novelty in terms of composition and dynamics, and their implications for
73 biodiversity and ecosystem services (Cramer, Hobbs, and Standish 2008; Hobbs, Higgs,
74 and Harris 2009; Perring, Standish, and Hobbs 2013; Dey et al. 2019). Despite their
75 conservation value as wildlife habitat, especially for songbirds, shrublands are widely
76 unpopular cover types in terms of their perceived aesthetic, recreational and economic
77 values (Askins 2001; King and Schlossberg 2014). Although long disregarded, these lands
78 are rapidly gaining attention in today's urgent push to implement natural climate solu-
79 tions (Fargione et al. 2018) and identify marginal or underutilized lands for renewable
80 energy generation.

81 However, current limitations to the classification and mapping of these cover types pose
82 obstacles to advancing both science and stewardship opportunities (Hobbs, Higgs, and
83 Harris 2009). Shrublands are a very challenging cover class to identify from imagery
84 alone, given the breadth of community types included (as noted above) and the high

variability in density and canopy cover that exists within and among those community types (King and Schlossberg 2014). In practical terms this means that, when relying solely on imagery, shrublands encompass a full gradient from resembling herbaceous or barren land to resembling closed-canopy conditions (Brown et al. 2020). Furthermore, shrublands are relatively rare in the US Northeast, making them particularly challenging to identify with standard classification methods (Bogner et al. 2018; Haibo He and Garcia 2009). As a result, imagery-based approaches tend to classify shrubland categories with substantially lower accuracy than other land cover classes (Wickham et al. 2021; Brown et al. 2020).

A solution for this problem might be to incorporate additional, non-imagery sources of remote sensing data into land cover classification methodologies. LiDAR data collected through airborne laser scanning can provide essential information for identifying low-statured vegetation such as early-successional forests (Falkowski et al. 2009). In combination with imagery, LiDAR data can enable continuous, broad-scale estimation of canopy heights and other structural traits which may greatly simplify the task of distinguishing between low-statured and taller closed-canopy cover types (Ruiz et al. 2018). Unfortunately, the cost and logistical challenges of airborne LiDAR collection have constrained its availability to smaller extents and with much longer return intervals than provided by satellite imagery. Yet if canopy structural estimates from airborne LiDAR could be used to label a training dataset in order to fit models using satellite imagery, it should be possible to produce models capable of identifying shrubland with greater accuracy than those trained on imagery alone, while being able to map/model a larger and more contiguous spatial extent than models relying on airborne LiDAR data as predictors. Similar methods have been used to automate labeling imagery used to train models for tree and building detection (Zarea and Mohammadzadeh 2016; Huang et al. 2019), but to our knowledge this data fusion approach has not yet been applied for shrubland mapping.

Here we explored this approach to map shrubland across New York State (NYS), a largely post-agricultural landscape containing extensive old-fields populated mostly by secondary forests or invasive shrub/scrub communities, which are difficult to parse based

on imagery alone. We leveraged a non-contiguous patchwork of existing large-footprint LIDAR data sets by creating very high resolution (1 m) canopy height models that covered approximately 60% of the study area (NYS). By sampling from these canopy height maps, we trained machine learning models with temporally segmented Landsat imagery and land cover data products, and created a stacked ensemble to map the probability of shrubland at a high resolution (30 m) across the study area. Models were evaluated across multiple sensitivity/specificity thresholds to generate a range of map outputs that can support future ground-truthing efforts as part of research and stewardship activities. We provide new maps of marginal cover types, such as invasive shrublands and degraded young forests, and demonstrate how to leverage an information-rich but geographically incomplete data source (LIDAR) for large-scale contiguous land cover classification and mapping based on widely available and standardized time-series imagery. Our results suggest that using LiDAR, where it exists, to label data for model training may help distinguish shrublands from other, optically similar land cover classes. We additionally demonstrate the utility of using targeted, regional land cover products to supplement national products for classes of regional interest.

2. Methods

In order to map potential shrubland across New York State, we first identified low-stature vegetation across a discontinuous temporal patchwork of 1m resolution canopy height models (CHMs) derived from 19 distinct LiDAR acquisitions. We then aggregated these 1m identifications into a 30m resolution raster surface, associated this surface with multiple data products derived from temporally matching remote sensing observations (including Landsat imagery as well as climate and topographic data), and used this data to produce a stacked ensemble model composed of a random forest, gradient boosting machine, and artificial neural network combined through a logistic regression. We used this ensemble model to predict shrubland for a spatiotemporal patchwork matching LiDAR acquisitions, as well as for the entire state in 2019. A flowchart of this process is included as Figure 1.

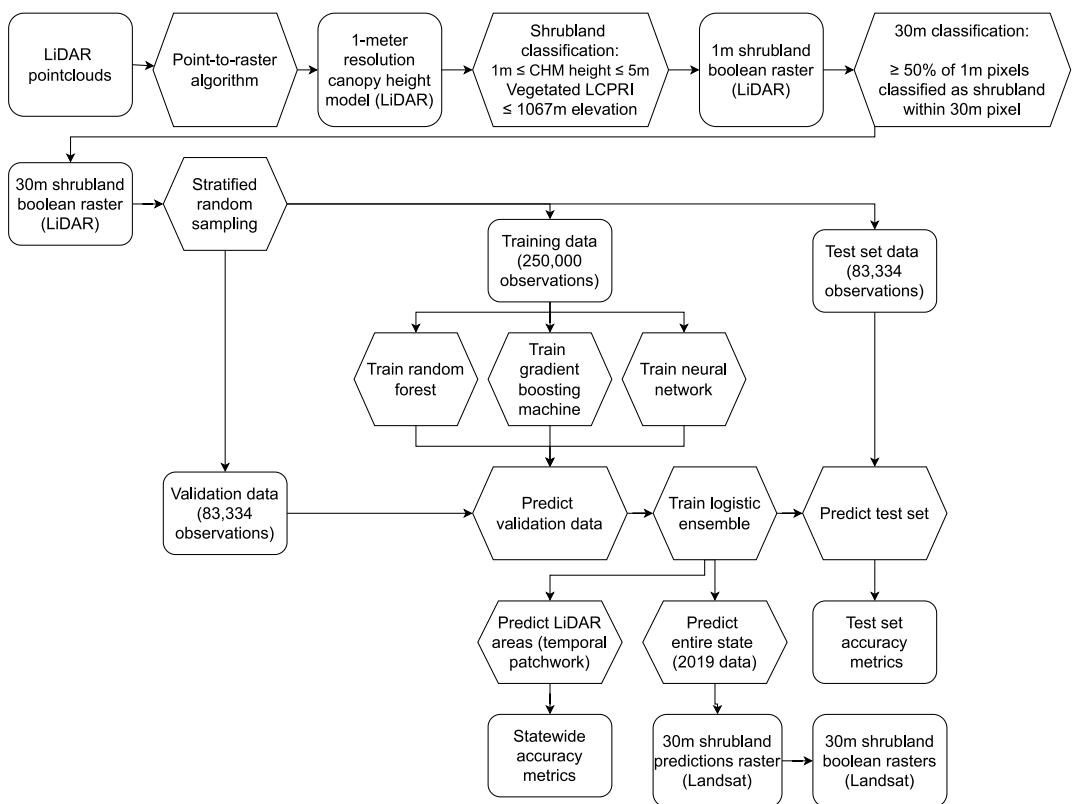


Figure 1.: A flowchart diagram showing the key elements of the identification methodology. Rectangular boxes represent data products and results, while hexagonal boxes represent methodological steps.

¹⁴³ **2.1. Study area**

¹⁴⁴ New York State spans an area of 141,297 km² of the northeastern United States. Ex-
¹⁴⁵ tensive land clearing for agriculture and industry in the 18th through 19th centuries
¹⁴⁶ decimated forests throughout the region, with forest cover dropping to 10-30% of the
¹⁴⁷ landscape by 1880 (Lorimer 2001). While total forest cover recovered rapidly at the turn
¹⁴⁸ of the 20th century, these forests were almost entirely young due to the combination
¹⁴⁹ of regeneration on abandoned farmland with the continual coppicing and harvesting of
¹⁵⁰ more established woodlots for fuelwood and other products (Whitney 1994). A variety
¹⁵¹ of factors, among them a decrease in the use of wood for residential heating and an
¹⁵² increase in forest conservation efforts, caused these forests to begin to mature in the
¹⁵³ 1930s, with the effect that the majority of forest stands across the Northeast are now
¹⁵⁴ over 100 years old. Two forest preserves, the Adirondack Park in the northeast and the
¹⁵⁵ Catskill Park in the southeast, have been protected in the New York State constitu-
¹⁵⁶ tion as “forever wild”; timber harvesting has been generally prohibited on state-owned
¹⁵⁷ parcels in these regions since they were incorporated into the Forest Preserve, a process
¹⁵⁸ beginning in 1885. As a result, there is very little shrubland in these preserves.

¹⁵⁹ Elevations across New York State range from -2 m to 1,584 m above sea level (U.S. Geo-
¹⁶⁰ logical Survey 2019), with daily temperatures in 2019 ranging from -17 °C to 28 °C and
¹⁶¹ monthly precipitation for the same period ranging from 5.0 cm to 16.8 cm (NOAA Na-
¹⁶² tional Centers for Environmental Information 2022). The majority of the state occupies
¹⁶³ the northern hardwoods-hemlock forest region, though there are important inclusions of
¹⁶⁴ beech-maple-basswood and Appalachian oak communities in the western and southern
¹⁶⁵ reaches of the state, respectively (Dyer 2006).

¹⁶⁶ **2.2. LiDAR Data and Shrubland Identification**

¹⁶⁷ Although a distinction exists between early-successional and young forest habitats in
¹⁶⁸ eastern North America, for simplicity we have followed King and Schlossberg (2014) in
¹⁶⁹ combining these categories into a single shrubland classification due to their structural
¹⁷⁰ similarity. This terminology aligns with Anderson (1976), who described shrubland in
¹⁷¹ the eastern United States as “former croplands or pasture lands (cleared from original

172 forest land) which now have grown up in brush in transition back to forest land to the
173 extent that they are no longer identifiable as cropland or pasture from remote sensor
174 imagery.”

