



1 Mapping Wetland Probability Across Massachusetts with Machine 2 Learning and Multiscale Predictors

3 Rafter S. Ferguson¹, Sebastian Gutwein¹, Eric Giordano¹, Keith Zaltzberg-Drezdahl¹

4
5 ¹Regenerative Design Group, Greenfield, 01301, United States

6 Correspondence to: Rafter S. Ferguson (rafterf@regenerativedesigngroup.com)

7
8 **Abstract.** Wetlands perform a vital array of ecosystem functions, but up to 50% of global wetlands have been lost and those
9 that remain are under ongoing threat from development pressures. Accurate and comprehensive maps are critical for the
10 management and protection of wetland resources. Conventional methods for wetland mapping are time consuming and
11 resource intensive, and the common mapping methods that rely on the inspection of aerial imagery often miss forested and
12 other wetland types that do not have a distinctive visual signature, i.e. cryptic wetlands. The use of machine learning and
13 spatial data to map wetlands is a growing field that promises a fast and efficient complement to conventional methods and
14 improved detection of forested and other cryptic wetlands. In this paper we demonstrate the use of a random forest model to
15 generate a large-scale, state-wide map of wetland probabilities in the Commonwealth of Massachusetts, using widely
16 available open source software and publicly accessible data. Through this model we also test the efficacy of multi-scale
17 predictors, including not only terrain derivatives used in previous research but also multi-scale implementations of soil,
18 vegetation, and spectral data. The random forest was trained on the official Massachusetts wetland inventory, and achieved
19 an overall accuracy rate of 92% relative to that dataset. The model showed particular promise in detecting cryptic wetlands
20 by identifying an additional 40% of probable wetland area statewide, and an additional 46% of forested wetland specifically.
21 The use of diverse multi-scale predictors was supported by model performance, variable importance measures, and the
22 feature selection process. This strategy for improving detection of cryptic wetlands and creating better estimates of wetland
23 extent, using non-proprietary software and data, will be a vital adjunct to conventional methods for wetland mapping and
24 monitoring.

25 1 Introduction

26 Wetlands perform an array of vital ecosystem functions, supporting biodiversity, regulating and stabilizing water flow,
27 removing pollutants from surface waters, and sequestering carbon (Mitsch and Gosselink, 2015). In many areas, much of
28 historic wetland area has been lost to development, and ongoing development continues to encroach further into wetlands.
29 Global wetland loss since 1900 may be as high as 71% (Davidson, 2014), and the US is estimated to have lost 39% of its
30 historic wetlands (Fluet-Chouinard et al., 2023). The protection of wetlands becomes increasingly important in the face of



31 climate change, as wetlands represent some of the most intense concentrations of sequestered carbon in the biosphere
32 (Nahlik and Fennessy, 2016) as well as providing critical buffers against extreme weather events. Keeping sequestered
33 carbon in wetlands and out of the atmosphere requires vigorous protection from persistent development pressures.

34
35 Accurate and comprehensive maps are critical for the management and protection of wetland resources. Development of
36 wetland inventories, however, is time-consuming and expensive, and wetland mapping efforts are constrained by limitations
37 on the capacity of state agencies and non-governmental organizations. Recent advances in remote sensing and machine
38 learning provide techniques and data sets for the rapid and efficient identification of wetlands and can complement and
39 amplify investments in time-intensive field mapping and even conventional desk methods for wetland identification.

40
41 The use of machine learning to map wetlands is a dynamic and fast-growing topic in research and planning at the state,
42 national, and international level (Felton et al., 2019; Gale, 2021; Halabisky et al., 2023; Rapinel et al., 2023). Machine
43 learning offers the opportunity to leverage existing datasets to generate or improve assessments of the location and extent of
44 wetlands across large areas, quickly and efficiently. Not intended to replace other methods of wetland identification,
45 machine learning can provide estimates of cryptic wetland area and help identify priority areas for closer inspection,
46 including field inspection and delineation.

47
48 In this paper we discuss an implementation of the machine learning algorithm random forest to produce a wetland
49 probability map of the Commonwealth of Massachusetts. In order to estimate the extent of cryptic wetlands that do not
50 appear on existing inventories, we developed a machine learning model making use of statewide spatial data and widely
51 available open-source software. The wetland probability map can be combined with existing inventories to create a more
52 accurate and inclusive estimate of total wetland area.

53 **1.1 Digital wetland mapping with random forest**

54 The growing field of digital wetland mapping has drawn on an array of machine learning algorithms. Any suitable algorithm
55 must be able to handle large and complex datasets while making few assumptions about the distribution of variables.
56 Random forest is an ensemble learning algorithm that meets these criteria, and is increasingly recognized as one of the most
57 powerful and effective machine learning strategies for wetland detection (Jafarzadeh et al., 2022). Random forest creates an
58 ensemble of decision trees, each trained on a random bootstrap sample of the data and a random subset of predictors
59 (Breiman, 2001). The model prediction is the averaged prediction of all individual trees, with the prediction appearing as a
60 probability value between 0 and 1. In addition to meeting the criteria outlined above, the bootstrapping and aggregating
61 strategy reduces bias, avoids overfitting, accommodates correlated predictors, has been shown to produce highly accurate
62 predictions without extensive tuning of hyperparameters, i.e. “out of the box,” with default settings, accounts for complex



63 multi-way interactions, and has a built-in variable importance mechanism that is useful for feature selection—all of which
64 have led to extensive use of random forest in remote sensing (Belgiu and Drăguț, 2016). Recent years have seen growing use
65 of the random forest algorithm in state and regional wetland detection projects in the US (Maxwell et al., 2016; O’Neil et al.,
66 2018) and incorporation into a toolbox in popular proprietary GIS software (The Wetland Identification Model (WIM) – A
67 New Arc Hydro Functionality for Predicting Wetland Locations Using LiDAR Elevation Data and Machine Learning, 2024).

68
69 Digital mapping of wetlands can draw on a wide range of potential predictors. The presence of wetlands in the landscape is
70 determined by the interaction of terrain and hydrology, and is in turn reflected in a range of soils, vegetation, hydrological
71 characteristics, and spectral signatures, all of which have been used as inputs in various wetland mapping projects
72 (Jafarzadeh et al., 2022). For all predictors, spatial scale is a key question. Wetlands are produced by the interaction of forces
73 operating and combining across multiple spatial scales, and there is no clear cut method for determining which scale is most
74 relevant across contexts. Choice of spatial scale for predictors is especially pertinent for approaches like random forest,
75 which uses point-based intersections of predictors and does not natively account for contextual information. Using multi-
76 scale predictors, i.e. calculating, resampling, or smoothing spatial predictors at multiple scales, can therefore help address
77 several challenges at once. In the case of terrain derivatives, it is clear that topography can control the accumulation and
78 retention of water at the scale of local undulations in grade (e.g. 1-10 meter) or broader topographic depressions (e.g. 10-100
79 meter, or larger). Other classes of predictors, such as vegetation height, soil type, or spectral data, can show similar multi-
80 scale neighborhood effects, where what matters at a given sample point is not the value of the predictor at that point, but the
81 average value in the area that forms its context. In the case of soil-based predictors, multi-scale smoothing can also help by
82 creating fuzzy boundaries between soil units, which has the potential to ameliorate some of the known challenges with the
83 crisp and sometimes arbitrary delineations that soil maps impose on underlying heterogeneity and varying transition zones
84 (Hunter et al., 2009; Nikiforova et al., 2020). As there is no way to know a priori which scales will be relevant for which
85 predictors, the variable importance scores generated by the random forest algorithm can be used to select the most relevant
86 predictors and scales.

87
88 In the related domain of digital soil mapping, the use of such multi-scale terrain derivatives is becoming increasingly well
89 established (Behrens et al., 2019), and has been shown to increase prediction accuracy dramatically (Miller et al., 2015).
90 Application to other classes of predictors aside from terrain, however, remains somewhat neglected. Recent work has
91 introduced the use of multi-scale terrain features specifically for digital wetland mapping (Halabisky et al., 2022).

92 **1.3 Study Area**

93 Massachusetts is a coastal state in the northeastern United States with a diverse range of soils, topography, and vegetation.
94 While Massachusetts was among the first states to implement wetland protections in the mid-20th century, it’s estimated



95 Massachusetts has lost a third of its wetlands since colonial times (Protecting Wetlands in Massachusetts | Mass.gov, 2024).
96 Roughly 30% of the state by area, in the western and north-central zones, is in the EPA Level III Ecoregion the Northeastern
97 Highlands, characterized by high elevations, rugged topography, mixed deciduous and coniferous forests, and soils built on
98 glacial till (Griffith et al., 2009). From the central through the eastern coastal zone, the Northeastern Coastal Zone makes up
99 roughly 60% of the commonwealth. The Northeastern Coastal Zone is marked by lower elevations, less topographical
00 variation, and a prevalence of sandy loams inland and sandy soils with low organic content in coastal areas. The remaining
01 9% of the state is in the Atlantic Coastal Pine Barrens, a low-lying area with diverse coastal landforms including dunes,
02 marshes, and bays, with vegetation dominated by scrubby oak-pine forests.

03 **1.4 Objectives**

04 Our principal objective was to demonstrate a method for generating a large-scale, state-wide map of wetland probabilities,
05 using widely available open source software and publicly available data. Such a method could serve as a fast and efficient
06 adjunct to conventional field and desk methods, and in turn support many goals for wetland monitoring, management, and
07 conservation. A secondary objective was to use machine learning to leverage the utility of an existing wetland inventory by
08 using it as training data that (together with the array of predictors) could identify additional probable wetland area missing
09 from the same inventory. Our tertiary objective was to assess the multi-scale predictor approach using a broader swath of
10 environmental variables than has previously been demonstrated in wetland detection efforts.

