

C. Wade Ross¹, Sabine Grunwald¹, Jason Vogel², Allan Bacon¹, Eric J. Jokela³, Rosvel Bracho-Garrilo³, Daniel Markewitz⁴, Madison Akers⁴, Joshua Cucinella³, Andy Laviner⁵, Tom Fox⁵, Tim Martin³, Carlos Gonzalez⁶

¹ University of Florida, Department of Soil and Water Science, ²Texas A&M – Ecosystem Science and Management, ³University of Florida, School of Forest Resources and Conservation, ⁴University of Georgia, Warnell School of Forestry and Natural Resources, ⁵Virginia Polytechnic Institute and State University, ⁶Oregon State University

Rationale

Soils of the US Southeast are *estimated* to store between 8.9 and 51.2 Pg carbon (C), accounting for more than 1/3 of the total soil C storage for the conterminous US when using the median value. The large range of estimates is due to 1) high spatial variance of soil C, 2) difficulties associated with collecting and analyzing enough samples to capture this variance across large regions, and 3) estimating soil C from highly aggregated legacy datasets with limited sample support. Additionally, many studies do not attempt to quantify C below the top soil (0-20 cm). To improve soil C estimates to 1m depth in forested ecosystems of the US Southeast, we used machine learning (Random Forests) with measured soil properties from the Tier 2 network in conjunction with a large suite of environmental data. This *data mining* approach revealed relationships between soil C and environmental covariates.

Materials & Methods

The concentration of soil C was analyzed across the Tier 2 network (N ~ 322) at four depth profiles: 0-10, 10-20, 20-50, and 50-100 cm. To predict soil C to 1m, we aggregated these depth-based measurements to 1m with *smoothing splines*. We acquired nearly **7 terabytes of data** from a multitude of agencies to represent the 8 factors in the STEP-AWBH modeling framework. After spatially extracting this information to Tier 2 sites, we applied machine learning (Random Forests) algorithms to this large suite of environmental covariates (N ~ 300) to predict soil C to 1m.

Table 1: Factors representing the STEP-AWBH modeling framework used with Random Forests

Factors	Property	Factors	Property		
Soil	Soil order	Ecology	Monthly LAI		
	Soil suborder		Annual min, max and mean NDVI		
	Soil great group		NDVI green-up & brown-down		
	Soil texture		NDVI greenup and browndown rate		
	Cation exchange capacity		NDVI Season length		
	Bulk density		NDVI amplitude and base NDVI level		
	Organic matter content		Max peak NDVI		
	Site index		Canopy coverage and Imperviousness		
	Phosphorus		Basal area weighted canopy height		
	Nitrogen		Aboveground live dry biomass		
	Potassium	Gross and net primary production			
	Calcium	Woody biomass			
	Zinc	Root biomass			
	Manganese	Basal area			
	Copper	Leaf area index			
	pH	Stem volume over bark			
	Base saturation	Coarse root production			
	Boron	Foliage biomass production			
	Magnesium	Woody biomass production			
	gamma absorbed dose	Total net primary production (NPP)			
gammaPotassium	Above ground NPP				
gammaThorium	Site index				
gammaUranium	Parent Material	Age			
Iron		Predominant lithology			
Topography	Slope	Secondary lithology			
	Aspect	Atmospheric	Precipitation		
	Curvature		Maximum temperature		
	Profile curvature		Minimum temperature		
	Flow direction		Maximum relative humidity		
	Flow accumulation		Minimum relative humidity		
	Latitude / longitude		Solar radiation		
	Topographic wetness index		Water	Available water capacity	
	Ecology			Major Land Resource Area	Drainage class
				Biophysical settings	Hydric rating
Disturbance			Saturated hydraulic conductivity		
Land cover		Soil moisture (seasonal & annual)			
Environmental site potential		Biota	Foliage		
Monthly MODIS FPAR			Vegetation type		
Vegetation height			Vegetation type system group 1		
Vegetation cover			Vegetation type system group 2		
Forest canopy properties			Vegetation type order		
Landsat ETM β bands			Vegetation type class		
Landsat ETM β tasseled cap indices	Vegetation type subclass				
Landsat ETM β principal components	Human	Land use			
Monthly NDVI		Trees per hectare			
Monthly EVI		Stand / Plot age			