175 Shrubland was identified using 19 distinct leaf-off LiDAR acquisitions, collected and
176 made freely available as part of the US Geological Survey’s 3D Elevation program. All
177 LiDAR used in this study was collected between 2014 and 2020, with point densities
178 ranging from 1.98 to 3.24 points per square meter (Figure 2). All LiDAR data had a ver-
179 tical accuracy RMSE of \leq 10 cm. While horizontal accuracy was not typically reported
180 in provided LiDAR metadata, horizontal RMSE for all data sets is expected to be \leq
181 68 cm (ASPRS 2014). More information about individual LiDAR coverages is available
182 as Supplementary Materials 1. These point clouds were converted to 1 m canopy height
183 models (CHMs) using a point-to-raster algorithm implemented in the lidR R package
184 (Roussel et al. 2020). To reflect the 2D nature of a LiDAR return footprint, and mit-
185 igate potential voids in the resulting CHM, each return was replaced with a circle of
186 returns with a diameter equal to the pulse width present in the metadata (default 0.5
187 m; Supplementary Materials 1). These CHMs were then masked to exclude any pixel
188 assigned a non-vegetation primary land cover classification by the temporally-matching
189 USGS LCMAP land cover product (namely the developed, water, ice and snow, and
190 barren classes) using the terra R package (Brown et al. 2020; Hijmans 2021; R Core
191 Team 2021). CHMs were also masked to exclude any pixel with an elevation above
192 1067m (3500 ft), as shorter canopy heights at these elevations likely represent krumh-
193 holz (stunted trees near elevational treeline) instead of shrubs or regenerating forest.
194 Shrubland was then defined as any 1 m pixel with a CHM height between 1-5 m. The
195 lower threshold was defined to avoid classifying cropland and human structures (such as
196 low walls) as shrubland, and the upper threshold defined to match the USGS NLCD def-
197 inition of shrubland, which was derived from the Anderson classification system (Yang
198 et al. 2018; Anderson et al. 1976). While this coarse definition allows for efficient and
199 automated labeling of data, it also has its limitations, including being reliant upon the
200 accuracy and unbiasedness of LCMAP class predictions and potentially incorrectly clas-
201 sifying some cropland, including orchards, as shrublands. To evaluate these limitations,
202 we validated model predictions using manual inspection of optical imagery, described

203 further in Section 2.4.

204 These 1 m pixels were then aggregated into 30 m resolution pixels using GDAL
205 (GDAL/OGR contributors 2021), with each 30 m pixel defined as shrubland if more
206 than 50% of its constituent 1 m subpixels were identified as shrubland.

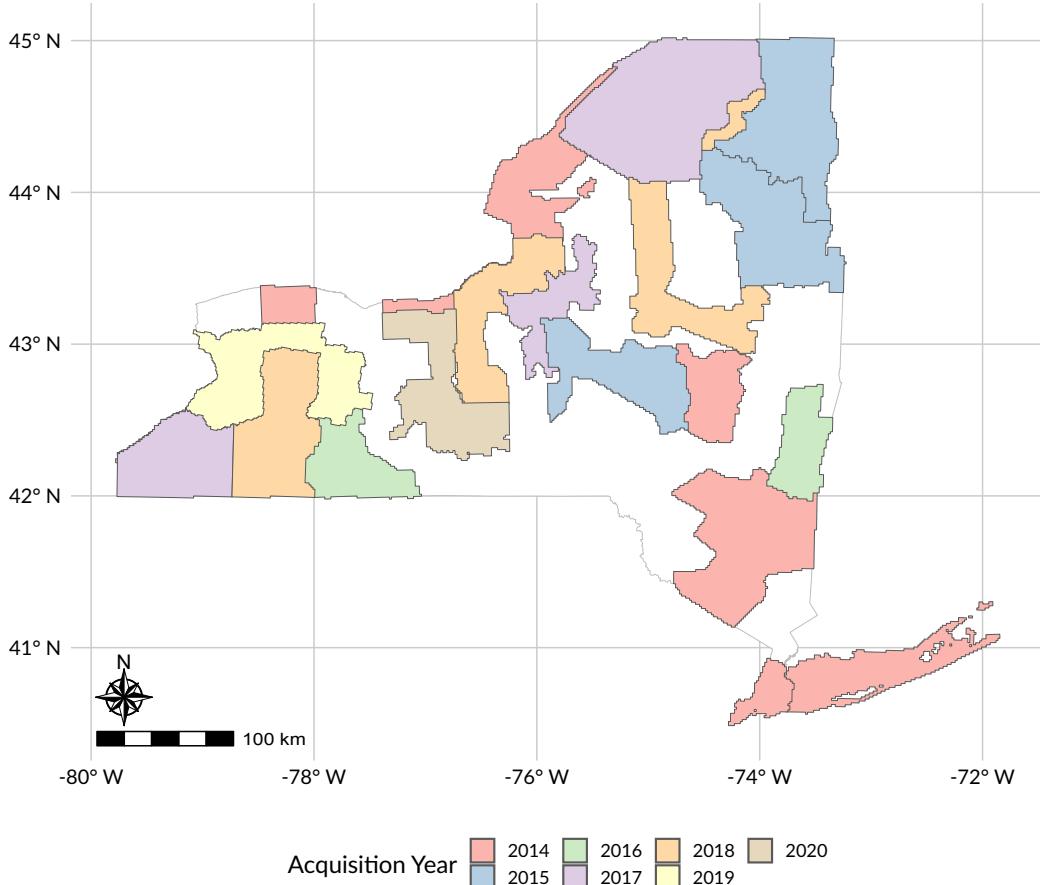


Figure 2.: Boundaries for all LiDAR coverages used in this project, colored by year of data acquisition. More information about each coverage is included as Online Resource 1.

207 2.3. Predictor Creation

208 We produced a set of 10 annual Landsat-derived predictors by processing Landsat anal-
209 ysis ready data (ARD, Dwyer et al. 2018) in Google Earth Engine (GEE, Gorelick et al.
210 2017), using a medoid composite of imagery acquired between July 1st and September
211 1st of the same year LiDAR was acquired at each pixel. Landsat ARD was processed us-
212 ing the LandTrendr implementation in GEE (hereafter LT-GEE) to fill gaps (e.g., clouds

and shadows) in the annual time series, smooth interannual variations (noise), and quantify the disturbance history for each pixel (Kennedy, Yang, and Cohen 2010; Kennedy et al. 2018). The LT-GEE predictors included three tasseled cap indices, (brightness - TCB, greenness - TCG, and wetness - TCW) and their respective deltas computed with a 1-year lag, all fit to Normalized Burn Ratio (NBR) temporally segmented vertices, as well as an NBR index and respective 1-year delta (Kauth and Thomas 1976; Cocke, Fulé, and Crouse 2005; Kennedy et al. 2018). We also processed a separate NBR segmented time-series with LT-GEE parameters tailored to be more sensitive to the timing of discrete disturbances using LT-GEE code to produce two predictors describing disturbances at an individual pixel: year-of-most-recent-disturbance (YOD; 1985-2020) and associated magnitude-of-most-recent-disturbance (MAG; unitless measure of change in NBR value) (Kennedy et al. 2018). All predictors are described in Table 1. The 8 annual indices and associated deltas aimed to capture the surface reflectance for a given pixel at a given time, while the two disturbance predictors aimed to describe disturbance history for a given pixel. NBR was chosen as the base predictor, providing disturbance history information and temporal break-points to which tasseled cap predictors were fit, as it has been shown to be the most sensitive for capturing disturbance events (Kennedy, Yang, and Cohen 2010). Information on LT-GEE parameters used is included as Supplementary Materials 2.

Table 1.: Definitions of predictors used for model fitting.

Predictor	Definition
TCB, TCW, TCG	Tassled cap brightness, wetness, and greenness, with noise removed using LT-GEE
NBR	Normalized burn ratio with noise removed using LT-GEE
MAG, YOD	Magnitude and year of most recent disturbance, as identified using LT-GEE
PRECIP, TMAX, TMIN	30-year normals for precipitation, maximum temperature, and minimum temperature, derived from annual PRISM climate models
ASPECT, ELEVATION, SLOPE, TWI	Aspect, elevation, slope, and topographic wetness index derived from a 30-meter digital elevation model
LCSEC	LCMAP secondary land cover classification

232 In addition to Landsat-derived predictors, a set of steady-state ancillary predictors was
233 included to represent geospatial variation in climate and topography (Kennedy et al.
234 2018). These predictors included precipitation and temperature 30 year normals derived
235 from PRISM Climate Group data (PRISM Climate Group 2022), the secondary land
236 cover classification prediction from LCMAP (Brown et al. 2020), and elevation, aspect,
237 slope, and topographic wetness indices derived from a 30 meter digital elevation model
238 (Mahoney, Beier, and Ackerman 2022; U.S. Geological Survey 2019; Beven and Kirkby
239 1979). In total, models were fit using 14 separate predictors (Table 1).

240 **2.4. Model Fitting and Evaluation**

241 Each LiDAR-derived shrubland layer was combined with a set of temporally matching
242 predictors, which were then merged into a single temporal patchwork data set repre-
243 senting each region of the state during its year of LiDAR acquisition. A total of 416,668
244 pixels were then sampled from this patchwork set (Figure 1), stratified so that half
245 (208,334) represented shrubland and half other land cover types. These pixels were then
246 split at random into a training set of 250,000 pixels, a validation set of 83,334 pixels,
247 and a hold-out evaluation set of 83,334 pixels. We then fit three separate models, a
248 random forest (Breiman 2001), stochastic gradient boosting machine (Friedman 2002),
249 and deep neural network (LeCun, Bengio, and Hinton 2015), against the training data
250 set to estimate the probability of a given pixel representing shrubland. Models were fit
251 using hyperparameters chosen to minimize out-of-sample binary cross entropy (Good
252 1952).

253 The random forest was fit using 3,000 trees, each using a random bootstrap sample of
254 20% of the data, sampled with replacement, a minimum node size of 6 observations, a
255 single random variable per split, and splitting to minimize Gini impurity. The gradient
256 boosting machine was fit using 2,500 trees, each with unlimited depth, a maximum of 14
257 leaves, a learning rate of 0.01, a minimum of 10 observations per leaf and 3 observations
258 per bin, an L1 regularization constant of 0 and a L2 regularization constant of 0.5. Each
259 tree was fit using a new bootstrap sample with half the number of observations as the
260 training data, with access to 90% of predictors. The neural net was a simple additive

261 neural network with seven layers: five densely connected layers with numbers of nodes
262 halving with each additional layer, decreasing from 256 to 128 to 64 to 32 to 16, then
263 feeding into a 20% dropout layer before the final densely connected output layer with
264 a single node. All dense layers used rectified linear activation functions save for the
265 output layer, which used a sigmoid activation function. The model was trained using
266 1,000 epochs, but final weights used the epoch which maximized the area under the
267 precision-recall curve of the validation data set.

268 We then used each of these models to predict the probability of a pixel representing
269 shrubland for each observation in the validation data set. Next, we fit a logistic regres-
270 sion to the validation set predictors and predicted probabilities to combine our three
271 models into a single stacked ensemble model (Wolpert 1992; Dormann et al. 2018). This
272 ensemble model was then used to generate predictions for the test set, for the temporal
273 patchwork data set, and for data reflecting the entire state for 2019 (chosen in order to
274 compare predictions to the 2019 NLCD land cover map). The same model was used for
275 predicting each of these data sets.