11 **2 Methods**

12 The data and modelling processes are represented in Fig. 1. All operations in this project were performed with free open
13 source software. Acquisition and preprocessing of geospatial data were performed using QGIS (QGIS Development Team,
14 2023), and data wrangling and modelling were conducted in the R statistical programming language (R Core Team, 2023).

15 **2.2 Data**

16 The Massachusetts Department of Environmental Protection (MassDEP) has produced an official state wetland inventory
17 that represents a distinctively high-quality resource of state-level wetland information for planning and research in
18 Massachusetts (Baker et al., 2019). The first version of the MassDEP Wetlands map was completed in 2006, based on visual
19 assessment of 1:12,000 stereo imagery taken between 1990 and 2000. This original map featured field validation for up to
20 10% of identified wetlands. Starting in 2007, MassDEP created an updated layer, using 2005 imagery at 1:5,000 scale.
21 Production of the 2005 MassDEP Wetlands layers (MDW) did not feature any field validation. The layer compares favorably
22 with the National Wetland Inventory (NWI), which is produced with similar methods but using older and larger-scale
23 imagery. The present utility of the MDW is constrained by both the intrinsic limitations of the method that it shares with the
24 NWI, and with the age of the imagery used to produce it. Numerous studies have confirmed that the NWI misses many

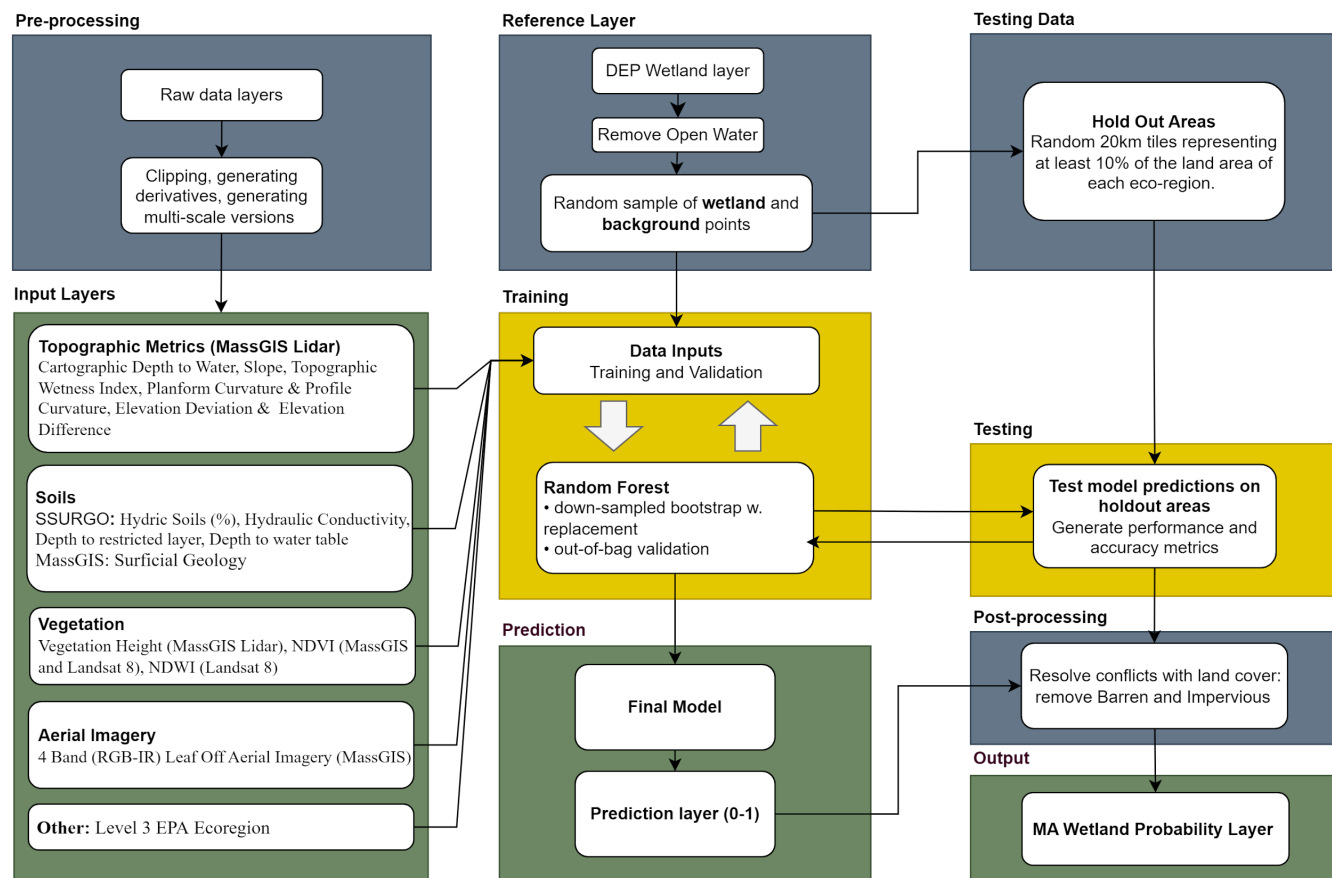


25 wetlands, particularly but not only forested and smaller wetlands (each of which may in aggregate make up a significant
26 portion of total wetlands) (Kudray and Gale, 2000; Martin et al., 2012; Matthews et al., 2016; Stolt and Baker, 1995; Tiner,
27 1990). While no comparison study has been conducted in Massachusetts, reports from field delineators, together with
28 comparison with acquired field delineations, confirm that similar limitations apply to the MDW. Additionally, trends in
29 wetland loss and gain are ongoing and dynamic, and the imagery on which the MDW is based is now nearly two decades old
30 (Baker et al., 2019).

31

32 Training and testing data were developed from the MDW. The layer was edited to remove open water according to the DEP
33 Wetland designation. Given the scale of the study area, researchers opted for a moderately dense sampling approach,
34 creating a training dataset of 30,000 points. Reference wetland data were assembled from 10,000 random points drawn from
35 within wetland polygons, and background data were assembled from 20,000 random points not within wetland polygons.
36 Non-wetland points are treated as ‘background’ or ‘pseudo-absence’ rather than true absence points, as it is expected that
37 some of the 20,000 background points were placed within cryptic wetlands that do not appear in the MDW. This distinction
38 was accommodated in our downsampling approach described in section **2.3 Model Fitting** below.

39



40

41

42 **Figure 1. Predictors and process flow for generating a statewide wetland probability layer.** Using machine learning to

43 generate a Wetland Probability Score at 4m resolution will support a more inclusive estimation of wetland extent across the

44 state.

45

46

47 We extend the logic of the multi-scale approach in digital soil mapping, and the application of multi-scale terrain derivatives

48 to wetland mapping, by aggregating all predictors to multiple larger resolutions (whenever feasible based on the starting

49 resolution of available layers), and including all resolutions in our initial model. Figure 2 illustrates three of the most

50 important predictors rendered at multiple scales for the same area. Approaches to resampling included specification of

51 neighborhood size for terrain derivatives and mean smoothing or filtering for other variables.

52

53 Terrain variables were developed from the MassGIS LiDAR point cloud data. The point cloud was used to generate a digital

54 elevation model (DEM) for the entire state, resampled to 4m resolution. The DEM was in turn used to generate derivative



55 measures at multiple scales, including Slope, Cartographic Depth to Water (DTW), Planoform Curvature, Profile Curvature,
56 and Topographic Wetness Index. Cartographic Depth to Water is a metric that estimates the difference in elevation between
57 the soil surface and the closest open water (Murphy et al., 2011). Topographic Wetness Index (TWI) is a predecessor of
58 DTW that is used, like DTW, to model soil moisture conditions (Murphy et al., 2009). Profile and planform curvature
59 quantify the curvature of the earth's surface in the direction of slope, and perpendicular to slope, respectively (Maxwell et
60 al., 2016).

