## Automated global soil mapping: discovering spatial soil patterns using machine learning



T. (Tom) Hengl <<u>tom.hengl@isric.org</u>>

DailyMeteo.org/2014 Conference, Belgrade 23-27 June





"Where meteorologists and climatologists meet R developers..."







ScienceDire	ect Journ	nals Books			Shopping
Search all fields		Author name	This Journal/Book	Volume Issue Page	Advanced search
	Spatial Statist	tics			
SPATIAL STATISTICS	Supports Open Access	About this Journal	Sample Issue Online	Submit your Article	
Annual Annual Market Market	Set new article feed	b			
Manufacture Constant Property State Constant Constant Property Engine Constant of Constant Property	Set new Open Acce	ess article feed			
	Subscribe to new vo	olume alerts			
	Add to Favorites				
and the second second second	Copyright © 2015 Elsevie	er B.V. All rights reserved			











## Software outsourcing to Eastern Europe: Which countries work best?

Eastern Europe is building up a head of steam as a software outsourcing powerhouse, as companies focus more on cost-effectiveness and proximity, rather than just low cost.

1

By Andrada Fiscutean for Central European Processing | April 27, 2016 -- 11:13 GMT (04:13 PDT) | Topic: Developer

Recommended Content:

0

f 177

#### White Papers: How to prevent security breaches in your retail network

This paper describes the multi-vector nature common to most high-profile security breaches in retail enterprises. IT directors and security officers will learn about several of the characteristics common to recent breaches..

♦ Find out more



in 558



finds



## SOIL GRIDS A system for automated global soil mapping

www.soilgrids.org





## Automated (predictive) global soil mapping

## We aims at implementing (all) knowledge of soil science through robust algorithms (BUP), and putting the system in operational use...

# so that you do not need any soil mappers any more :)





#### SoilGrids1km — Global Soil Information Based on Automated Mapping

Tomislav Hengl , Jorge Mendes de Jesus, Robert A. MacMillan, Niels H. Batjes, Gerard B. M. Heuvelink, Eloi Ribeiro, Alessandro Samuel-Rosa, Bas Kempen, Johan G. B. Leenaars, Markus G. Walsh, Maria Ruiperez Gonzalez

Published: August 29, 2014 • DOI: 10.1371/journal.pone.0105992

Article	About the Authors	Metrics	Comments	Related Content	Download PDF		-
*					Print	Shar	re
<ul> <li>Abstract Introduction</li> </ul>	Abstract				🔘 CrossMark		
Materials and Methods	Backgroun	d			Subject Area	as	?
Results	5						1
Discussion	Soils are widely	recognized as a non-renev	wable natural resource and	l as biophysical carbon	Agricultural	oil scie	-
Acknowledgments	sinks. As such, t	here is a growing requiren	nent for global soil informat	ion. Although several	Clay mineral	ogy	1
Author Contributions	limited spatial de	etail.	st, these tend to suffer from	n inconsistencies and	Deserts		
References	Methodolog	gy/Principal Findi	ngs		Forecasting		P
					Remote sens	ing	
Reader Comments (0)	We present Soil	Grids1km — a global 3D s	oil information system at 1	km resolution —	Surface temp	erature	in-
Figures	organic carbon (	g kg-1), soil pH, sand, silt	and clay fractions (%), bul	k density (kg m-3), cation-			-
	exchange capac	ity (cmol+/kg), coarse frag	ments (%), soil organic ca	rbon stock (t ha-1), depth	Taxonomy		12
	to bedrock (cm), predictions are b	World Reference Base so	oil groups, and USDA Soil T ediction models which we fit	Faxonomy suborders. Our	Urban areas		

## SoilGrids250m: global gridded soil information based on Machine Learning

Tomislav Hengl<sup>1</sup>, Jorge Mendes de Jesus<sup>1</sup>, Gerard B.M. Heuvelink<sup>1</sup>, Maria Ruiperez Gonzalez<sup>1</sup>, Milan Kilibarda<sup>2</sup>, Aleksandar Blagotić<sup>3</sup>, Wei Shangguan<sup>4</sup>, Marvin N. Wright<sup>5</sup>, Xiaoyuan Geng<sup>6</sup>, Bernhard Bauer-Marschallinger<sup>7</sup>, Mario Antonio Guevara<sup>8</sup>, Rodrigo Vargas<sup>8</sup>, Robert A. MacMillan<sup>9</sup>, Niels H. Batjes<sup>1</sup>, Johan G.B. Leenaars<sup>1</sup>, Eloi Ribeiro<sup>1</sup>, Ichsani Wheeler<sup>10</sup>, Stephan Mantel<sup>1</sup>, and Bas Kempen<sup>1</sup>