276 Probability thresholds used to classify individual pixels were chosen using model pre-
277 dictions for the validation set. Four separate thresholds were identified: the one that
278 maximized the model's summed sensitivity (Equation 1) and specificity (Equation 2),
279 calculated using Youden's J statistic (hereafter 'Youden Optimal') (Youden 1950); and
280 three that maximized sensitivity while keeping specificity above 90%, 95%, and 99%. All
281 accuracy metrics were assessed using the temporal patchwork data set, with Landsat-
282 derived predictors temporally matched to the LiDAR data set. All metrics were calcu-
283 lated considering shrubland pixels as positive cases; higher specificity targets reflected
284 the relatively rare abundance of shrubland throughout the state (approximately 2.5% of
285 mapped pixels) necessitating low levels of false positives. Predictions against both the
286 test set and the temporal patchwork data set were classified using each of these thresh-
287 olds, then assessed using sensitivity (Equation 1), specificity (Equation 2), precision
288 (Equation 3), and F1 score (Equation 4).

$$\text{Sensitivity} = \frac{T_{\text{Positives}}}{T_{\text{Positives}} + F_{\text{Negatives}}} \quad (1)$$

$$\text{Specificity} = \frac{T_{\text{Negatives}}}{F_{\text{Positives}} + T_{\text{Negatives}}} \quad (2)$$

$$\text{Precision} = \frac{T_{\text{Positives}}}{T_{\text{Positives}} + F_{\text{Positives}}} \quad (3)$$

$$F1 = \frac{2 * \text{Precision} * \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}} \quad (4)$$

289 Where $T_{\text{Positives}}$ is the number of pixels correctly classified as shrubland, $F_{\text{Positives}}$ the
 290 number of pixels incorrectly classified as shrubland, $T_{\text{Negatives}}$ the number of pixels
 291 correctly classified as not being shrubland, and $F_{\text{Negatives}}$ the number of pixels incorrectly
 292 classified as not being shrubland.

293 Models were additionally assessed using the area under the receiver operating charac-
 294 teristic curve (AUC) (Austin and Steyerberg 2012), calculated for the test set using all
 295 observations and for the temporal patchwork set using a random sample of 1,000,000
 296 pixels due to computational limitations. Given the imbalance of our classes, we do
 297 not report overall accuracy or balanced accuracy for either the test set or the temporal
 298 patchwork, given that approximately 97.5% overall accuracy could be achieved by never
 299 predicting shrubland.

300 A final assessment involved comparing a stratified sample of model predictions to 1m
 301 resolution 2019 imagery from the National Agricultural Inventory Program (US Depart-

ment of Agriculture 2019). Predictions from the 2019 ensemble model were stratified spatially into 51 16,974 km² hexagons (equaling an apothem of 70 kilometers). Within each hexagon, five pixels were randomly selected from each of 12 probability bins (predicted probabilities from the ensemble model between 0% and 5%, 5% and 10%, 90% and 95%, and 95% and 100%, as well as each decile between 10% and 90%). Samples taken for two additional bins, ranging from 0% to 1% and 99% to 100%, were merged into the 0% to 5% and 95% to 100% bins due to the rarity of these extreme probabilities, resulting in these bins having slightly more observations. When regions of a hexagon extended beyond the mapped area, the number of pixels selected per bin was scaled in proportion to the mapped region of the hexagon. Pixels were classified as either clearly representing shrubland, clearly representing a non-shrubland class, or as “unknown” if the correct classification could not be ascertained from imagery. A single rater classified all validation pixels.

All models were fit using either R version 4.1.2 (R Core Team 2021) or Python version 3.9.10 (Python Core Team 2022). Random forests were fit using the ranger R package (Wright and Ziegler 2017), gradient boosting machines using lightgbm (Ke et al. 2017), and neural nets using keras (Chollet 2015).

3. Results

3.1. LiDAR Classification

Based on LIDAR-derived CHMs, approximately 2.5% of the study area was initially mapped as shrubland (1-5 m tall), representing about 1.83×10^6 ha of the 73.3×10^6 ha land area that remained after LCMAP masking removed non-vegetated cover types (Figure 3). Shrubland was identified in every LiDAR coverage, with proportions ranging from 0.3% (Great Gully) to 7.6% (Great Lakes) of the coverage footprint area (Figure 3).

Shrubland cover was present in each vegetated LCMAP primary classification classes, but was weighted more heavily towards areas classified as cropland or tree cover (Table 2). Approximately 7.3% of land classified by LCMAP as “grassland/shrub” was classified as shrubland through this process, though this result should not be over-

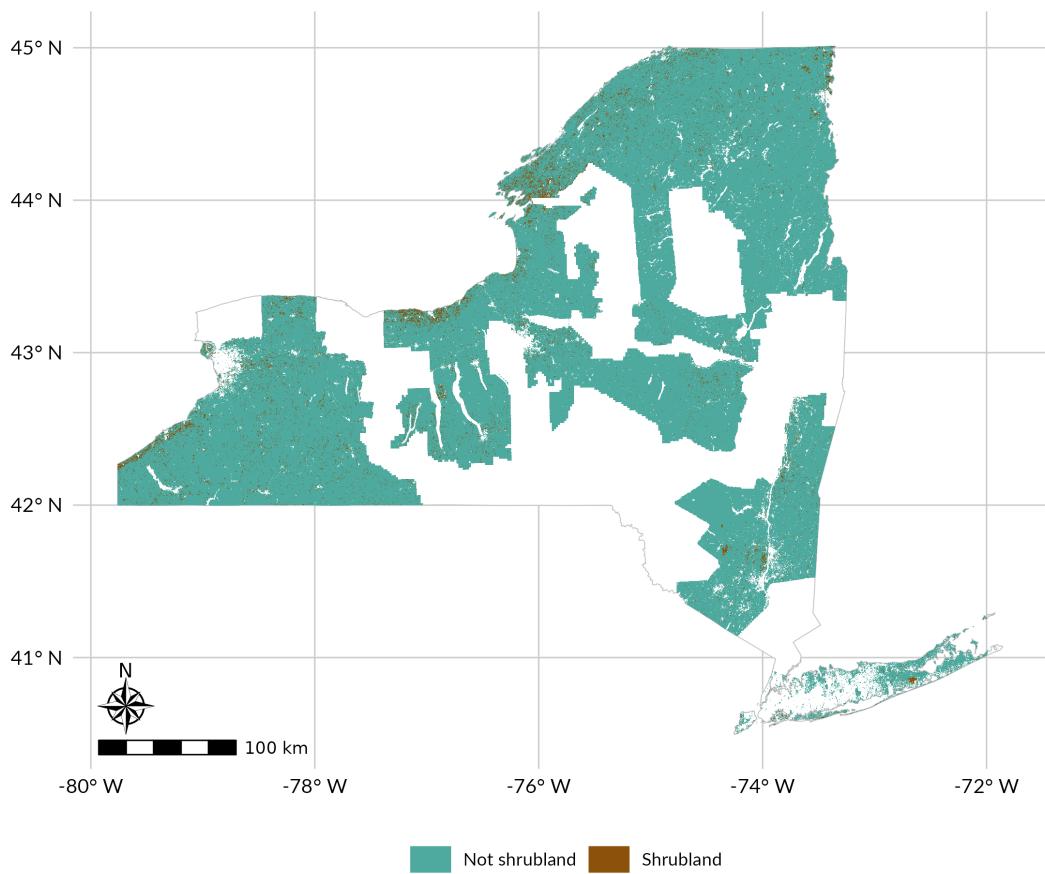


Figure 3.: Identified shrubland areas within each available LiDAR coverage. Shrubland was defined at a 1 meter resolution as being any area within a vegetated LCPRI land cover class and below 1067 meters elevation with a LiDAR-derived height between 1 and 5 meters. 30 meter pixels, used for analysis and modeling, were then defined as shrubland if more than 50% of their contained 1 meter pixels were classified as shrubland. In total, approximately 2.5% of 30 meter pixels were classified as shrubland.

Table 2.: Total area and amount of shrubland within each LCMAP primary land cover classification (“LCPRI”), in square kilometers, for each map surface (bolded text). Both the temporal patchwork and statewide models were classified using a 95% specificity threshold.

LCPRI	Total area	Shrubland area	% Shrubland
LiDAR Classification			
Cropland	20 077.6	679.5	3.4%
Grass/Shrub	2 085.3	151.8	7.3%
Tree Cover	44 879.5	597.2	1.3%
Wetland	6 635.0	467.9	7.1%
Temporal Patchwork			
Cropland	20 077.6	1 451.2	7.2%
Grass/Shrub	2 085.3	416.2	20.0%
Tree Cover	44 879.5	1 428.3	3.2%
Wetland	6 635.0	1 166.3	17.6%
2019 Statewide			
Cropland	29 922.9	2 005.9	6.7%
Grass/Shrub	2 994.2	549.3	18.3%
Tree Cover	69 419.0	1 694.6	2.4%
Wetland	9 008.3	1 425.0	15.8%

330 interpreted given the LCMAP “grassland/shrub” category includes herbaceous land
 331 covers alongside shrublands (Brown et al. 2020).

332 **3.2. Model Accuracy**

333 The ensemble model had a higher AUC than the individual component models on both
 334 the test set (Ensemble 0.893; random forest 0.843; gradient boosting machine 0.889;
 335 neural network 0.883) and a random sample of 1×10^6 pixels from the LiDAR temporal
 336 patchwork data set (Ensemble 0.904; random forest 0.844; gradient boosting machine
 337 0.890; neural network 0.901) (Figure 4; Supplementary Materials 3). As a result, we
 338 focus here on accuracy assessments for the ensemble model; accuracy assessments for
 339 the component models are included as Supplementary Materials 3.

340 When evaluating models against the balanced test set, the highest F1 score (0.816)
 341 was achieved using the Youden-optimal classification threshold (Table 3). The model
 342 retained its high AUC, sensitivity, and specificity when predicting the LiDAR temporal
 343 patchwork data set, precision was lower due to the large imbalance between shrubland

Table 3.: Model accuracy metrics for logistic ensemble model with predictions classified using various thresholds, calculated using both the balanced test set and the LiDAR patchwork surface. AUC for the LiDAR patchwork was calculated using a random sample of 1,000,000 pixels, while all other metrics used all predicted pixels. Thresholds were selected using a separate validation set, using values chosen to maximize the Youden J statistic ('Youden optimal') or to target a certain minimum specificity ('% specificity').