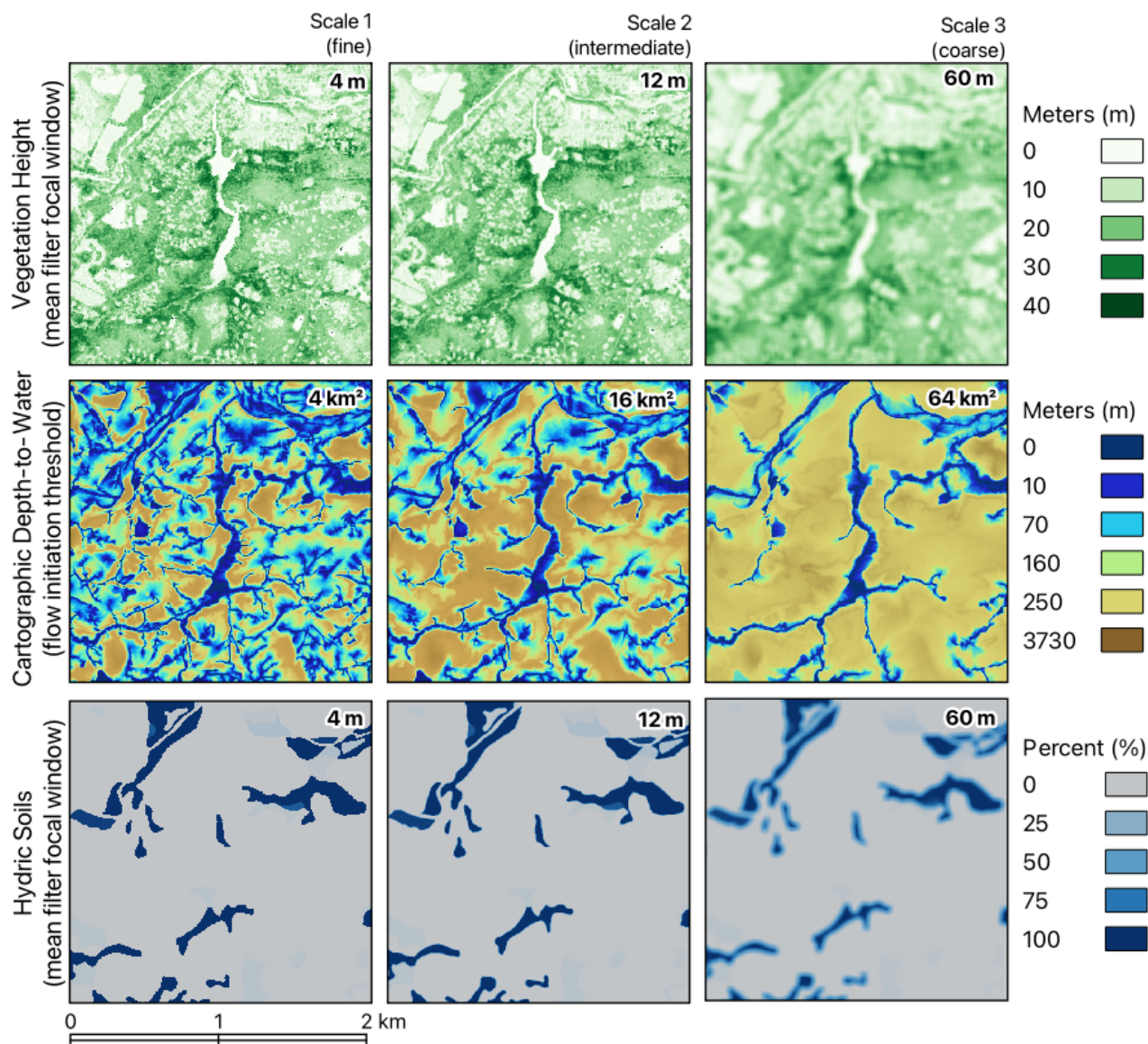
61

62 With one exception, soil variables were drawn from the Soil Survey Geographic Database (SSURGO) maintained by the US
63 Natural Resources Conservation Service (Soil Survey Staff, n.d.). These variables included Depth to a Restrictive Soil Layer,
64 Depth to Water Table, Percent Hydric Soils, and Hydraulic Conductivity (Soil Survey Staff, 2015). Depth to Restrictive Soil
65 Layer is the estimated distance from the soil surface to a layer that impedes the movement of water (often bedrock). Percent
66 Hydric Soils refers to the percent of the SSURGO map unit that is composed of hydric components. Hydraulic Conductivity
67 is the volume of water that would move through the soil under saturation conditions in a standardized area and unit of time.
68 The vector-based variables from SSURGO were modified to multiple scales by generating a raster layer for each of the
69 relevant fields of the SSURGO vector layer, at 4m resolution, and then applying a mean filter at 12m and 60m. An additional
70 variable, Surficial Geology, was acquired from the MassGIS Surficial Geology layer, which was in turn acquired from the
71 US Geological Survey (Stone et al., 2018).

72

73 Vegetation variables were acquired from multiple sources. The mean and median of NDVI Composite and Normalized
74 Difference Water Index (NDWI) Composite were each acquired, based on 30m resolution LANDSAT 8 imagery, directly
75 via Google Earth Engine. Normalized Difference Vegetation Index (NDVI) was calculated from 2021 MassGIS Leaf-Off
76 Aerial Imagery, and aggregated to multiple scales. NDVI is a widely used metric of vegetation health that uses the red and
77 near-infrared bands to monitor the level of photosynthetic activity. NDWI similarly uses the green and near-infrared bands to
78 assess moisture levels. Vegetation height was calculated based on first returns from the MassGIS LiDAR point cloud data
79 sets, and aggregated to multiple scales. The EPA US Level 3 Ecoregions for Massachusetts were also included as a three-
80 value categorical predictor: Atlantic Coastal Pine Barrens, Northeastern Coastal Zone, and Northeastern Highlands. In
81 addition, bands one through four (red, green, blue, and near infrared, respectively) from 2021 MassGIS Leaf-Off Aerial
82 Imagery were acquired and smoothed to multiple scales.

83



84

85 **Figure 2. Multi-scale predictors.** Three of the predictors included in the final model are shown for the same area, an
 86 approximately 2 km square centrally located in the Northeastern Coastal Zone. Predictors are shown across the range of
 87 scales that were each included in the model. The method of multi-scaling for each predictor is indicated in parentheses below
 88 the predictor name on the left side of the figure. Model predictions for the area shown here can be seen in callout (C) in Fig.

89 4.

90



91 **2.3 Training and testing**

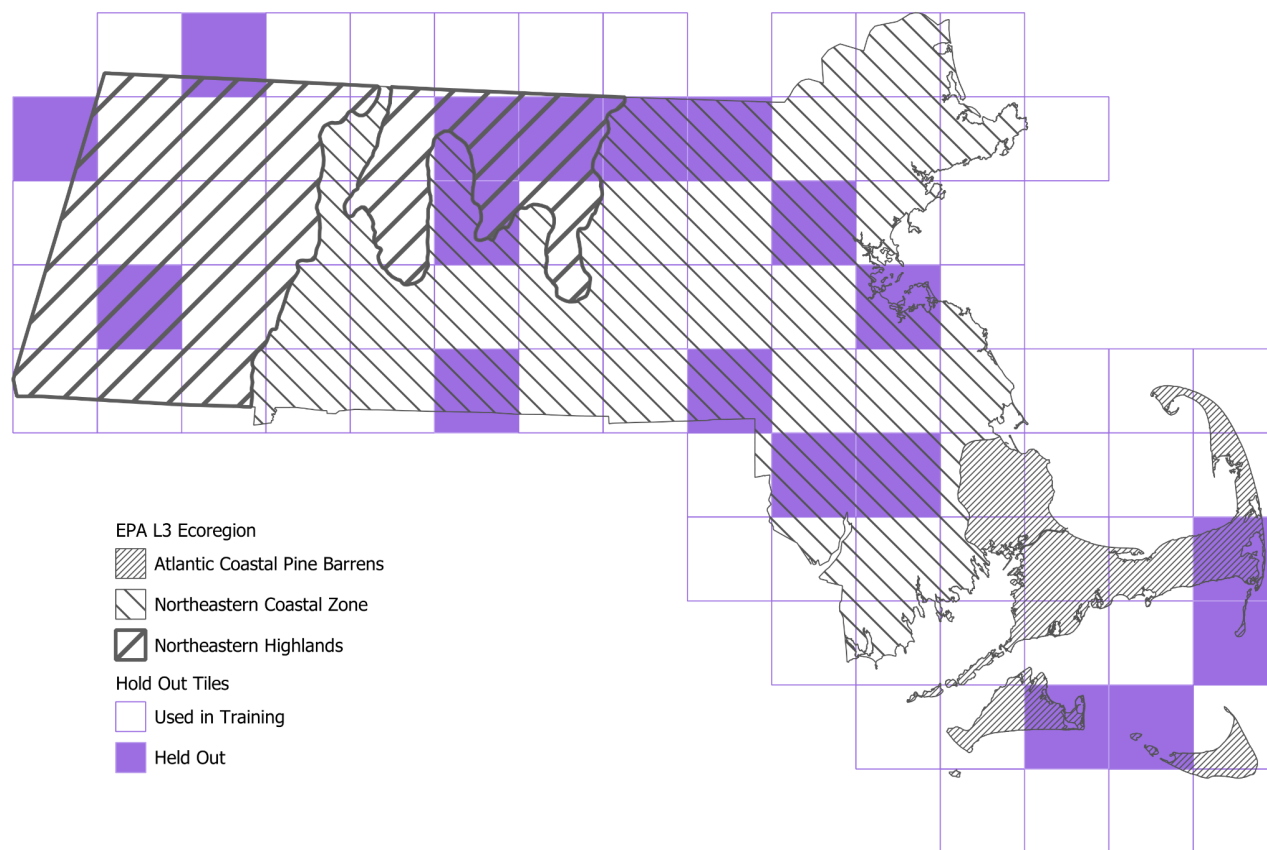
92 Model testing was conducted with two different spatial configurations of holdout data. In the interest of a more rigorous
 93 assessment of model performance and generalizability, holdout data for testing the final model was based on the random
 94 selection of 20 km by 20 km tiles within each EPA Level 3 ecoregion (Fig. 3). Within each of the three ecoregions in
 95 Massachusetts, tiles were randomly selected for holdout until at least 10% of land area within that ecoregion was designated
 96 as holdout (Table 1). Points in these tiles were held back from the training data, comprising 23% of the 30,000 points. By
 97 using tile-based holdout, predictive performance could be assessed based on nearby but unsampled landscapes, rather than
 98 on holdout points interspersed within the same landscapes on which the model was trained. In order to assess the impact of
 99 the tile-based holdout on model performance, we fit a model with identical parameters but with testing data withheld fully
 :00 randomly, ignoring tile structure, comprising 25% of sample points across the state.

:01

:02 **Table 1. Holdout areas for model testing**

:03

Ecoregion	acres	Holdout area			Holdout points	
		%	acres	tiles	wetland	background
Atlantic Coastal Pine Barrens	488,314 ac	17%	80,647	4	181	293
Northeastern Coastal Zone	3,134,117 ac	26%	813,232	9	1637	2959
Northeastern Highlands	1,567,251 ac	21%	332,372	5	459	1484
Total	5,189,682 ac	21%	1,226,251	18	2277	4736



.04

.05 **Figure 3. EPA Level 3 Ecoregions for Massachusetts with Hold Out Areas**

.06 **2.3 Model fitting**

.07 The wetland probability score was generated using a random forest (RF) model, using the canonical randomForest package
.08 in R (Liaw and Wiener, 2002). Before fitting the statewide model, researchers generated a succession of preliminary models
.09 covering progressively larger areas. Preliminary models were used as sanity checks, by examining out-of-bag (OOB, see
.10 section **2.4 Model validation and testing** below) metrics, visual inspection of overlays with the DEP wetland map, aerial
.11 imagery, and field wetland delineations, and comparison with landscapes known to the researchers.

