
Predicting forest carbon stocks in the contiguous U.S.

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 PLACEHOLDER

2 1 Introduction

3 An important step in addressing the risks of climate change is to protect and properly manage our
4 natural resources. Forests are one such resource with a portfolio of benefits from temperature and
5 moisture attenuation to carbon capture and storage [source]. This resource has been directly impacted
6 by humans – through urbanization and agricultural expansion – and indirectly through droughts, fires,
7 and other risks posed by climate change. Monitoring our forests and properly accounting for their
8 extent and health could help improve their management and strengthen their protection. We can grow
9 this pool through afforestation and deferred deforestation, oftentimes supported through carbon offset
10 payments by sectors emitting carbon to the landowners growing these forests.

11 In order for a payment system like this to work, we need precise, low-cost methods for estimating the
12 amount of carbon stored in a forest. These methods can be used to establish baseline estimates of
13 forest biomass, providing a counterfactual scenario to how forest growth would have changed in the
14 absence of offset payments. These methods can also be used to monitor changes in forest biomass
15 after offset payments (the difference between these scenarios is precisely the amount of carbon offset
16 by payments).

17 Decades of research and data collection by US government departments offers an opportunity for
18 such low-cost estimation methods. In this project I propose to train a model to use remote sensing
19 data collected by NASA satellites to predict forest biomass, collected through surveys by the USDA
20 FIS. The model will aim to predict 5-year survey-period forest biomass at the county-level across
21 the contiguous United States (CONUS). I will also aim to generate bootstrapped standard errors for
22 my predictions to properly account for uncertainty in forest size, as provided by sample errors in the
23 USDA FIS dataset. Standard errors can be a critical ingredient in accurately pricing carbon offsets.

24 1.1 Related work

25 Given the economic and biologic importance of forests, there is already much established literature
26 on forest biomass prediction.

27 Han et al. and Li et al. both evaluate several machine learning algorithms to estimate forest biomass,
28 using Sentinel-1 and Landsat 8 satellites [2, 3]. [[SUMMARY OF APPROACHES]]. [[SUMMARY
29 OF FINDINGS]]. Bjork et al. adopts deeper learning techniques for a similar aim: to estimate
30 forest biomass from radar backscatter [1]. [[SUMMARY OF APPROACH]]. [[SUMMARY OF
31 FINDINGS]].

32 Saarela et al. and Naik et al. evaluate newer remote sensing layers for incremental benefits to
33 forest biomass estimates [5, 4]. Saarela uses LiDAR and field data to estimate aboveground biomass
34 and associated uncertainty, while Naik uses multispectral remote sensing data [[SUMMARY OF
35 APPROACHES]]. [[SUMMARY OF FINDINGS]].

36 And Vorster et al. contributes to this body of research by exploring how allometric equations
37 can improve estimates of uncertainty in biomass estimates [6]. [[SUMMARY OF APPROACH]],
38 [[SUMMARY OF FINDINGS]].

39 **2 Dataset**

40 **3 Technical approach**

41 In this paper I aim to train a model with the lowest prediction error as determined by root mean
42 squared error (RMSE). This metric best aligns with the ultimate objective of this model, which is
43 to make the best county-level predictions of forest biomass from historic remote sensing data. For
44 robustness, I also evaluate mean absolute error (MAE) and the coefficient of determination (R^2)
45 across all models. Among a subset of the highest performing models, I use bootstrapping techniques
46 to estimate standard errors of mean predictions across these models for the final test dataset.

47 **3.1 Model selection**

48 I use cross validation to evaluate all models and arrive at the performance metrics described above.
49 I choose to hold out 10% of my data for final testing, and to use the remaining 90% for model
50 selection. Because of the potential spatial and temporal correlation of my data, I evaluate these
51 models across three different splitting regimes: First, within each fold of my cross validation, I
52 hold out [[25%]] of the county-years from each state at random and evaluate the performance of
53 out-of-county-year prediction. Second, I hold out the last forest biomass report for each county and
54 evaluate the performance of future predictions using a model trained on past data (each county has
55 between 1 and 4 total reports). And third, I perform leave-one-state-out cross validation where, in
56 each fold, I use all available time and county data from all but one state and evaluate the performance
57 of out-of-state prediction.

58 In selecting the highest-performing model, I consider three categories of models: baseline models,
59 machine learning models, and neural net models. I first train two baseline models. These models
60 serve as a standard against which I can aim to improve predication error through the use of machine
61 learning and neural net models. First, I consider using the simple global mean at each prediction
62 location. And second, I consider a simple linear regression of all predictor variables on county-level
63 forest biomass predictions.