	Threshold	Sensitivity	Specificity	Precision	F1
Test set (AUC: 0.893)					
Youden optimal	0.489	0.842	0.780	0.791	0.816
90% specificity	0.755	0.659	0.900	0.867	0.807
95% specificity	0.840	0.496	0.949	0.906	0.641
99% specificity	0.907	0.218	0.989	0.952	0.355
LiDAR patchwork (AUC: 0.904)					
Youden optimal	0.489	0.858	0.783	0.094	0.169
90% specificity	0.755	0.689	0.896	0.149	0.245
95% specificity	0.840	0.514	0.951	0.219	0.307
99% specificity	0.907	0.247	0.989	0.376	0.298

344 and other cover classes across the state. As a result, the model attained its highest F1
 345 score of 0.307 when using a classification threshold to that targeted 95% specificity.

346 *3.2.1. LiDAR Patchwork Predictions*

347 When predicting the temporal patchwork, our model predicted the highest probabilities
 348 of shrubland along the northern reaches of the state, matching the distribution of shrub-
 349 land in the true LiDAR-derived surface (Figure 5). However, the model also predicted
 350 higher than average probabilities of shrubland in the southwestern and central regions
 351 of the state, neither of which were reflected in the original LiDAR-derived surface.

352 Boolean surfaces classified using the Youden optimal probability threshold classified
 353 23% of pixels as shrubland, an order of magnitude greater than the 2.5% identified from
 354 the true LiDAR-derived surface (Figure 6). More conservative predictions based on
 355 specificity thresholds predicted a lower proportion of shrubland (90% specificity: 11.9%;
 356 95% specificity: 6.1%, 99% specificity: 1.7%) with a higher precision, resulting in more
 357 accurate overall predictions and higher F1 scores.

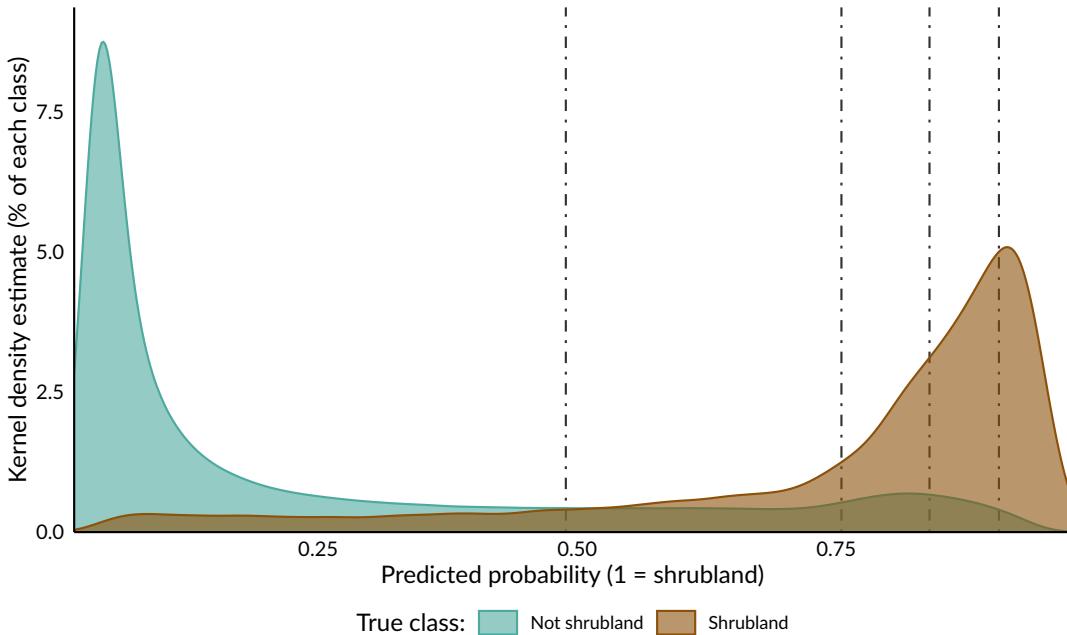


Figure 4.: Smoothed kernel density estimates of predicted probability of shrubland for both shrubland and non-shrubland pixels, calculated using a random sample of 1,000,000 pixels taken from the LiDAR patchwork prediction surface using the logistic ensemble model. Vertical lines indicate each of the four probability thresholds used to classify pixels. Colors represent the correct classification of the pixel.

358 3.2.2. 2019 Statewide Predictions

359 Model predictions reflected a similar geographical distribution of shrubland when ex-
 360 trapolating beyond the spatiotemporal boundaries of available LiDAR data to map
 361 shrubland across the entire state for 2019. Areas throughout the Adirondack Park and
 362 the Catskill Park, montane regions with mostly contiguous forest cover, showed notably
 363 less shrubland than more heavily populated areas (Figure 7). As expected, predictions
 364 in areas included in the LiDAR patchwork data set resembled the predictions for the
 365 LiDAR patchwork surface (Figure 5), though with some variation due in part to the
 366 temporal mismatch.

367 Shrubland probabilities were highest in areas classified by 2019 LCMAP as shrubland,
 368 as well as in areas classified as wetlands by either LCMAP or NLCD (Figure 8). Areas
 369 classified as tree cover were assigned extremely low probabilities.

370 Predictions for 2019 classified using the Youden optimal probability threshold classified
 371 22.3% of the state as shrubland (Figure 9), in line with the Youden optimal classified

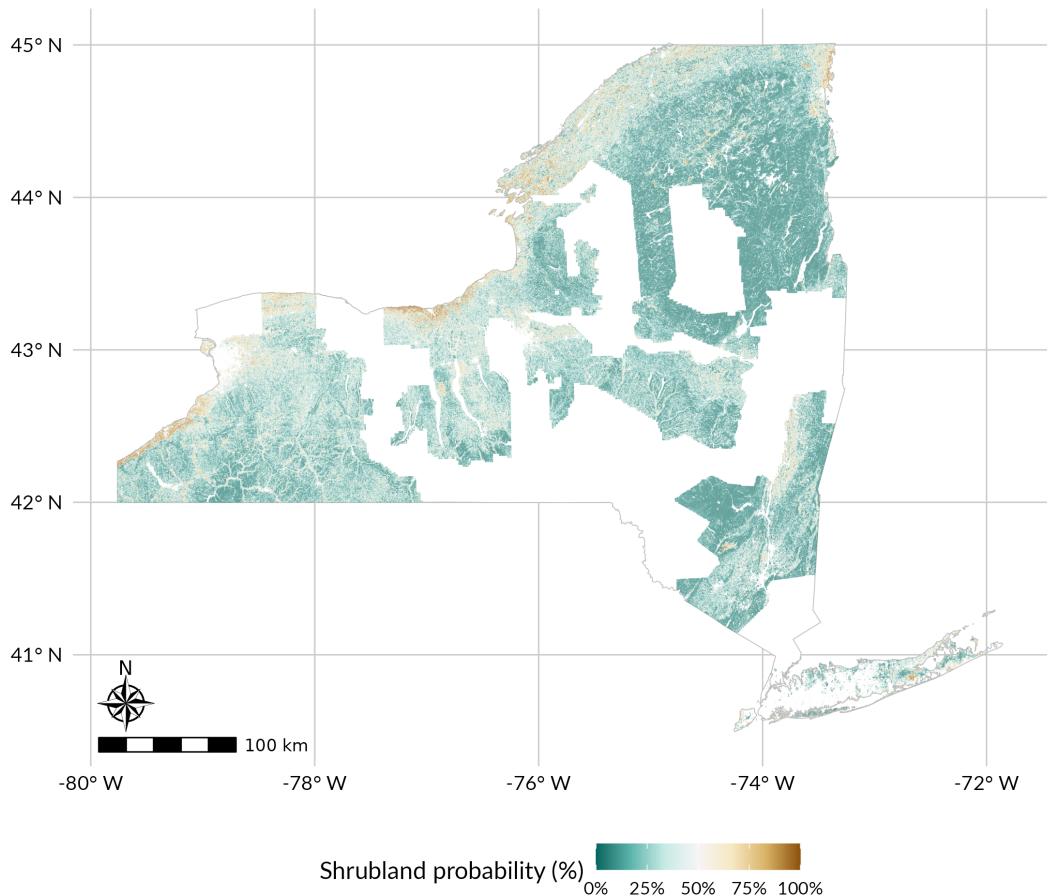


Figure 5.: Predicted probability of shrubland for the boundaries of all used LiDAR coverages, from the logistic ensemble model. Predictions were made using data reflecting the same year as LiDAR acquisition; the map therefore represents a temporal patchwork of predictions. Pixels in non-vegetated LCPRI land cover classes (developed, water, ice/snow, and barren) or above 1067 meters in elevation were not mapped and are shown in white.

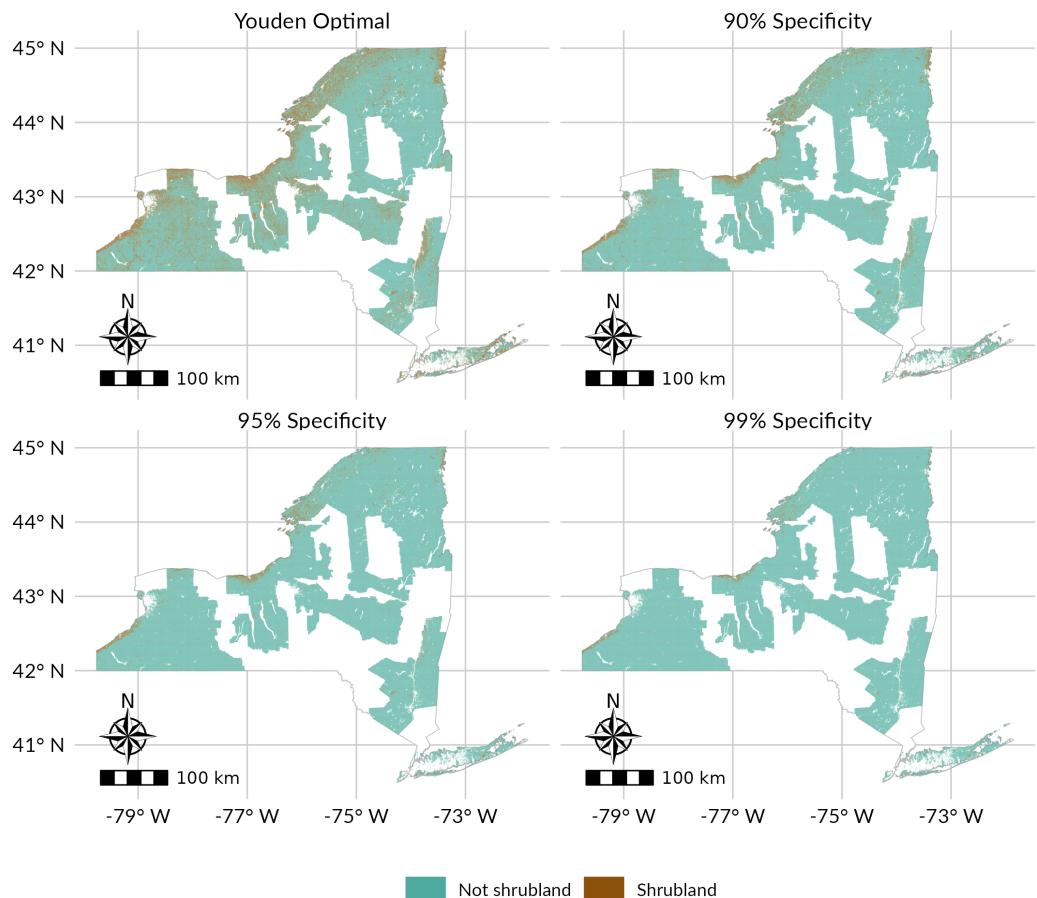


Figure 6.: Predicted shrubland locations within each LiDAR coverage, from the logistic ensemble model. Predicted pixel probabilities were classified using either the Youden-optimal threshold (which maximizes both sensitivity and specificity) or a threshold chosen to target a certain level of specificity, using thresholds derived from the validation data set. Predictions were made using data reflecting the same year as LiDAR acquisition; the map therefore represents a temporal patchwork of predictions. Pixels in non-vegetated LCPRI land cover classes (developed, water, ice/snow, and barren) or above 1067 meters in elevation were not mapped and are shown in white.