.12

.13 The first statewide model was fit using all 63 predictors (23 variables, each at 1-4 scales). In order to reduce model
.14 complexity and processing time, feature selection was carried out using a holdout variable importance procedure that
.15 compares full ensembles fit with and without the predictor, using the randomForestSRC package (Ishwaran and Kogalur,
.16 2007). Through this process 34 predictors were found to be making no contribution to model accuracy. The reduced model



.17 was fit with the remaining 29 predictors (Table 2), and compared to the full model to assess any impact on out-of-bag
.18 accuracy.

.19

.20 While the more conservative and time-intensive holdout variable importance procedure was used for variable selection, the
.21 more widely used permutational variable importance measure was used to understand the relative influence of the 29
.22 variables included in the final model (Belgiu and Drăguț, 2016). Permutational variable importance compares the OOB
.23 prediction accuracy using the original variable against the accuracy with the variable randomly permuted. The variable
.24 importance is calculated as the mean decrease in accuracy under permutation. As the stochastic nature of the random forest
.25 algorithm makes variable importance rankings shift across multiple runs of the model, nine additional models were fit with
.26 identical parameters, and variable importance measures extracted from each, in order to provide a sense of range and
.27 stability of the importance measure.

.28

.29 The use of the DEP wetland map as training data meant that the training data lacked true absence points, necessitating the
.30 assumption that some of the ostensibly non-wetland points would fall in unmapped, cryptic wetlands. The use of background
.31 or pseudo-absence points is a common situation in the domain of species distribution modelling (SDM), which generally
.32 assumes that background points do not reliably represent true absence. Random forest is also used extensively in SDM, and
.33 researchers in that field have evaluated multiple approaches for dealing with the lack of true absence points. Among the best
.34 performing and simplest of these approaches is down sampling, or the downward adjustment of the bootstrap sample to
.35 match the number of presence points (Valavi et al., 2021, 2022). The default bootstrap size for the randomForest function
.36 used in this study is 0.632 of the number of observations in the training data. Use of the default setting would have produced
.37 a bootstrap sample size of 14,438. To implement down sampling, the bootstrap sample of the random forest model was set to
.38 7586, the number of wetland points in the training data. In order to assess any effect of down sampling on model
.39 performance, the full model was also fit without down sampling, i.e. with default bootstrap settings.

.40 **2.4 Model validation and testing**

.41 Model performance was assessed in several ways, using both out-of-bag prediction and test data. Out-of-bag prediction is a
.42 strength of the random forest approach, made possible by the bootstrap sampling process. At each iteration, a decision tree is
.43 grown using a subsample of the training data. Out-of-bag refers to the data which was not part of a given bootstrap sample,
.44 and therefore not used to grow the associated decision tree. Each decision tree can therefore be validated by using it to
.45 predict the OOB data, as they have not been used in its training. Averaging the error rate over all trees in the forest gives the
.46 OOB error rate for the model. While OOB metrics are sometimes used in place of external holdout testing data, for this
.47 project testing with external data was performed as well.

.48



.49

.50

.51

.52

Table 2. Predictors used in the final model

Name	Description	Type	Source Data	Scales	In Final Model	Scale Unit
Depth_ Restrictive	Estimated Depth to a Restrictive Soil Layer	Soil	NRCS SSURGO	4, 12, 60	12	meters
Depth_ Water	Estimated Depth to Mean Water Table	Soil	NRCS SSURGO	4, 12, 60	4, 12	meters
Hydric_ Soils	Percent Hydric Soils in Map Unit	Soil	NRCS SSURGO	4, 12, 60	4, 12, 60	meters
KSat	Estimated Hydraulic Conductivity	Soil	NRCS SSURGO	4, 12, 60	60	meters
Surf_Geo	Surficial Geology	Soil	MassGIS Surficial Geology	4	4	meters
B2	Aerial Imagery: Green Band	Spectral	MA 2021 Leaf Off Aerial Imagery	4, 12, 60	12	meters
B4	Aerial Imagery: Near Infrared Band	Spectral	MA 2021 Leaf Off Aerial Imagery	4, 12, 60	60	meters
DEV	Deviation in Elevation Relative to the Neighborhood Pixels	Terrain	LiDAR Derived Elevation	4, 12, 60, 300, 1000	12, 60	meters
DIF	Difference in Elevation Relative to the Neighborhood Pixels	Terrain	LiDAR Derived Elevation	4, 12, 60, 300, 1000	12, 60	meters
DTW	Cartographic Depth to Water	Terrain	LiDAR Derived Elevation	250, 1000, 4000, 16000	250, 1000, 4000	catchment area (pixels)
Slope	Slope of the Topography in a Neighborhood	Terrain	LiDAR Derived Elevation	4, 12, 60, 300, 1000	12, 60	meters



NDVI_ Mean	Mean of the Landsat 8 Normalized Difference Vegetation Index Composite	Vegetation	LANDSAT 8 from Google Earth Engine	30	30	meters
NDVI_ Median	Median of the Landsat 8 Normalized Difference Vegetation Index Composite	Vegetation	LANDSAT 8 from Google Earth Engine	30	30	meters
NDWI_ Mean	Mean of the Landsat 8 Normalized Difference Water Index Composite	Vegetation	LANDSAT 8 from Google Earth Engine	30	30m	meters
Veg_ Height	Vegetation Height	Vegetation	LiDAR Derived Vegetation Height	4, 12, 60	4m, 12m, 60m	meters

.53

.54

.55 As the model produces a 0-1 probability score and not a classification, in order to calculate accuracy metrics it is necessary
 .56 to collapse the probability score into a dichotomous outcome (i.e. wetland/upland). The standard probability threshold of .5
 .57 was used to bin pixels into wetland and non-wetland categories. Predictions can then be compared to the MDW. Model
 .58 accuracy was explored through several metrics: overall accuracy, and true/false positive and negative rates. It is important to
 .59 remember that true and false are here relative to the reference layer, and that some amount of over-prediction relative to that
 .60 layer (i.e. identification of additional wetlands) is a fundamental goal of the project. The confusion matrix captures values
 .61 for true and false negatives and positives in a table that crosses the values of the reference data with the model predictions
 .62 (Ting, 2010). Accuracy and confusion matrices were also calculated for alternate parameterizations and sampling strategies
 .63 in order to understand any impacts of different strategies on model performance.

.64

.65 Additionally, the model was evaluated with a metric designed specifically for probabilistic binary classifiers, which does not
 .66 require collapse of continuous probability scores into dichotomous outcomes. The measure known as *area under the receiver*
 .67 *operating characteristics curve* (AUC) provides a metric that has been found useful for binary classifiers, including in digital
 .68 wetland mapping (Maxwell et al., 2016). The curve from which the AUC is derived plots the true positive rate against the
 .69 false positive rate across a spectrum of probability thresholds. For the purposes of interpretation, the AUC is equivalent to
 .70 the probability that the model will assign a higher probability score to a randomly chosen wetland point than to a randomly
 .71 chosen background point (Fawcett, 2006). An AUC of 1 indicates perfect prediction, while an AUC of .5 would mean model
 .72 predictions are totally random. For our purposes, its utility is in generating a single measure of probabilistic model



.73 performance. While there are no hard and fast thresholds for AUC scores in remote sensing, scores above 0.70 are often
.74 judged acceptable, above .80 good, and above .90 excellent ([Hosmer et al., 2013](#)).

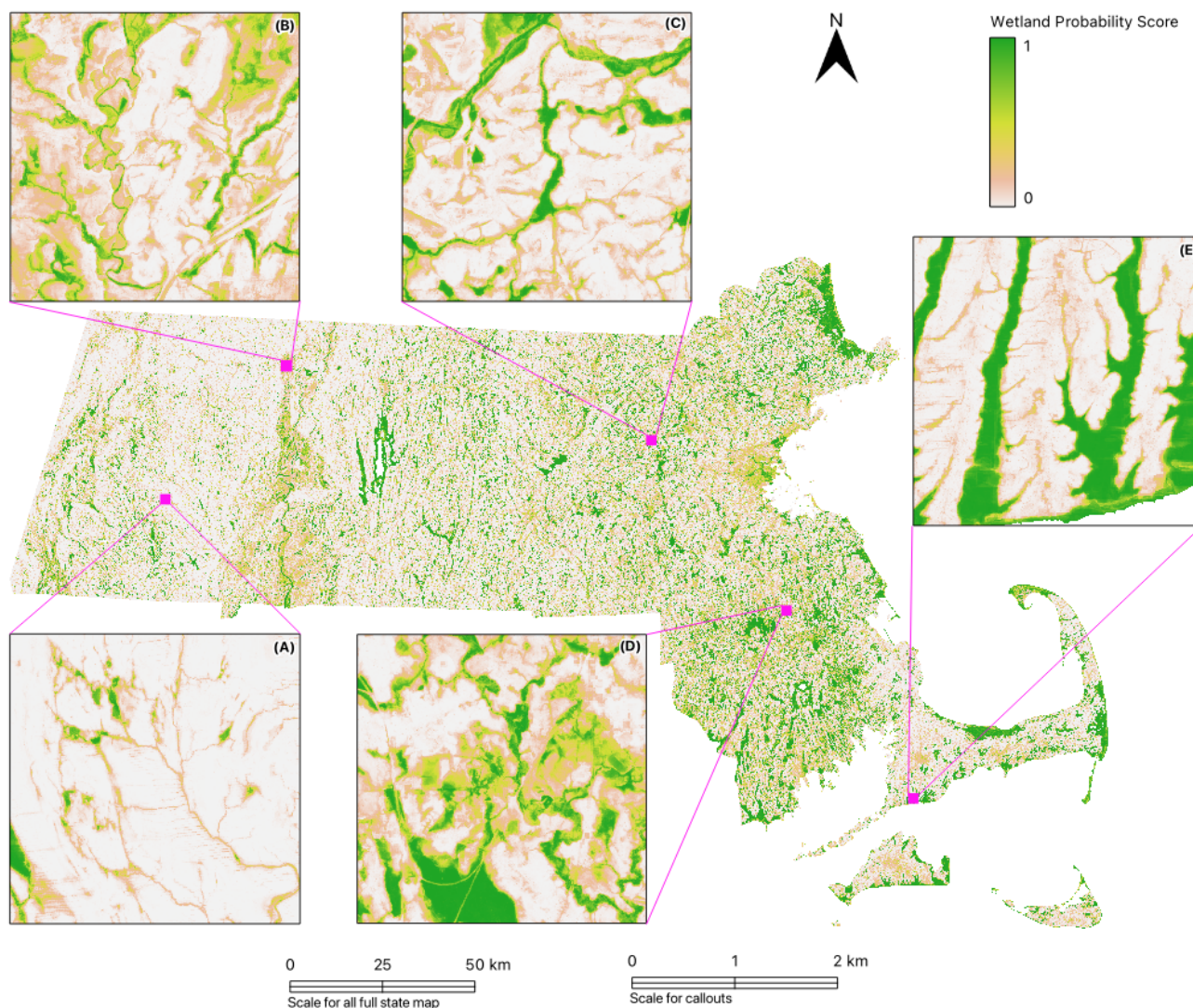
.75
.76 In the interest of understanding how the model performed on different wetland types and land covers, DEP wetland types
.77 were crosswalked with categories of land cover in the MassGIS 2016 Land Cover/Land Use dataset. The crosswalk was
.78 supplemented with coastal wetlands data accessed through the Northeast Oceans Data Portal, which were used to assign
.79 predicted additional wetlands to the category of salt marsh (Northeast Ocean Data Portal, 2024). In this way, model
.80 predictions of previously mapped DEP wetland classes could be compared with model predictions of additional wetlands
.81 that, as they do not appear in the DEP wetland inventory, have not been assigned a wetland type.