64 I next optimize 5 machine learning models: Regularized linear regression (with Elastic Net), Simple
65 Decision Trees, AdaBoost, XGBoost, and Random Forests. I optimize the hyperparameters associated
66 with each of these models through a grid search approach. And lastly, I implement a fully connected
67 neural net model, selecting optimization parameters, as well as the number and size of hidden layers,
68 through a grid search approach.

69 My model selection is based on the highest performing model among those described in this section.

70 **3.2 Bootstrapping prediction error**

71 With the [[3]] highest performing models, I also bootstrap mean prediction values for each county
72 and the standard errors of those mean predictions. I do this by generating new training datasets by
73 taking random draws of a normal distribution centered at the USDA FIS provided estimate, with
74 standard deviation equal to the standard error provided by the USDA FIS. I then evaluate the same
75 scores I describe above, but can generate ranges of uncertainty about each. These results can be used
76 to further support the final model selection.

Split	Model	RMSE	MAE	R^2
County-year	Global mean			
County-year	Linear regression			
County-year	Elastic Net			
County-year	Decision Tree			
County-year	AdaBoost			
County-year	XGBoost			
County-year	Random Forest			
County-year	FC NNet			
Out-of-state	Global mean			
Out-of-state	Linear regression			
Out-of-state	Elastic Net			
Out-of-state	Decision Tree			
Out-of-state	AdaBoost			
Out-of-state	XGBoost			
Out-of-state	Random Forest			
Out-of-state	FC NNet			
Forward-in-time	Global mean			
Forward-in-time	Linear regression			
Forward-in-time	Elastic Net			
Forward-in-time	Decision Tree			
Forward-in-time	AdaBoost			
Forward-in-time	XGBoost			
Forward-in-time	Random Forest			
Forward-in-time	FC NNet			

Table 1: Performance results of models across splitting regimes

4 Preliminary results

After tuning each of the models, I find that [[MODEL]] results in the lowest RMSE in [[SPLITTING REGIME]]. This performance is [[XXX]] percentage points lower than the next highest model, [[MODEL]], within this splitting regime. Overall, I find that [[MODEL]] performs best overall across all splitting regimes. Results are included in Table 1

Broader impact

This research contributes to the larger field of forest biomass quantification. On the whole, improved approaches to forest monitoring will support humans' care for this natural resource.

While a central component of this paper focuses on estimate uncertainty, it is often the case for downstream users to use estimates without consideration of their uncertainty. When these forest biomass estimates are used in carbon markets or other in governmental or corporate statements without uncertainty qualifications, this can lead and has led to inflated claims of climate impact (a.k.a., "greenwashing").

Additionally, this research and the models in this paper directly benefit those countries and regions with better surveyed data. While remote sensing layers tend to have complete coverage of the Earth, the surveyed forest biomass as well as other potential predictors like soil types and crop layers tend to only be available in wealthier countries. This translates to wealthier nations having access to more precise models and downstream economic and environmental benefits. That said, the contiguous United States has a variety of climactic regions and may be representative of large parts of the world, meaning these models could support developing countries.

Beyond what is mentioned above, leveraging biases in the data is not applicable to this research.

98 References

- 99 [1] S. Björk, S. N. Anfinsen, E. Næsset, T. Gobakken, and E. Zahabu. Constructing forest biomass prediction
100 maps from radar backscatter by sequential regression with a conditional generative adversarial network.
101 *CoRR*, abs/2106.15020, 2021.
- 102 [2] H. Han, R. Wan, and B. Li. Estimating forest aboveground biomass using gaofen-1 images, sentinel-1
103 images, and machine learning algorithms: A case study of the dabie mountain region, china. *Remote Sensing*,
104 14(1), 2022.
- 105 [3] Y. Li, M. Li, C. Li, and Z. Liu. Forest aboveground biomass estimation using landsat 8 and sentinel-1a data
106 with machine learning algorithms. *Scientific Reports*, 10(1):9952, 2020.
- 107 [4] P. Naik, M. Dalponte, and L. Bruzzone. Prediction of forest aboveground biomass using multitemporal
108 multispectral remote sensing data. *Remote Sensing*, 13(7), 2021.
- 109 [5] S. Saarela, A. Wästlund, E. Holmström, A. A. Mensah, S. Holm, M. Nilsson, J. Fridman, and G. Ståhl.
110 Mapping aboveground biomass and its prediction uncertainty using lidar and field data, accounting for
111 tree-level allometric and lidar model errors. *Forest Ecosystems*, 7(1):43, 2020.
- 112 [6] A. G. Vorster, P. H. Evangelista, A. E. L. Stovall, and S. Ex. Variability and uncertainty in forest biomass es-
113 timates from the tree to landscape scale: the role of allometric equations. *Carbon Balance and Management*,
114 15(1):8, 2020.