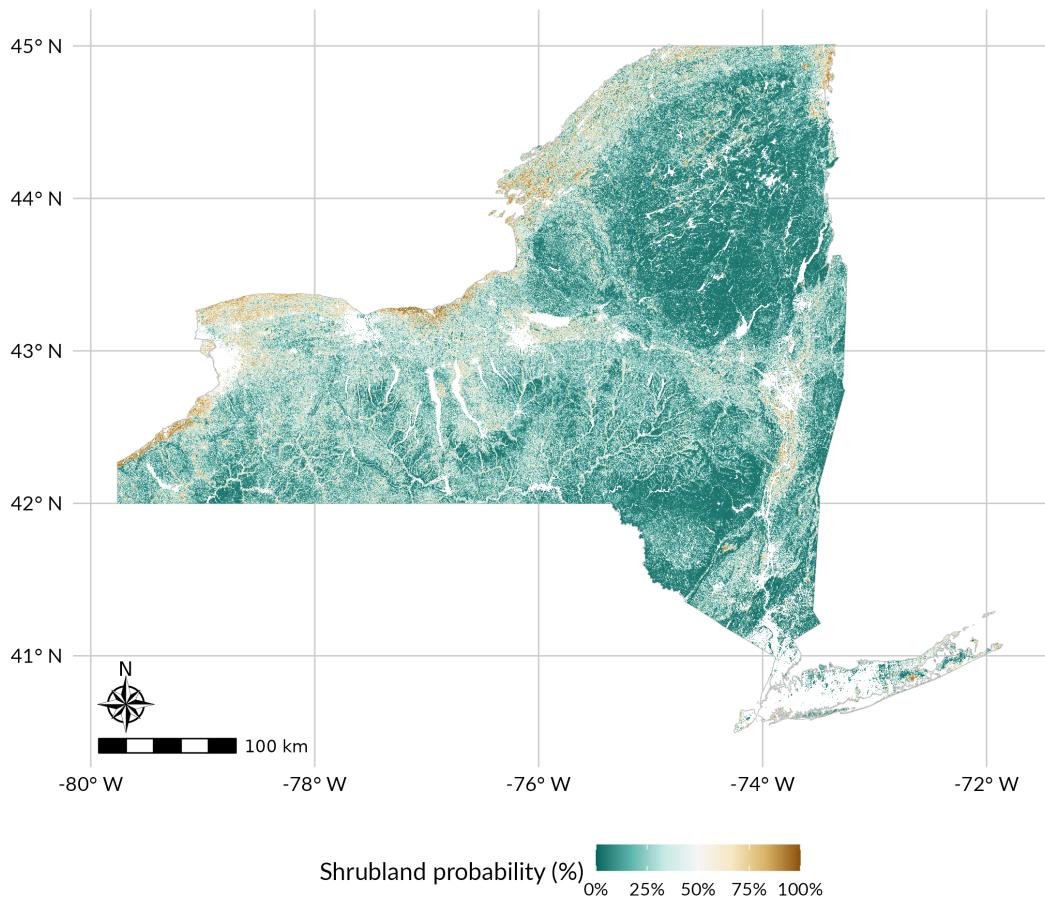


Figure 7.: Predicted probability of shrubland for 2019 across all mapped areas within New York State, from the logistic ensemble model. Pixels in non-vegetated LCPRI land cover classes (developed, water, ice/snow, and barren) or above 1067 meters in elevation were not mapped and are shown in white.

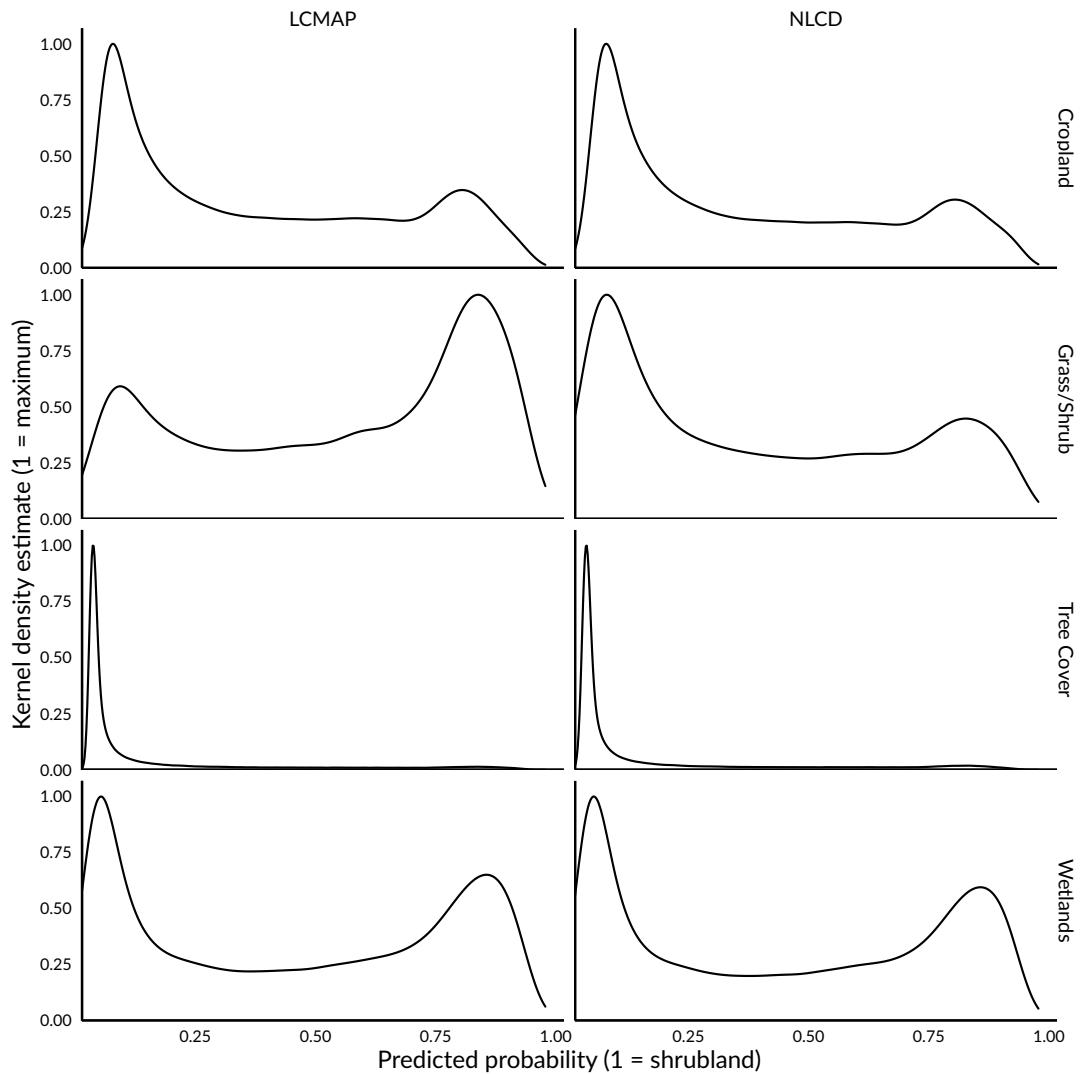


Figure 8.: Smoothed kernel density estimates of predicted probability of shrubland for the included LCMAP classes, calculated using a random sample of 1,000,000 pixels taken from the LiDAR patchwork prediction surface using the logistic ensemble model. Density estimates have been rescaled so that the most common probability for each panel is assigned a value of 1. NLCD land cover classes were remapped to LCMAP classes using LCMAP-defined translations.

372 LiDAR patchwork data set. The target-specificity thresholds classified more realistic
 373 proportions (90% specificity: 10.7%; 95% specificity: 5.1%, 99% specificity: 1.2%).