.82
.83 In order to assess impact of model predictions on state-level wetland extents, the wetland probability map was post-
.84 processed using the MA 2016 High Resolution Land Cover dataset to remove the land cover classes of impervious and
.85 barren, which were identified as areas where the prediction was likely misidentifying wetlands. Additionally, the areas
.86 identified as Open Water in the MDW were removed.

.87 **3 Results**

.88 Initial models were fit with 500 trees. After plots of error rate against number of trees showed that the model achieved
.89 stability after approximately 100 trees, subsequent models were fit with 400 trees to reduce processing time. Initial models
.90 were fit with 63 predictors, and after feature selection using holdout variable importance, a reduced model was fit with 29
.91 predictors. After confirming equivalent performance with the reduced model, the latter was selected as the final model.
.92 Model output of the Wetland Probability Score for the entire state is shown in Fig. 4.

.93
.94



.95

.96 **Figure 4. Wetland Probability Score for Massachusetts.** This map reflects the raw WPS without post-processing to
.97 remove conflicting land use designations, e.g. open water and impermeable surfaces. From west to east, callout (A) is in the
.98 Berkshires Hills of the Northeastern Highlands ecoregion; callouts (B-C) are in the Northeastern Coastal Zone, with (B) in
.99 the Connecticut River Valley, (C) in the Gulf of Maine Coastal Plain, and (D) in the Narragansett/Bristol Lowland; and
.00 callout (E) is in Atlantic Coastal Pine Barrens. Callout (C) shows the same area represented in Fig. 2. Multi-scale predictors.

.01 3.1 Accuracy and Prediction Outcomes

.02 Model performance met or exceeded all expectations for prediction. Generating predictions on the held-out testing data, and
.03 a dichotomized prediction (wetland/upland at a probability threshold of .5), the model achieved an overall accuracy rate of



04 92%, with a true positive rate of 87% and a true negative rate of 94%. Prediction accuracy based on the testing data was
 05 nearly identical to the out-of-bag accuracy score generated through the model fitting process. Performance was similar
 06 among all comparison models (Table 3). The probabilistic metric AUC was 0.971, i.e. there is a 97% chance that the model
 07 would assign a higher probability score to a randomly selected wetland point than to a randomly selected non-wetland point.
 08 While there are no hard and fast thresholds for acceptability and excellence for AUC scores, by any standard 0.971 indicates
 09 outstanding performance (Hosmer et al., 2013).

10

11 **Table 3. Confusion matrices and performance metrics for final model and comparison models**

	Train/test	Bootstrap	Type of test	Model prediction	Reference data (DEP Wetlands)				Metrics	
					Counts		Rates		Accuracy	AUC
					0	1	0	1		
Final model 29 predictors	Tile-based holdout	Down sampled	Holdout	0	4476	289	95%	13%	0.922	0.971
				1	259	1977	5%	87%		
			Out-of-bag	0	14459	775	93%	11%	0.918	
				1	1068	6155	7%	89%		
Comparison models 63 predictors	Random holdout	Down sampled	Holdout	0	4770	351	94%	13%	0.919	0.973
				1	244	2010	6%	87%		
			Out-of-bag	0	14210	745	93%	11%	0.918	
				1	1061	6067	7%	89%		
	Tile-based holdout	Default sampling	Out-of-bag	0	14457	790	93%	11%	0.918	
				1	1077	6339	7%	89%		