Workflow

- 1) Aggregate measured soil C to 1m using a spline function
- 2) Compile data
- 3) Extract data to Tier 2 sites
- 4) Perform data munging / cleaning / aggregating
- 5) Create a single 'master file' containing all relevant data
- 6) Split data into training (70%) and validation (30%) sets
- 7) Random Forest generates an ensemble of *decision trees* for soil C predictions
- 8) Finally, model performance is tested using the validation dataset

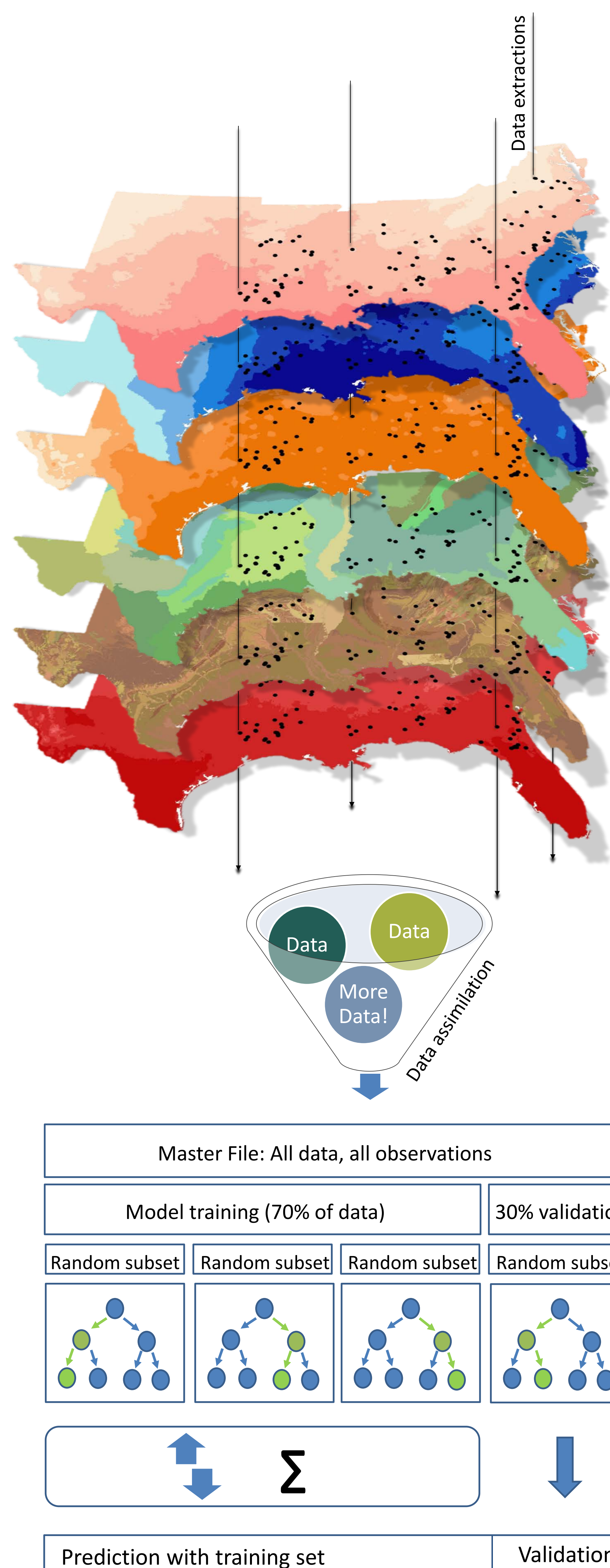


Figure 1: Conceptual model of workflow.

Results

- **Top predictors** - Major Land Resource Areas, geology, biomass, climate data, and soil properties
- **Model performance** - adjusted R² of 0.93 and 0.52 for training and validation sets using all data
- **Model performance** - adjusted R² of 0.90 and 0.36 (training & validation) when using only ancillary data

The best models were achieved using a comprehensive set of STEP-AWBH predictors, which consists of measured soil properties in conjunction with a wide variety of publicly available data sources from various foreign and domestic agencies, which resulted in an adjusted R² of 0.93 and 0.52 for training and validation sets, respectively. The top five predictors using this approach include eco-regions, soil great group, concentration of soil nitrogen, major land resource areas, and age of parent material (geology). Of these predictors, four represent categorical data, while the remaining represents continuous data.