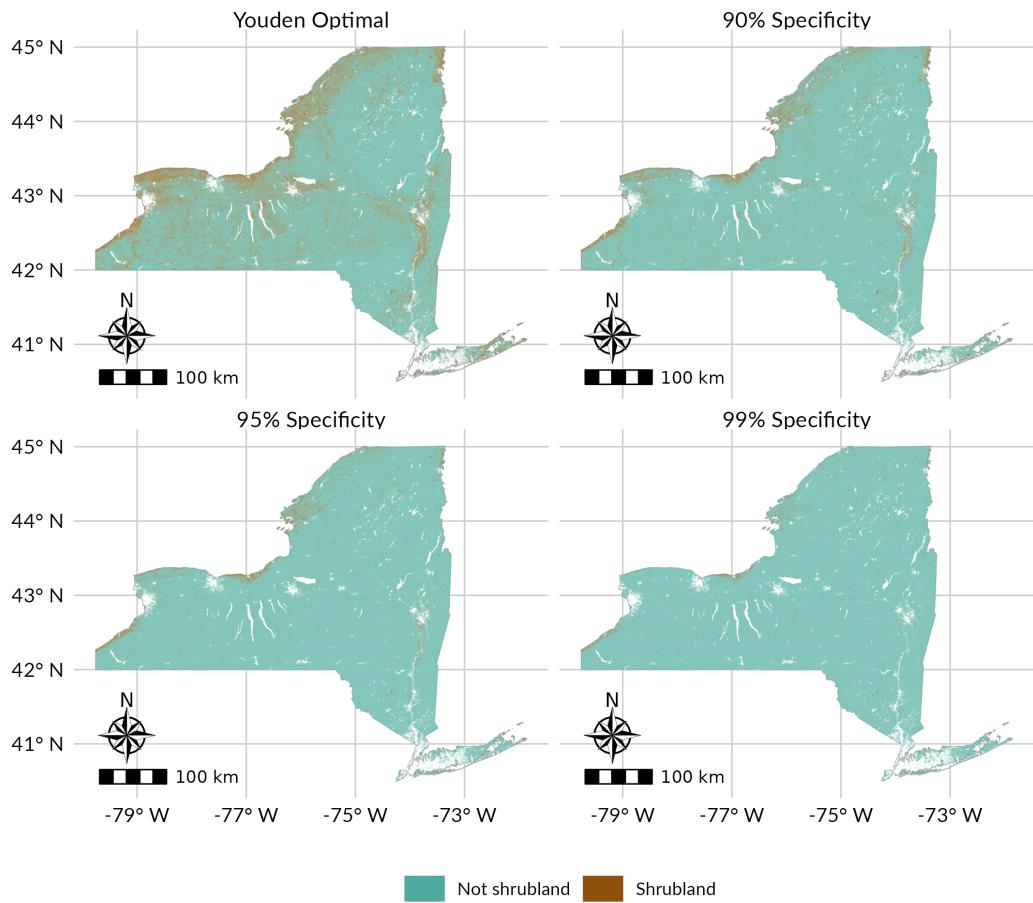


Figure 9.: Predicted shrubland locations for 2019 across the entire state, from the logistic ensemble model. Predicted pixel probabilities were classified using either the Youden-optimal threshold (which maximizes both sensitivity and specificity) or a threshold chosen to target a certain level of specificity, using thresholds derived from the validation data set. Pixels in non-vegetated LCPRI land cover classes (developed, water, ice/snow, and barren) or above 1067 meters in elevation were not mapped and are shown in white.

374 3.2.3. Validation Using 2019 Imagery

375 Pixels with a higher predicted probability of shrubland were generally more likely to
 376 represent shrubland, as determined using predictions and imagery from 2019 (Table 4).
 377 Pixels with a nominal predicted probability of shrubland at or above 0.95 were less
 378 likely to actually represent shrubland than pixels with probabilities between 0.8 and
 379 0.95, suggesting a failure to generalize to pixels with predictor values outside the ranges
 380 reflected in the training data (Efron 2020) (Table 4).

Table 4.: True classification of map pixels at various predicted shrubland probabilities. Numbers in parentheses reflect the percentage of pixels in each probability bracket in a given classification. Samples taken for two additional bins, ranging from 0% to 1% and 99% to 100%, were merged into the 0% to 5% and 95% to 100% bins due to the rarity of these extreme probabilities, resulting in these bins having slightly more observations.

Predicted probability of shrubland	# in sample	True classification (from imagery)		
		Not shrubland	Unknown	Shrubland
(0,0.05]	163	158 (96.9%)	5 (3.1%)	0 (0.0%)
(0.05,0.1]	159	144 (90.6%)	11 (6.9%)	4 (2.5%)
(0.1,0.2]	159	137 (86.2%)	18 (11.3%)	4 (2.5%)
(0.2,0.3]	159	130 (81.8%)	17 (10.7%)	12 (7.5%)
(0.3,0.4]	159	121 (76%)	19 (12%)	19 (12%)
(0.4,0.5]	159	118 (74.2%)	24 (15.1%)	17 (10.7%)
(0.5,0.6]	158	114 (72.2%)	26 (16.5%)	18 (11.4%)
(0.6,0.7]	159	104 (65.4%)	25 (15.7%)	30 (18.9%)
(0.7,0.8]	159	100 (62.9%)	24 (15.1%)	35 (22.0%)
(0.8,0.9]	159	70 (44.0%)	25 (15.7%)	64 (40.3%)
(0.9,0.95]	159	48 (30%)	30 (19%)	81 (51%)
(0.95,1]	240	154 (64%)	16 (7%)	70 (29%)

381 The majority of shrubland pixels, as identified using imagery, were classified as either
 382 tree cover or wetland pixels in both NLCD and LCMAP (Figure 10). The majority of
 383 non-shrubland pixels with a high predicted probability of shrubland were classified as
 384 cropland by LCMAP (Figure 10).

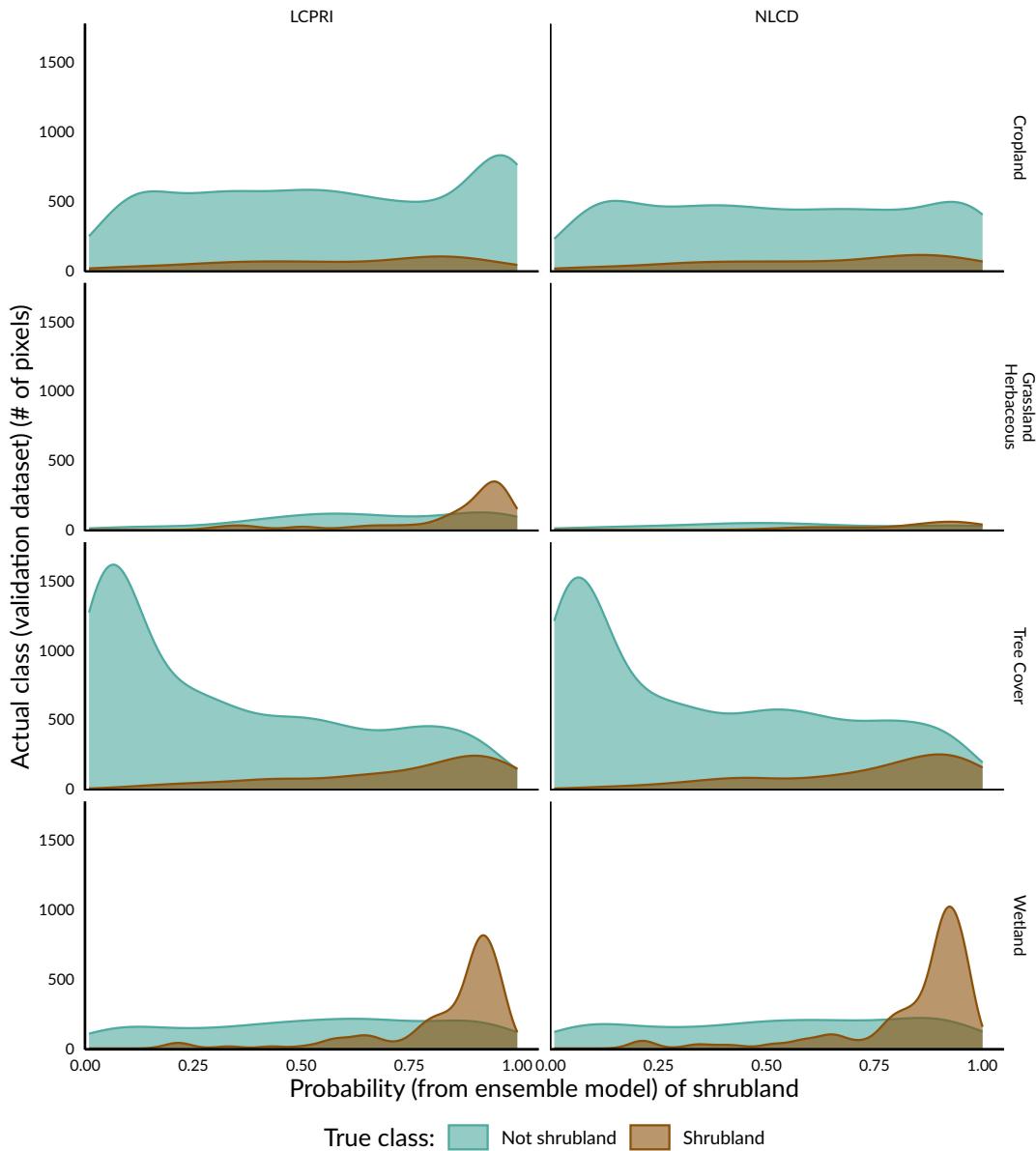


Figure 10.: Number of pixels (kernel density estimate) representing shrubland (determined from NAIP imagery) by LCMAP and NLCD classification. Pixels which could not be definitively identified as shrubland are classed as “not shrubland”.

385 **4. Discussion**

386 In this study, we produced a model capable of predicting shrubland locations, as delin-
387 eated using airborne LiDAR data and LCMAP land cover classifications, across a large,
388 fragmented and heterogenous landscape (New York State). Overall, we found that our
389 models were effective at distinguishing between shrubland and other land cover classes,
390 and produced qualitatively sensible map outputs, even when extrapolating beyond the
391 original training data. Our results serve to demonstrate that incorporating airborne Li-
392 DAR data can improve land cover classifications, particularly for marginal, transitional
393 and emergent cover types that may not be well represented in land cover class defini-
394 tions. More practically, we provide new maps of marginal cover types, such as invasive
395 shrub/scrub and degraded young forests, that have emerged mostly on post-agricultural
396 and post-industrial lands during the last century. While our model has limitations in
397 the labeling of training data and the inherent difficulty in predicting rare events, it was
398 effective at distinguishing between shrubland and other cover types. This modeling ap-
399 proach addresses a persistent measurement gap and enables monitoring, research, and
400 stewardship of these emerging and novel communities at a landscape scale.

401 **4.1. *Model predictions reflect known patterns***

402 Our model predictions, both for areas included in the LiDAR patchwork and for the
403 2019 surface, reflect known patterns in land cover throughout New York State. Although
404 areas classified as developed by LCMAP were excluded from predictions, areas of higher-
405 intensity human land use – such as the Hudson Valley in the eastern region of the state,
406 the I-90 highway corridor running East-West (along roughly 43N latitude), and the
407 northern border of the state, particularly around the Great Lakes – were consistently
408 classified as having a higher probability of shrubland (Figure 5, Figure 7). These areas
409 have likely been more recently impacted by human activity, with cropland only more re-
410 cently being left to natural regeneration. By the same pattern, areas of lower population
411 density and less intensive land use history such as the Adirondack and Catskill Parks
412 have consistently low probabilities of shrubland, reflecting the older, mid-successional
413 forests that characterize these areas. When extrapolating beyond the spatiotemporal

414 boundaries of available LiDAR data, our 2019 statewide model predicted a similar abundance
415 of shrubland with a similar regional distribution as the LiDAR patchwork data set, with similar areas of shrubland along the Great Lakes and adjacent to human population centers. This pattern reflects the continuing decline in agricultural land across
416 New York State (USDA National Agricultural Statistics Service 2019); while much agricultural land is being converted to developed land classes, a large proportion of former
417 cultivated lands have also been allowed to regenerate.
418

421 ***Predicting rare events is challenging***

422 As previously noted, shrubland is rare in New York due to the state's history of deforestation and regeneration combined with the relatively mild disturbance regime of the
423 northern forest (Lorimer 2001). One challenge in predicting rare events is that even a
424 low percentage of false positives can quickly drown out true positives; a model with
425 90% specificity predicting a data set with only 1% positive cases will produce 9 false
426 positives for every true positive it generates.
427

428 For this reason, we classified our predictions using a range of thresholds targeting increasingly high specificities (Table 3). These more stringent thresholds successfully increased
429 model precision (sometimes referred to as positive predictive power), making it more
430 likely that a positive prediction represents a true positive, though at the cost of lower
431 sensitivity. Of these three thresholds, the 95% specificity target best balanced sensitivity
432 and precision, based on its F1 score. Using this threshold, more than half of true
433 shrubland pixels were correctly classified and 22% of positive predictions reflect true
434 shrubland (as defined in the LiDAR analysis), approximately 10 times better than simply
435 guessing using the shrubland occurrence rate of 2.5%. Although there is room for
436 improvement, this level of accuracy is sufficient to identify potential areas for further
437 research and stewardship.
438

439 Additionally, many of the areas identified as shrubland by our model and confirmed via
440 NAIP imagery were classified as tree cover or wetlands by LCMAP and NLCD (Figure
441 10). This suggests that incorporating LiDAR data in land cover mapping workflows
442 may help to distinguish shrublands from optically similar classes, even if LiDAR is only

443 used to improve the quality of training set labels. It also potentially speaks to the ben-
444 efits of more targeted, regional land cover products to supplement the well-established
445 national models; regional efforts that can take advantage of regional data sets to improve
446 accuracy on cover types of regional importance.