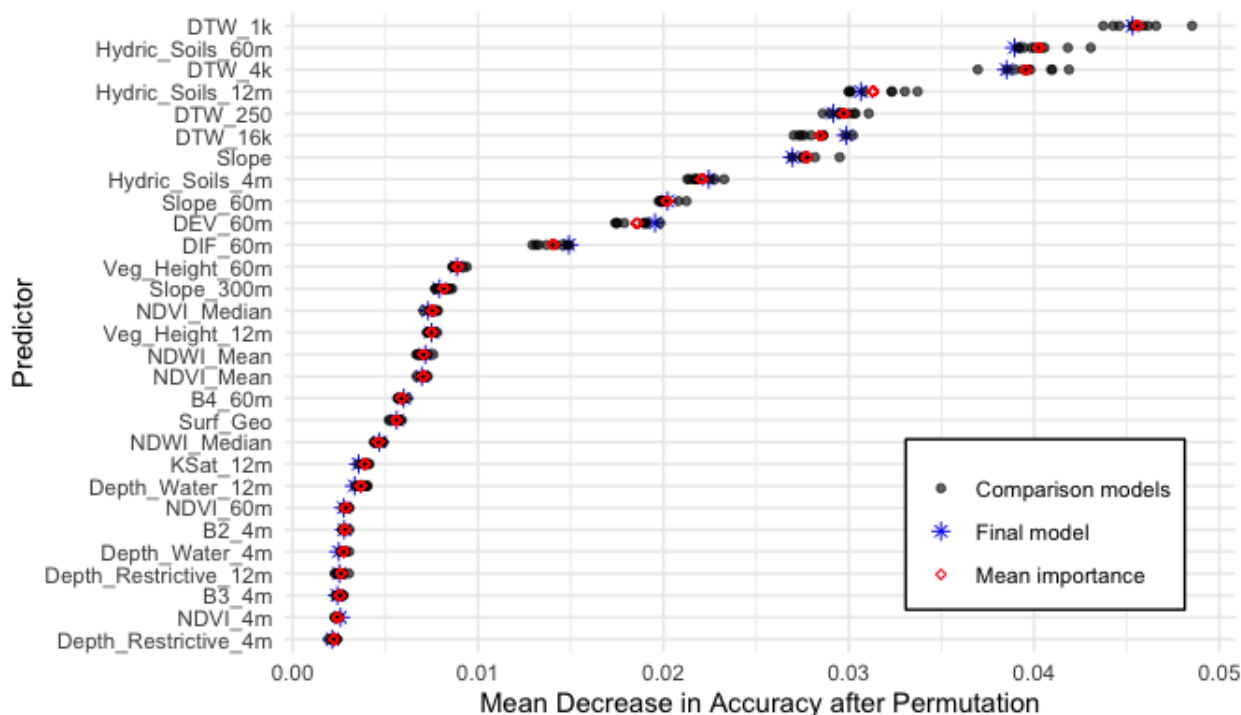
True Neg	False Neg
False Pos	True Pos

12



13 **3.2 Variable Importance**

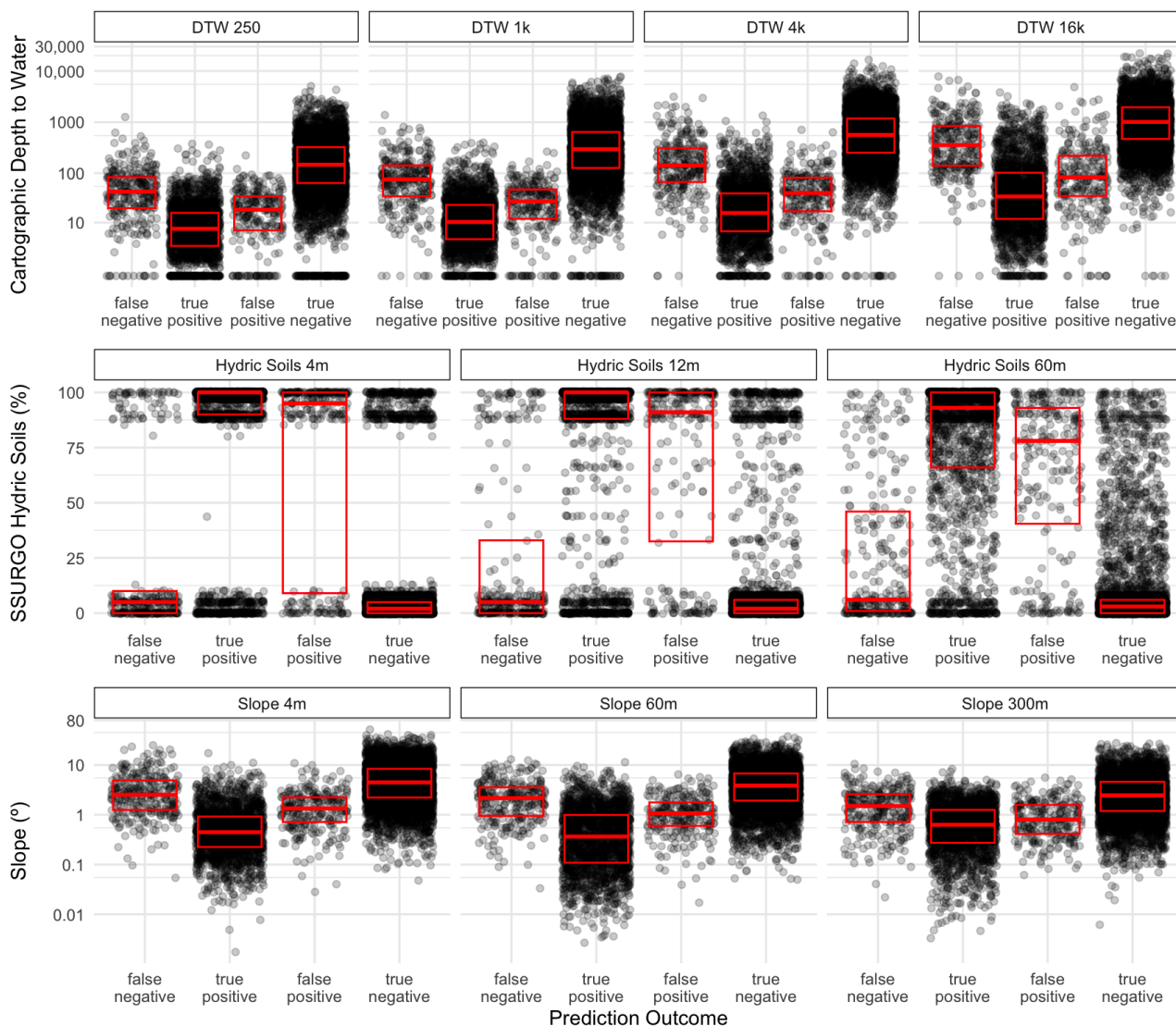
14 Terrain and soil variables were the most decisive predictors of wetland presence (Fig. 5). The LIDAR-derived cartographic
 15 Depth to Water (DTW), and the estimated percent of hydric soils from SSURGO, at various scales comprised the top six
 16 predictors and seven out of the top eight. Examination of the range of importance scores (mean decrease in accuracy) for
 17 each predictor, across 10 model fits, shows clusters of stability within rankings that shifted from iteration to iteration.
 18 DTW_1k held the top rank in every iteration, while Hydric_Soils_60m and DTW_4k each consistently held either the 2nd or
 19 3rd positions. Future investigations may focus primarily on the highest-performing variables from this predictor set.
 20
 21
 22



23
 24
 25 **Figure 5. Variable importance in the final model.** Permutational variable importance of the 29 variables used in the final
 26 model, across 10 models fit with identical parameterization. Plotting importance across multiple models allows for the
 27 assessment of relative importance while accounting for stochastic variability across model fittings. The red diamond
 28 indicates the mean importance for that variable across all 10 models, and the blue asterisk shows the importance within the
 29 final model itself.
 30



31 Examining the relationship between important variables and prediction outcomes yields some understanding of patterns in
32 over- and under-prediction. Figure 6 shows the distribution of predictor values across prediction outcomes, i.e. true and false
33 positive, and true and false negative, for the three most important multi-scale predictors: cartographic depth to water, percent
34 hydric soils, and slope. DEP wetland and background points show consistent separation across important predictors, with
35 little or no overlap of interquartile ranges between true positives and true negatives. Unsurprisingly, atypical values for these
36 predictors lead to divergent classifications. For example, background (nominally non-wetland) points typically have
37 relatively higher cartographic depth to water, so background points with lower DTW values are more frequently over-
38 predicted by the model. However it must be remembered that divergent classification is relative to the MDW and some
39 amount of divergence—especially over-prediction—is a fundamental goal of the analysis.



40

41 **Figure 6. Distribution of predictor values across prediction outcomes for the three most important multi-scale**
 42 **predictors: cartographic depth to water, percent hydric soils, and slope.** Prediction outcomes are based on binning the
 43 wetland probability score into a dichotomous (wetland/upland) outcome at a .5 threshold and comparing with the 2005
 44 MassDEP Wetland layer. Red boxes show the median and interquartile range of predictor values. Results show that atypical
 45 values for important predictors lead to classifications that diverge from the reference layer.



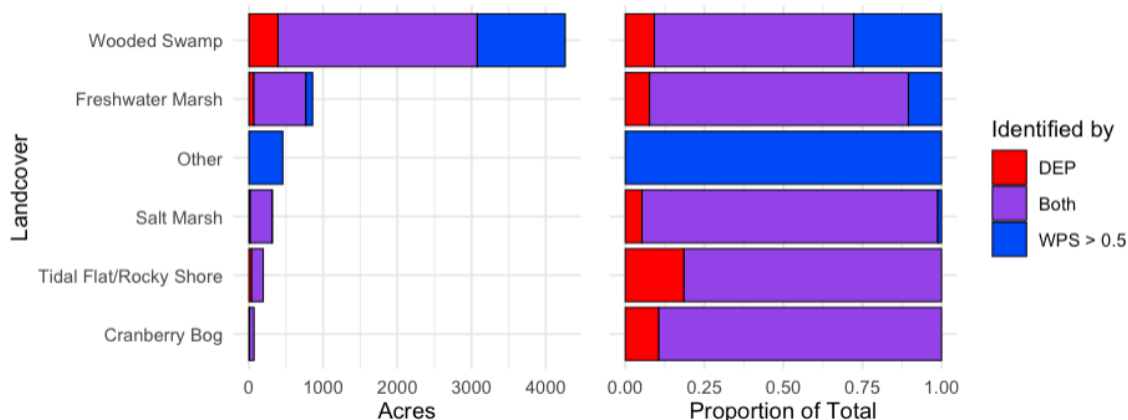
46

47 **3.3 Model performance by wetland type and landcover**

48 DEP wetland types (as given in the MDW) were crosswalked with land cover categories in the MassGIS 2016 Land
 49 Use/Land Cover layer to facilitate comparison with model predictions (Fig. 7). To supplement the land cover crosswalk,
 50 predicted wetlands that overlapped with coastal wetlands appearing in the Northeast Oceans dataset were assigned to the Salt
 51 Marsh category. The model excelled at identifying additional forested wetland areas (i.e. ‘Wooded Swamp’), adding 167,329
 52 acres of probable wetland, representing a 46% increase over the state inventory. Conversely, the model only identified 84%
 53 (57,518 acres) of the mapped forested wetland area appearing in the DEP inventory, suggesting that the even combined DEP
 54 and model results might seriously underestimate the extent of some cryptic wetlands. To existing estimates of freshwater
 55 marsh, model predictions added 96,443 acres, a 12% increase, and correctly identified 91% of DEP-designated freshwater
 56 marsh area. For each of the other crosswalked categories (excluding ‘Other’), the model identified some but not all of the
 57 DEP-designated wetland area: Salt Marsh (94%), Cranberry Bog (93%), and Tidal Flat/Rocky Shore (46%). Results in the
 58 ‘Other’ category include land use designations *Impervious*, *Developed Open Space*, *Bare Land*, *Pasture/Hay*, and *Cultivated*,
 59 and likely encompass many historic wetland areas that have been lost to agricultural, residential, or urban development.

60

61



62

63

64 **Figure 7. Comparison of DEP wetland layer and model predictions across wetland types.** Comparison of raw Wetland
 65 Probability Score (dichotomized at 0.5) with DEP designation. Landcovers in the 'Other' category include, in descending
 66 order of acreage, Impervious, Developed Open Space, Bare Land, Pasture/Hay, Cultivated, Water, Palustrine Aquatic Bed,
 67 and Unconsolidated Shore.



68 **3.4 Impact on estimates of wetland extent**

69 After post-processing to remove impervious surface and barren lands, the model identified an additional 225,781 acres of
 70 probable wetland compared to the MDW inventory—an increase of 40% (Table 4, ecoregions shown in Fig. 3). The
 71 Northeastern Highlands had the highest proportional increase at 49%, likely due to the dominance of the forested highland
 72 landscapes for which wetland detection via visual assessment methods is difficult. The Atlantic Coastal Pine Barrens,
 73 dominated by sandy soils, had the smallest absolute and proportional increase in wetland area. The 15% of additional
 74 wetland area here represents approximately 1/3rd of the proportional increase of the other two ecoregions.

77 **Table 4. Additional wetland area identified across the ecoregions of Massachusetts**

Ecoregion	2005 DEP	2023 Additional predicted (at .5 threshold)	Total	Percent increase
Atlantic Coastal Pine Barrens	78,402 ac	11,430 ac	89,832 ac	15%
Northeastern Coastal Zone	396,695 ac	171,085 ac	567,780 ac	43%
Northeastern Highlands	88,153 ac	43,266 ac	131,419 ac	49%
Total	563,250 ac	225,781 ac	789,031 ac	40%

80 **4 Discussion**

81 **4.1 Overview**

82 A random forest model was trained on the 2005 MassDEP Wetlands layer (MDW), using widely available open source
 83 software and publicly available data, and generated a Wetland Probability Score for the entire state of Massachusetts at a 4m
 84 resolution. The model has performed well by all metrics, and provides support for the use of machine learning to quickly and
 85 efficiently generate large-scale maps of wetland probabilities that can serve as an adjunct to conventional field and desk
 86 methods. The balance of over- and under-prediction relative to the MDW suggests that model is fit for the purpose of
 87 providing a more inclusive estimate of wetland extent in Massachusetts by modeling the extent of additional wetlands that
 88 are not captured in the MDW. Our results suggest a promising avenue for bootstrapping existing wetland inventories by
 89 using them as training data to identify additional probable wetland acres.