Random Forests also performed well when using only ancillary data (not explicitly measured at Tier 2) as predictors, and resulted in an adjusted R² of 0.90 and 0.36 for training and validation datasets, respectively. The top five predictors using this approach include eco-regions, soil great group, major land resource areas, parent material age, and soil suborder. The top five predictors all represent categorical data.

Refer to Figure 2 and 3 below for the full list of important variables.

Important variable using measured Tier 2 soil data

Predictors	Data source	Importance
Eco-region level 4	Eco-regions	0
Soil great group	gSSURGO	0
Nitrogen concentration	Tier 2	0
Major land resource area	MLRA	0
Parent material age	US geology	0
Organic matter (1m)	Tier 2	0
Cation exchange capacity	Tier 2	0
Magnesium ppm	Tier 2	0
Soil suborder	gSSURGO	0
Calcium ppm	Tier 2	0
Secondary lithology	US geology	0
Bulk density	Tier 2	0
Predominant lithology	US geology	0
Soil order	gSSURGO	0
Soil moisture (April)	SMOS	0
Manganese ppm	Tier 2	0
Trees per hectare	Tier 2	0
Minimum solar radiation	METDATA	0
gammaPotassium	Gamma Ray	0
GammaAbsorbedDose	Gamma Ray	0
Clay content	gSSURGO	0
Drainage class	gSSURGO	0
Potassium ppm	Tier 2	0
Land use / Land cover	NASS	0
Temperature (March-May)	METDATA	0
Disturbance	LANDFIRE	0
Phenology production	AVHRR	0
Foliage	Tier 2	0
Basal area	Tier 2	0
Precipitation (June-August)	METDATA	0

Figure 2. Variables of importance when predicting soil C to 1m depth using measured soil data collected from the Tier 2 network in addition to ancillary data.

Important variables without measured Tier 2 soil data

Predictors	Data source	Importance
Eco-region level 4	Eco-regions	0
Soil great group	gSSURGO	0
Major land resource area	MLRA	0
Age of parent material	US geology	0
Soil suborder	gSSURGO	0
Predominant lithology	US geology	0
Secondary lithology	US geology	0
Topographic curvature	SRTM	0
Land use / land cover	NASS	0
Trees per hectare	Tier 2	0
Sand content	gSSURGO	0
Stem volume over bark	Tier 2	0
Leaf are index	Tier 2	0
Organic matter	gSSURGO	0
Drainage class	gSSURGO	0
Minimum solar radiation	METDATA	0
Soil moisture (April)	SMOS	0
Soil moisture (July)	SMOS	0
Clay content	gSSURGO	0
Disturbance	LANDFIRE	0
Foliage production	Tier 2	0
Site index	Tier 2	0
gammaPotassium	Gamma Ray	0
Phenology	AVHRR	0
Soil moisture (Jan.-March)	SMOS	0
Biomass	NBCD	0
Precipitation (June-August)	METDATA	0
Max. temp. (Dec. Feb.)	METDATA	0
Flow accumulation	SRTM	0
Basal area	Tier 2	0

Figure 3. Variables of importance when predicting soil C to 1m depth using only ancillary data (soil data not explicitly measured at the Tier 2 network).

Concluding Remarks

Although climate is an important factor in determining soil C at global and regional scales, other factors proved more important for predicting soil C at the regional scale. This is explained by 1) scaling considerations and 2) ecosystem properties. In regards to scaling, climate varies much more at global scales compared to smaller, regional scales, therefore reducing its prediction power in this model. Additionally, many of the top predictors, such as MLRAs and eco-regions, are ecological factors, and are not only functions of their biotic components, but also abiotic components such as soil, water, and climate. Therefore, it would be wrong to assume that climate is not an important factor for soil C at regional scales. However, ecosystem properties, which include climate information, proved to be most relevant for predicting soil C.