447 **5. Conclusion**

448 This study aimed to predict the locations of shrubland across New York State, in or-
449 der to improve both understanding and stewardship of these plant communities. Using
450 a stacked ensemble model combining multiple machine learning models fit to data la-
451 beled using a combination of airborne LiDAR data and national land cover products,
452 we generated predictions of shrubland occurrences for both all spatiotemporal extents
453 with matching LiDAR data and for the entirety of New York State for 2019. Our model
454 was highly effective at distinguishing between shrubland and other cover types on both
455 the test set (AUC 0.893) and the LiDAR temporal patchwork (AUC 0.904), and bal-
456 anced sensitivity and precision effectively given the rarity of shrubland across the state.
457 These results suggest that combining remote sensing data from multiple sources may
458 improve land cover models, that regionally focused land cover models may complement
459 national products, and that shrubland may be effectively identified and monitored using
460 spaceborne remote sensing data.

461 **6. Acknowledgments**

462 Funding was provided by the Environmental Protection Fund via the NYS Department
463 of Environmental Conservation.

464 **7. Disclosure Statement**

465 The authors report there are no competing interests to declare.

⁴⁶⁶ **8. Data Availability Statement**

⁴⁶⁷ Data is available online from <https://doi.org/10.5281/zenodo.6519232>.

468 **References**

- 469 Alexander, Jake M, Jeffrey M Diaz, and Jonathan M Levine. 2015. "Novel Competitors
470 Shape Species' Responses to Climate Change." *Nature* 525: 515–18. <https://doi.org/10.1038/nature14952>.
- 472 Anderson, James R, Ernest E Hardy, John T Roach, and Richard E Witmer. 1976. *A
473 Land Use and Land Cover Classification System for Use with Remote Sensor Data.*
474 Vol. 964. US Government Printing Office.
- 475 Askins, Robert A. 2001. "Sustaining Biological Diversity in Early Successional Commu-
476 nities: The Challenge of Managing Unpopular Habitats." *Wildlife Society Bulletin*
477 29 (2): 407–12.
- 478 ASPRS. 2014. "ASPRS Positional Accuracy Standards for Digital Geospatial Data."
479 *Photogrammetric Engineering and Remote Sensing* 81 (3): 53. <https://doi.org/10.14358/PERS.81.3.A1-A26>.
- 481 Austin, Peter C, and Ewout W Steyerberg. 2012. "Interpreting the Concordance Statis-
482 tic of a Logistic Regression Model: Relation to the Variance and Odds Ratio of
483 a Continuous Explanatory Variable." *BMC Medical Research Methodology* 12 (82).
484 <https://doi.org/10.1186/1471-2288-12-82>.
- 485 Benjamin, Karyne, Gerald Domon, and Andre Bouchard. 2005. "Vegetation Composi-
486 tion and Succession of Abandoned Farmland: Effects of Ecological, Historical and
487 Spatial Factors." *Landscape Ecology* 20 (6): 627–47. <https://doi.org/10.1007/s10980-005-0068-2>.
- 489 Beven, Keith J., and Mike J. Kirkby. 1979. "A Physically Based, Variable Contributing
490 Area Model of Basin Hydrology." *Hydrological Sciences Bulletin* 24 (1): 43–69. <https://doi.org/10.1080/02626667909491834>.
- 492 Bogner, Christina, Bumsuk Seo, Dorian Rohner, and Björn Reineking. 2018. "Classifi-
493 cation of Rare Land Cover Types: Distinguishing Annual and Perennial Crops in an
494 Agricultural Catchment in South Korea." Edited by Krishna Prasad Vadrevu. *PLOS
495 ONE* 13 (1): e0190476. <https://doi.org/10.1371/journal.pone.0190476>.
- 496 Breiman, Leo. 2001. "Random Forests." *Machine Learning* 45: 5–32. <https://doi.org/10.1023/A:1010933404324>.
- 498 Brown, Jesslyn F., Heather J. Tollerud, Christopher P. Barber, Qiang Zhou, John L.

- 499 Dwyer, James E. Vogelmann, Thomas R. Loveland, et al. 2020. “Lessons Learned
500 Implementing an Operational Continuous United States National Land Change
501 Monitoring Capability: The Land Change Monitoring, Assessment, and Projec-
502 tion (LCMAP) Approach.” *Remote Sensing of Environment* 238: 111356. <https://doi.org/10.1016/j.rse.2019.111356>.
- 503
- 504 Chollet, François. 2015. “Keras.” <https://keras.io>.
- 505 Cocke, Allison E., Peter Z. Fulé, and Joseph E. Crouse. 2005. “Comparison of
506 Burn Severity Assessments Using Differenced Normalized Burn Ratio and Ground
507 Data.” *International Journal of Wildland Fire* 14: 189–98. <https://doi.org/10.1071/WF04010>.
- 508
- 509 Cramer, Viki A, Richard J Hobbs, and Rachel J Standish. 2008. “What’s New about
510 Old Fields? Land Abandonment and Ecosystem Assembly.” *Trends in Ecology &*
511 *Evolution* 23 (2): 104–12. <https://doi.org/10.1016/j.tree.2007.10.005>.
- 512 Dey, Daniel C, Benjamin O Knapp, Mike A Battaglia, Robert L Deal, Justin L Hart,
513 Kevin L O’Hara, Callie J Schweitzer, and Thomas M Schuler. 2019. “Barriers to
514 Natural Regeneration in Temperate Forests Across the USA.” *New Forests* 50 (1):
515 11–40. <https://doi.org/10.1007/s11056-018-09694-6>.
- 516 Dormann, Carsten F., Justin M. Calabrese, Gurutzeta Guillera-Arroita, Eleni Mate-
517 chou, Volker Bahn, Kamil Bartoń, Colin M. Beale, et al. 2018. “Model Averag-
518 ing in Ecology: A Review of Bayesian, Information-Theoretic, and Tactical Ap-
519 proaches for Predictive Inference.” *Ecological Monographs* 88 (4): 485–504. <https://doi.org/10.1002/ecm.1309>.
- 520
- 521 Dwyer, John L., David P. Roy, Brian Sauer, Calli B. Jenkerson, Hankui K. Zhang, and
522 Leo Lymburner. 2018. “Analysis Ready Data: Enabling Analysis of the Landsat
523 Archive.” *Remote Sensing* 10 (9). <https://doi.org/10.3390/rs10091363>.
- 524 Dyer, James M. 2006. “Revisiting the Deciduous Forests of Eastern North Amer-
525 ica.” *BioScience* 56 (4): 341–52. [https://doi.org/10.1641/0006-3568\(2006\)56%5B341:RTDFOE%5D2.0.CO;2](https://doi.org/10.1641/0006-3568(2006)56%5B341:RTDFOE%5D2.0.CO;2).
- 526
- 527 Efron, Bradley. 2020. “Prediction, Estimation, and Attribution.” *Journal of the Amer-
528 ican Statistical Association* 115 (530): 636–55. <https://doi.org/10.1080/01621459.2020.1762613>.
- 529

- 530 Falkowski, Michael J., Jeffrey S. Evans, Sebastian Martinuzzi, Paul E. Gessler, and
531 Andrew T. Hudak. 2009. "Characterizing Forest Succession with Lidar Data: An
532 Evaluation for the Inland Northwest, USA." *Remote Sensing of Environment* 113
533 (5): 946–56. [https://doi.org/https://doi.org/10.1016/j.rse.2009.01.003](https://doi.org/10.1016/j.rse.2009.01.003).
- 534 Fargione, Joseph E, Steven Bassett, Timothy Boucher, Scott D Bridgman, Richard T
535 Conant, Susan C Cook-Patton, Peter W Ellis, et al. 2018. "Natural Climate Solutions
536 for the United States." *Science Advances* 4 (11): eaat1869. <https://doi.org/10.1126/sciadv.aat1869>.
- 538 Flinn, Kathryn M, and Mark Vellend. 2005. "Recovery of Forest Plant Communities in
539 Post-Agricultural Landscapes." *Frontiers in Ecology and the Environment* 3 (5): 243–
540 50. [https://doi.org/10.1890/1540-9295\(2005\)003%5B0243:ROFPCI%5D2.0.CO;2](https://doi.org/10.1890/1540-9295(2005)003%5B0243:ROFPCI%5D2.0.CO;2).
- 541 Flinn, Kathryn M, Mark Vellend, and PL Marks. 2005. "Environmental Causes and
542 Consequences of Forest Clearance and Agricultural Abandonment in Central New
543 York, USA." *Journal of Biogeography* 32 (3): 439–52. <https://doi.org/10.1111/j.1365-2699.2004.01198.x>.
- 545 Foster, David R, Glenn Motzkin, and Benjamin Slater. 1998. "Land-Use History as Long-
546 Term Broad-Scale Disturbance: Regional Forest Dynamics in Central New England."
547 *Ecosystems* 1 (1): 96–119. <https://doi.org/10.1007/s100219900008>.
- 548 Fridley, Jason. 2012. "Extended Leaf Phenology and the Autumn Niche in Deciduous
549 Forest Invasions." *Nature* 485: 359–62. <https://doi.org/10.1038/nature11056>.
- 550 Friedman, Jerome H. 2002. "Stochastic Gradient Boosting." *Computational Statistics
551 and Data Analysis* 38 (4): 367–78. [https://doi.org/10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2).
- 552 GDAL/OGR contributors. 2021. *GDAL/OGR Geospatial Data Abstraction Software
553 Library*. Open Source Geospatial Foundation. <https://gdal.org>.
- 554 Good, Irving J. 1952. "Rational Decisions." *Journal of the Royal Statistical Society.
555 Series B (Methodological)* 14 (1): 107–14. <http://www.jstor.org/stable/2984087>.
- 556 Gorelick, Noel, Matt Hancher, Mike Dixon, Simon Ilyushchenko, David Thau, and Re-
557 becca Moore. 2017. "Google Earth Engine: Planetary-Scale Geospatial Analysis for
558 Everyone." *Remote Sensing of Environment* 202: 18–27.
- 559 Haibo He, and E. A. Garcia. 2009. "Learning from Imbalanced Data." *IEEE Transac-
560 tions on Knowledge and Data Engineering* 21 (9): 1263–84. <https://doi.org/10.1109/>