91 Adding 46% to existing estimates of forested wetlands in the state, results suggest that models such as this may be especially
92 useful for identifying cryptic, previously unmapped forested wetlands. These wetlands are especially difficult to identify via
93 the methods used by both NWI and the MDW. Efforts to combat climate change make it critical to understand the extent and
94 location of forested wetlands, as they represent greater stores of aboveground carbon than other wetland types and greater
95 stores of soil organic carbon than other forests. The model also predicted 69,761 acres of wetland in developed and
96 agricultural areas (i.e. the 'Other' category of Fig. 7), entirely absent from the MDW, suggesting that this model may also be
97 of use to understand long-term land use change, by estimating the extent and location of wetlands that have been lost to
98 development of various kinds. The identification of topographically suitable areas with a high wetland probability may also
99 be useful for the identification of potential sites for wetland restoration.

.00

.01 Examination of model performance and variable importance together provide support for the use of multi-scale terrain
.02 derivatives in wetland detection, and further expansion of multi-scalar approaches to other types of predictors. Variable
.03 importance scores show that model performance was dependent on the inclusion of terrain derivatives (e.g. cartographic
.04 depth to water and slope) at multiple scales, reinforcing the efficacy of multi-scale terrain derivatives established in previous
.05 work on wetland detection (Halabisky et al., 2023). More novel results include the positive performance of non-terrain multi-
.06 scale predictors, principally the SSURGO-derived Hydric Soils. The model was improved by inclusion of Hydric Soils at
.07 three different scales, and the larger-scale versions (smoothed at 12m and 60m resolutions) scored higher on importance than
.08 the baseline 4m scale. This finding provides support for the use of multi-scale smoothing with soil map data, both to account
.09 for neighborhood effects and to ameliorate the issues caused by crisp boundaries and their tendency to obscure underlying
.10 heterogeneity and varying transition zones.

.11 **4.2 Limitations**

.12 The use of the MDW for training and testing has implications for interpretation of model results and accuracy metrics: all
.13 accuracy metrics are relative to the reference dataset. Over-prediction, or false positives, should therefore be assumed to be a
.14 mix of false positives and real, cryptic wetlands that are missing from MDW. Likewise, under-prediction or false negatives
.15 are inevitably a mix of the model failing to identify an existing wetland, and false positives in the MDW where the model
.16 correctly identifies an upland. We expect the MDW itself is much more prone to false negatives than false positives, and that
.17 over-prediction by our model is therefore more informative than under-prediction.

.18

.19 Results show that the model failed to identify significant portions of DEP-designated wetlands: Tidal Flat/Rocky Shore
.20 (54%), Wooded Swamp (16%), Freshwater Marsh (9%), Cranberry Bog (7%), and Salt Marsh (6%). Researchers were also
.21 able to identify some patterns in over- and under-prediction through visual inspection of results and comparison with
.22 familiar sites and a small number of digitized field delineations. The model is prone to identifying low lying urban areas and



.23 drained agricultural fields as wetlands. As noted above, these areas are topographically suited to be wetlands and prior to
.24 development they likely were. The model is also under-predicting certain wetland configurations, including sloped wetlands
.25 and small wetlands high in the drainage network. Future efforts to better capture these wetland types and/or environmental
.26 configurations could include oversampling these areas or fitting a weighted model.

.27 **4.3 Next Steps for Machine Learning and Wetland Mapping**

.28 The strong performance and relative speed of this approach suggest many useful avenues for further development. Releasing
.29 a publicly available GIS layer based on the Wetland Probability Score could be an aid to researchers as well as a useful
.30 screening tool for the town conservation commissions who, in Massachusetts, are charged with assessing the veracity of
.31 wetland delineations submitted by project proponents, as well as with desktop evaluation of potential wetlands requiring
.32 field delineation. Moving forward, this model or ones very similar to it are likely to be critical in large-scale assessments of
.33 soil carbon across the landscape—not only in wetlands (Stewart et al., 2024). It's also possible that a very similar model, with
.34 additional development, could be used to predict not only wetland extent but also wetland type.

.35
.36 In the domain of model performance and overall value, this project opens up several pathways for exploration and further
.37 improvement. First and foremost, the integration of field data into the training and/or validation process would almost
.38 certainly improve the performance of the model as well as reinforcing its validity for researchers, policymakers, and the
.39 public. One potential avenue for acquiring data at the necessary scale would be creating a centralized repository and unified
.40 format for field delineation data gathered for wetland permit applications. Conversely, a map of wetland probabilities could
.41 be used to inform a field sampling plan for ongoing carbon inventory and accounting (Stewart et al., 2024).

.42
.43 The development of this model was focused narrowly on the goal of estimating the extent of cryptic wetlands in aggregate
.44 (and by implication providing proof of concept for digital wetland mapping using machine learning), rather than producing
.45 the most accurate predictions on a pixel-by-pixel basis, or producing the fastest or most parsimonious model. As such,
.46 results suggest several low hanging fruits for further investigation and improvement of model efficiency and performance.
.47 For example, it is possible to generate variable importance metrics on a pixel-by-pixel basis in the training and/or testing
.48 data. These case-level metrics would support more detailed understanding of key topics such as which predictors support the
.49 identification of atypical wetlands currently underpredicted by the model, and potentially improving predictive performance
.50 with variable weighting and/or oversampling strategies. Additionally, there are several other ways to further refine the
.51 feature selection process, such as testing the importance of highly correlated predictors by fitting models with only the single
.52 most important predictor from each family, i.e. terrain, vegetation, soils, and so on. This could also be applied to a deeper
.53 investigation of the impact of multi-scale approach, by fitting models with only the single most important scale for each
.54 predictor and assessing model performance. In the interest of identifying the most lightweight and parsimonious model that



.55 can achieve required levels of accuracy, these investigations should be accompanied by assessment of other strategies for
.56 tuning model hyperparameters such as number of trees, size of bootstrap sample, and number of variables selected for each
.57 decision tree.

.58 **5 Conclusion**

.59 The status of wetlands is a critical issue that is only becoming more important as the climate crisis intensifies. Monitoring,
.60 managing, and protecting these vital ecosystems, safeguarding the many functions they perform, all depend on a strong
.61 understanding of their extent and location. This is a critical issue for cryptic wetlands that are missing from most inventories,
.62 often due to forest land cover, which represent especially dense pools of organic carbon. Fast and efficient processes for
.63 detecting probable wetlands, using free, open source software and publicly available data, are a needed adjunct to
.64 conventional methods for wetland mapping and monitoring.

.65
.66 *Competing interests:* The authors declare that they have no conflict of interest.

.67
.68 *Author contributions.* RSF and SG designed the data pre-processing and sampling protocol. SG and EG conceptualized and
.69 carried out all operations in QGIS: data acquisition, pre-processing, sampling, and post-processing. RF designed and carried
.70 out all operations in R, including data wrangling and model training, testing, and prediction processes, and data
.71 visualizations. KZD provided project administration and support. RSF prepared the paper with contributions from all
.72 authors.

.73
.74 *Acknowledgments.* The authors would like to acknowledge Jenny Watts, Gillian Davies, and Matthew von Walde for their
.75 input and support early in the development of this project. Initial development of the wetland probability map was conducted
.76 as part of a project for the Massachusetts Department of Environmental Protection. Contents of this article are solely the
.77 responsibility of the authors; the views and conclusions contained in this document are those of the authors and should not be
.78 interpreted as representing the opinions or policies of the government of Massachusetts or any of its representatives.

.79
.80 *Data and code availability:* All datasets used in this analysis are publicly available. R code is available on request.

.81 **References**

.82 Northeast Ocean Data Portal: <https://www.northeastoceandata.org/>, last access: 1 May 2024.
.83 Protecting Wetlands in Massachusetts | Mass.gov: <https://www.mass.gov/info-details/protecting-wetlands-in-massachusetts>,
.84 last access: 7 March 2024.



- 85 Baker, C. D., Polito, K. E., Beaton, M. A., and Suuberg, M.: Inland and Coastal Wetlands of Massachusetts: Status and
86 Trends, Commonwealth of Massachusetts, 2019.
- 87 Behrens, T., Viscarra Rossel, R. A., Kerry, R., MacMillan, R., Schmidt, K., Lee, J., Scholten, T., and Zhu, A.-X.: The
88 relevant range of scales for multi-scale contextual spatial modelling, *Sci. Rep.*, 9, 14800, [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-019-51395-3)
89 019-51395-3, 2019.
- 90 Belgiu, M. and Drăguț, L.: Random forest in remote sensing: A review of applications and future directions, *ISPRS J.*
91 *Photogramm. Remote Sens.*, 114, 24–31, <https://doi.org/10.1016/j.isprsjprs.2016.01.011>, 2016.
- 92 Breiman, L.: Random forests, *Mach. Learn.*, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- 93 Davidson, N. C.: How much wetland has the world lost? Long-term and recent trends in global wetland area, *Mar. Freshw.*
94 *Res.*, 65, 934, <https://doi.org/10.1071/MF14173>, 2014.
- 95 Fawcett, T.: An introduction to ROC analysis, *Pattern Recognit. Lett.*, 27, 861–874,
96 <https://doi.org/10.1016/j.patrec.2005.10.010>, 2006.
- 97 Felton, B. R., O’Neil, G. L., Robertson, M.-M., Fitch, G. M., and Goodall, J. L.: Using Random Forest Classification and
98 Nationally Available Geospatial Data to Screen for Wetlands over Large Geographic Regions, *Water*, 11, 1158,
99 <https://doi.org/10.3390/w11061158>, 2019.
- 100 Fluet-Chouinard, E., Stocker, B. D., Zhang, Z., Malhotra, A., Melton, J. R., Poulter, B., Kaplan, J. O., Goldewijk, K. K.,
101 Siebert, S., Minayeva, T., Hugelius, G., Joosten, H., Barthelmes, A., Prigent, C., Aires, F., Hoyt, A. M., Davidson, N.,
102 Finlayson, C. M., Lehner, B., Jackson, R. B., and McIntyre, P. B.: Extensive global wetland loss over the past three
103 centuries, *Nature*, 614, 281–286, <https://doi.org/10.1038/s41586-022-05572-6>, 2023.
- 104 Gale, S.: Automated Identification of Wetlands Using GIS in North Carolina, North Carolina Dept. of Environmental
105 Quality, Div. of Water Resources, 2021.
- 106 Griffith, G. E., Omernik, J. M., Bryce, S. A., Royte, J., Hoar, W. D., Homer, J. W., Keirstead, D., Metzler, K. J., and Hellyer,
107 G.: Ecoregions of New England (2 sided color poster with map, descriptive text, summary tables, and photographs), *Rest.*
108 *VA US Geol. Surv.* Scale 11325000, 2009.
- 109 Halabisky, M., Miller, D., Stewart, A. J., Lorigan, D., Brasel, T., and Moskal, L. M.: The Wetland Intrinsic Potential tool:
110 Mapping wetland intrinsic potential through machine learning of multi-scale remote sensing proxies of wetland indicators,
111 *EGUsphere*, 1–19, <https://doi.org/10.5194/egusphere-2022-665>, 2022.
- 112 Halabisky, M., Miller, D., Stewart, A. J., Yahnke, A., Lorigan, D., Brasel, T., and Moskal, L. M.: The Wetland Intrinsic
113 Potential tool: mapping wetland intrinsic potential through machine learning of multi-scale remote sensing proxies of
114 wetland indicators, *Hydrol. Earth Syst. Sci.*, 27, 3687–3699, <https://doi.org/10.5194/hess-27-3687-2023>, 2023.
- 115 Hosmer, D. W., Lemeshow, S., and Sturdivant, R. X.: *Applied Logistic Regression*, John Wiley & Sons, 528 pp., 2013.
- 116 Hunter, G. J., Bregt, A. K., Heuvelink, G. B. M., De Bruin, S., and Verrantaus, K.: Spatial Data Quality: Problems and
117 Prospects, in: *Research Trends in Geographic Information Science*, edited by: Navratil, G., Springer Berlin Heidelberg,
118 Berlin, Heidelberg, 101–121, https://doi.org/10.1007/978-3-540-88244-2_8, 2009.



- 19 Ishwaran, H. and Kogalur, U. B.: Random survival forests for R, *R News*, 7, 25–31, 2007.
- 20 Jafarzadeh, H., Mahdianpari, M., Gill, E. W., Brisco, B., and Mohammadimanesh, F.: Remote Sensing and Machine
21 Learning Tools to Support Wetland Monitoring: A Meta-Analysis of Three Decades of Research, *Remote Sens.*, 14, 6104,
22 <https://doi.org/10.3390/rs14236104>, 2022.
- 23 Kudray, G. M. and Gale, M. R.: Evaluation of National Wetland Inventory maps in a heavily forested region in the upper
24 Great Lakes, *Wetlands*, 20, 581–587, [https://doi.org/10.1672/0277-5212\(2000\)020\[0581:EONWIM\]2.0.CO;2](https://doi.org/10.1672/0277-5212(2000)020[0581:EONWIM]2.0.CO;2), 2000.
- 25 Liaw, A. and Wiener, M.: Classification and Regression by randomForest, *R News*, 2, 18–22, 2002.
- 26 Martin, G. I., Kirkman, L. K., and Hepinstall-Cymerman, J.: Mapping Geographically Isolated Wetlands in the Dougherty
27 Plain, Georgia, USA, *Wetlands*, 32, 149–160, <https://doi.org/10.1007/s13157-011-0263-7>, 2012.
- 28 Matthews, J. W., Skultety, D., Zercher, B., Ward, M. P., and Benson, T. J.: Field Verification of Original and Updated
29 National Wetlands Inventory Maps in three Metropolitan Areas in Illinois, USA, *Wetlands*, 36, 1155–1165,
30 <https://doi.org/10.1007/s13157-016-0836-6>, 2016.
- 31 Maxwell, A. E., Warner, T. A., and Strager, M. P.: Predicting Palustrine Wetland Probability Using Random Forest Machine
32 Learning and Digital Elevation Data-Derived Terrain Variables, *Photogramm. Eng. Remote Sens.*, 82, 437–447,
33 <https://doi.org/10.14358/PERS.82.6.437>, 2016.
- 34 Miller, B. A., Koszinski, S., Wehrhan, M., and Sommer, M.: Impact of multi-scale predictor selection for modeling soil
35 properties, *Geoderma*, 239–240, 97–106, <https://doi.org/10.1016/j.geoderma.2014.09.018>, 2015.
- 36 Mitsch, W. J. and Gosselink, J. G.: *Wetlands*, John Wiley & Sons, 744 pp., 2015.
- 37 Murphy, P., Ogilvie, J., and Arp, P.: Topographic modelling of soil moisture conditions: A comparison and verification of
38 two models, *Eur. J. Soil Sci.*, 60, 94–109, <https://doi.org/10.1111/j.1365-2389.2008.01094.x>, 2009.
- 39 Murphy, P. N. C., Ogilvie, J., Meng, F.-R., White, B., Bhatti, J. S., and Arp, P. A.: Modelling and mapping topographic
40 variations in forest soils at high resolution: A case study, *Ecol. Model.*, 222, 2314–2332,
41 <https://doi.org/10.1016/j.ecolmodel.2011.01.003>, 2011.
- 42 Nahlik, A. M. and Fennessy, M. S.: Carbon storage in US wetlands, *Nat. Commun.*, 7, 1–9,
43 <https://doi.org/10.1038/ncomms13835>, 2016.
- 44 Nikiforova, A. A., Fleis, M. E., Nyrtsov, M. V., Kazantsev, N. N., Kim, K. V., Belyonova, N. K., and Kim, J. K.: Problems
45 of modern soil mapping and ways to solve them, *CATENA*, 195, 104885, <https://doi.org/10.1016/j.catena.2020.104885>,
46 2020.
- 47 The Wetland Identification Model (WIM) – A New Arc Hydro Functionality for Predicting Wetland Locations Using
48 LiDAR Elevation Data and Machine Learning: [https://community.esri.com/t5/water-resources-blog/the-wetland-
49 identification-model-wim-a-new-arc/ba-p/884298](https://community.esri.com/t5/water-resources-blog/the-wetland-identification-model-wim-a-new-arc/ba-p/884298), last access: 23 January 2024.
- 50 O’Neil, G. L., Goodall, J. L., and Watson, L. T.: Evaluating the potential for site-specific modification of LiDAR DEM
51 derivatives to improve environmental planning-scale wetland identification using Random Forest classification, *J. Hydrol.*,
52 559, 192–208, <https://doi.org/10.1016/j.jhydrol.2018.02.009>, 2018.



- 53 QGIS Development Team: QGIS Geographic Information System, 2023.
- 54 R Core Team: R: A language and environment for statistical computing, 2023.
- 55 Rapinel, S., Panhelleux, L., Gayet, G., Vanacker, R., Lemercier, B., Laroche, B., Chambaud, F., Guelmami, A., and Hubert-
56 Moy, L.: National wetland mapping using remote-sensing-derived environmental variables, archive field data, and artificial
57 intelligence, *Heliyon*, 9, e13482, <https://doi.org/10.1016/j.heliyon.2023.e13482>, 2023.
- 58 Soil Survey Staff: SSURGO Metadata - Table Column Descriptions Report, Natural Resources Conservation Service, United
59 States Department of Agriculture, Washington, D.C, 2015.
- 60 Soil Survey Staff: Soil Survey Geographic (SSURGO) Database for Massachusetts, n.d.
- 61 Stewart, A. J., Halabisky, M., Babcock, C., Butman, D. E., D'Amore, D. V., and Moskal, L. M.: Revealing the hidden
62 carbon in forested wetland soils, *Nat. Commun.*, 15, 726, <https://doi.org/10.1038/s41467-024-44888-x>, 2024.
- 63 Stolt, M. H. and Baker, J. C.: Evaluation of national wetland inventory maps to inventory wetlands in the southern Blue
64 Ridge of Virginia, *Wetlands*, 15, 346–353, <https://doi.org/10.1007/BF03160889>, 1995.
- 65 Stone, J. R., Stone, B. D., DiGiacomo-Cohen, M. L., and Mabee, S. B.: Surficial materials of Massachusetts—A 1:24,000-
66 scale geologic map database, Scientific Investigations Map, U.S. Geological Survey, <https://doi.org/10.3133/sim3402>, 2018.
- 67 Tiner, R. W.: Use of high-altitude aerial photography for inventorying forested wetlands in the United States, *For. Ecol.*
68 *Manag.*, 33–34, 593–604, [https://doi.org/10.1016/0378-1127\(90\)90221-V](https://doi.org/10.1016/0378-1127(90)90221-V), 1990.
- 69 Ting, K. M.: Confusion Matrix, in: *Encyclopedia of Machine Learning*, edited by: Sammut, C. and Webb, G. I., Springer
70 US, Boston, MA, 209–209, https://doi.org/10.1007/978-0-387-30164-8_157, 2010.
- 71 Valavi, R., Elith, J., Lahoz-Monfort, J. J., and Guillera-Arroita, G.: Modelling species presence-only data with random
72 forests, *Ecography*, 44, 1731–1742, <https://doi.org/10.1111/ecog.05615>, 2021.
- 73 Valavi, R., Guillera-Arroita, G., Lahoz-Monfort, J. J., and Elith, J.: Predictive performance of presence-only species
74 distribution models: a benchmark study with reproducible code, *Ecol. Monogr.*, 92, e01486,
75 <https://doi.org/10.1002/ecm.1486>, 2022.