- 561 tkde.2008.239.
- 562 Hijmans, Robert J. 2021. *Terra: Spatial Data Analysis*. [https://CRAN.R-project.org/
package=terra](https://CRAN.R-project.org/package=terra).
- 564 Hobbs, Richard J, Eric Higgs, and James A Harris. 2009. “Novel Ecosystems: Impli-
565 cations for Conservation and Restoration.” *Trends in Ecology & Evolution* 24 (11):
566 599–605. <https://doi.org/10.1016/j.tree.2009.05.012>.
- 567 Huang, Jianfeng, Xinchang Zhang, Qinchuan Xin, Ying Sun, and Pengcheng Zhang.
568 2019. “Automatic Building Extraction from High-Resolution Aerial Images and Li-
569 DAR Data Using Gated Residual Refinement Network.” *ISPRS Journal of Photo-*
570 *grammetry and Remote Sensing* 151 (May): 91–105. <https://doi.org/10.1016/j.isprsjprs.2019.02.019>.
- 572 Johnson, Vanessa S, John A Litvaitis, Thomas D Lee, and Serita D Frey. 2006. “The
573 Role of Spatial and Temporal Scale in Colonization and Spread of Invasive Shrubs
574 in Early Successional Habitats.” *Forest Ecology and Management* 228 (1-3): 124–34.
575 <https://doi.org/10.1016/j.foreco.2006.02.033>.
- 576 Kauth, Richard J., and G. S. P. Thomas. 1976. “The Tasseled Cap - a Graphic De-
577 scription of the Spectral-Temporal Development of Agricultural Crops as Seen by
578 Landsat.” In *Symposium on Machine Processing of Remotely Sensed Data*.
- 579 Ke, Guolin, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei
580 Ye, and Tie-Yan Liu. 2017. “LightGBM: A Highly Efficient Gradient Boosting De-
581 cision Tree.” In *Advances in Neural Information Processing Systems*, edited by I.
582 Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
583 R. Garnett. Vol. 30. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf>.
- 585 Kennedy, Robert E, Zhiqiang Yang, and Warren B. Cohen. 2010. “Detecting Trends in
586 Forest Disturbance and Recovery Using Yearly Landsat Time Series: 1. LandTrendr
587 — Temporal Segmentation Algorithms.” *Remote Sensing of Environment* 114 (12):
588 2897–2910. <https://doi.org/10.1016/j.rse.2010.07.008>.
- 589 Kennedy, Robert E, Zhiqiang Yang, Noel Gorelick, Justin Braaten, Lucas Cavalcante,
590 Warren B. Cohen, and Sean Healey. 2018. “Implementation of the LandTrendr Al-
591 gorithm on Google Earth Engine.” *Remote Sensing* 10 (5). <https://doi.org/10.3390/>

- 592 [rs10050691](#).
- 593 King, David I., and Scott Schlossberg. 2014. “Synthesis of the Conservation Value of
594 the Early-Successional Stage in Forests of Eastern North America.” *Forest Ecology
595 and Management* 324: 186–95. <https://doi.org/10.1016/j.foreco.2013.12.001>.
- 596 Kulmatiski, Andrew, Karen H Beard, and John M Stark. 2006. “Soil History as a Pri-
597 mary Control on Plant Invasion in Abandoned Agricultural Fields.” *Journal of Ap-
598 plied Ecology* 43 (5): 868–76. <https://doi.org/10.1111/j.1365-2664.2006.01192.x>.
- 599 LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. “Deep Learning.” *Nature* 521:
600 436–44. <https://doi.org/10.1038/nature14539>.
- 601 Lorimer, Craig G. 2001. “Historical and Ecological Roles of Disturbance in Eastern
602 North American Forests: 9,000 Years of Change.” *Wildlife Society Bulletin (1973-
603 2006)* 29 (2): 425–39. <http://www.jstor.org/stable/3784167>.
- 604 Mahoney, Michael J., Colin M. Beier, and Aidan C. Ackerman. 2022. “terrainr: An
605 R Package for Creating Immersive Virtual Environments.” *Journal of Open Source
606 Software* 7 (69): 4060. <https://doi.org/10.21105/joss.04060>.
- 607 McCay, Timothy S., and Deanna H McCay. 2009. “Processes Regulating the Invasion
608 of European Buckthorn (*Rhamnus Cathartica*) in Three Habitats of the Northeast-
609 ern United States.” *Biological Invasions* 11 (8): 1835–44. <https://doi.org/10.1007/s10530-008-9362-7>.
- 610 NOAA National Centers for Environmental Information. 2022. “Climate at a Glance:
611 Statewide Mapping.” <https://www.ncdc.noaa.gov/cag/>.
- 612 Perring, Michael P, Rachel J Standish, and Richard J Hobbs. 2013. “Incorporating
613 Novelty and Novel Ecosystems into Restoration Planning and Practice in the 21st
614 Century.” *Ecological Processes* 2 (1): 1–8. <https://doi.org/10.1186/2192-1709-2-18>.
- 615 PRISM Climate Group. 2022. “PRISM Climate Data.” <https://prism.oregonstate.edu>.
- 616 Python Core Team. 2022. *Python: A dynamic, open source programming language*.
617 Python Software Foundation. <https://www.python.org/>.
- 618 R Core Team. 2021. *R: A Language and Environment for Statistical Computing*. Vienna,
619 Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- 620 Roussel, Jean-Romain, David Auty, Nicholas C. Coops, Piotr Tompalski, Tristan R. H.
621 Goodbody, Andrew Sánchez Meador, Jean-François Bourdon, Florian de Boissieu,

- 623 and Alexis Achim. 2020. “lidR: An r Package for Analysis of Airborne Laser Scanning
624 (ALS) Data.” *Remote Sensing of Environment* 251: 112061. <https://doi.org/10.1016/j.rse.2020.112061>.
- 625 Ruiz, Luis Ángel, Jorge Abel Recio, Pablo Crespo-Peremach, and Marta Sapena. 2018.
626 “An Object-Based Approach for Mapping Forest Structural Types Based on Low-
627 Density LiDAR and Multispectral Imagery.” *Geocarto International* 33 (5): 443–57.
628 <https://doi.org/10.1080/10106049.2016.1265595>.
- 629 Spiering, David J. 2019. “Brownfields and Old-Fields: Vegetation Succession in Post-
630 Industrial Ecosystems of Western New York.” PhD thesis, State University of New
631 York at Buffalo.
- 632 Stover, Marian E., and PL Marks. 1998. “Successional Vegetation on Abandoned Culti-
633 vated and Pastured Land in Tompkins County, New York.” *Journal of the Torrey
634 Botanical Society*, 150–64. <https://doi.org/10.2307/2997302>.
- 635 U.S. Geological Survey. 2019. “3D Elevation Program 1-Meter Resolution Digital Ele-
636 vation Model.” <https://www.usgs.gov/core-science-systems/ngp/3dep/data-tools>.
- 637 US Department of Agriculture. 2019. “National Agricultural Imagery Program.”
- 638 USDA National Agricultural Statistics Service. 2019. “2017 Census of Agriculture.”
639 www.nass.usda.gov/AgCensus.
- 640 Whitney, Gordon G. 1994. *From Coastal Wilderness to Fruited Plain: A History of
641 Environmental Change in Temperate North America from 1500 to the Present*. Cam-
642 bridge, United Kingdom: Cambridge University Press.
- 643 Wickham, James, Stephen V. Stehman, Daniel G. Sorenson, Leila Gass, and Jon A. De-
644 witz. 2021. “Thematic Accuracy Assessment of the NLCD 2016 Land Cover for
645 the Conterminous United States.” *Remote Sensing of Environment* 257: 112357.
646 <https://doi.org/https://doi.org/10.1016/j.rse.2021.112357>.
- 647 Williams, John W., and Stephen T. Jackson. 2007. “Novel Climates, No-Analog Com-
648 munities, and Ecological Surprises.” *Frontiers in Ecology and the Environment* 5 (9):
649 475–82. <https://doi.org/10.1890/070037>.
- 650 Wolpert, David H. 1992. “Stacked Generalization.” *Neural Networks* 5 (2): 241–59. [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1).
- 651 Wright, Marvin N., and Andreas Ziegler. 2017. “ranger: A Fast Implementation of Ran-
652 ger.” *Journal of Statistical Software* 80 (5): 1–27. <https://doi.org/10.18637/jss.v080.i05>.

- 654 dom Forests for High Dimensional Data in C++ and R.” *Journal of Statistical*
655 *Software* 77 (1): 1–17. <https://doi.org/10.18637/jss.v077.i01>.
- 656 Yang, Limin, Suming Jin, Patrick Danielson, Collin Homer, Leila Gass, Stacie M. Ben-
657 der, Adam Case, et al. 2018. “A New Generation of the United States National
658 Land Cover Database: Requirements, Research Priorities, Design, and Implemen-
659 tation Strategies.” *ISPRS Journal of Photogrammetry and Remote Sensing* 146: 108–23.
660 <https://doi.org/10.1016/j.isprsjprs.2018.09.006>.
- 661 Youden, W. J. 1950. “Index for Rating Diagnostic Tests.” *Cancer* 3 (1): 32–35.
662 [https://doi.org/10.1002/1097-0142\(1950\)3:1%3C32::AID-CNCR2820030106%3E3.0.CO;2-3](https://doi.org/10.1002/1097-0142(1950)3:1%3C32::AID-CNCR2820030106%3E3.0.CO;2-3).
- 663 Zarea, Asghar, and Ali Mohammadzadeh. 2016. “A Novel Building and Tree Detection
664 Method from LiDAR Data and Aerial Images.” *IEEE Journal of Selected Topics in*
665 *Applied Earth Observations and Remote Sensing* 9 (5): 1864–75. <https://doi.org/10.1109/jstars.2015.2470547>.