<sup>1</sup>ISRIC — World Soil Information, Wageningen, the Netherlands
<sup>2</sup>Faculty of Civil Engineering, University of Belgrade, Serbia
<sup>3</sup>GILab Ltd, Belgrade, Serbia
<sup>4</sup>College of Global Change and Earth System Science, Beijing Normal University, Beijing, China
<sup>5</sup>Institut für Medizinische Biometrie und Statistik, Lübeck, Germany
<sup>6</sup>Agriculture and Agri-Food Canada, Ottawa, Ontario, Canada
<sup>7</sup>Department of Geodesy and Geoinformation, Vienna University of Technology, Vienna, Austria
<sup>8</sup>University of Delaware, Newark DE, USA
<sup>9</sup>LandMapper Environmental Solutions Inc., Edmonton, Canada
<sup>10</sup>Envirometrix Inc., Wageningen, the Netherlands

Correspondence to: T. Hengl (tom.hengl@isric.org)

5

Abstract. This paper describes the technical development and accuracy assessment of the most recent and improved version of the SoilGrids system at 250 m resolution (June 2016 update). SoilGrids provides global predictions for standard numeric soil properties (organic carbon, bulk density, Cation Exchange Capacity (CEC), pH, soil texture fractions and coarse fragments) at seven standard depths (0, 5, 15, 30, 60, 100 and 200 cm), in addition to predictions of depth to bedrock and distribution of soil classes based on the World Reference Base (WRB) and USDA classification systems (ca. 280 raster layers in total). Predictions were based on ca. 150,000 soil profiles used for training and a stack of 158 remote sensing-based soil covariates (primarily derived from MODIS land products, SRTM DEM derivatives, climatic images and global landform and lithology maps), which were used to fit an ensemble of machine learning methods — random forest and gradient boosting and/or multinomial logistic regression — as implemented in the R packages ranger, xgboost, nnet and caret. The results of 10–fold



Near-global freshwater-specific

## **Global Soil Information Facilities**

#### SoilGrids



- Soil property-depth curves (with confidence intervals)
- Dominant, 2<sup>nd</sup> and 3<sup>rd</sup> soil type
- Nutrient status (e.g. soil organic carbon stock)
- Available water capacity
- Current and past land use
- Closest soil data provider / agricultural extension service



## SoilGrids inputs

- → "World's largest" compilation of soil profile / soil sample data sets (national and international datasets from over 45 countries) → these contain both soil classification data + soil analytical/laboratory data.
- → A 40TB repository of MODIS land products, climatic images, DEM derivatives, geological and geomorphological data (all at 250 m resolution)
- → ISRIC's international network that can cross-check and validate spatial prediction patterns / values.



## ca 150,000 points shown on this map





>850,000 measurements of soil organic carbon

# Machine learning





IS

## Some common misconceptions

- → MLA is a black box it is not!
- → Most of new implementations of random forest and similar will penalize for overfitting.
- → MLA is non-statistical / there are no probabilities and the output (prediction) uncertainty is unknown — it is not true. "MLA is a marriage between stats and knowledge representation".





USRIC World Soil Information

## Watch out from overfitting







**Figure 1:** The black curve shows how  $\operatorname{Err}_t$  varies as a function of t for a particular realization of  $Q_1, \ldots, Q_t$ . The red curve was obtained by generating many ensembles, and computing the average of the respective  $\operatorname{Err}_t$  values at each t. Similarly, the blue curves depict the 10th and 90th percentiles of  $\operatorname{Err}_t$  at each t.



Based on Lopez (2015)

## **Statistical Modeling: The Two Cultures**

#### Leo Breiman

Abstract. There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

#### 1. INTRODUCTION

Statistics starts with data. Think of the data as being generated by a black box in which a vector of input variables  $\mathbf{x}$  (independent variables) go in one side, and on the other side the response variables  $\mathbf{y}$ come out. Inside the black box, nature functions to associate the predictor variables with the response variables, so the picture is like this:



The values of the parameters are estimated from the data and the model then used for information and/or prediction. Thus the black box is filled in like this:



*Model validation.* Yes-no using goodness-of-fit tests and residual examination.

Estimated culture population. 98% of all statisti-



"The tagline is, *Do you want to be* a machine learning ninja?" says Christine Robson, a product manager for Google's internal machine learning efforts, who helps administer the program. "So we invite folks from around Google to come and spend six months embedded with the machine learning team, sitting right next to a mentor, working on machine learning for six months, doing some project,



getting it launched and learning a lot."

D B



# MLAs of interest to global soil mapping



## Some logical choices

- → Soils vary locally (there is a strong local component)
   tree-based models probably more suited.
- → Soils are shaped as a result of complex processes running over large periods of time... but there are also sudden disturbances (soil slides, erosion)... in short: soil genesis is complex.
- → Remote sensing (VISNIR) can not be used to directly map soils, BUT it can be used to represent soil forming factors.





- → An extensive stack of carefully selected covariates (158)
- → Random forest + Gradient boosting
- → Put more effort to optimize the system so that it can be updated relatively fast





#### UNIVERSITÄT ZU LÜBECK INSTITUT FÜR MEDIZINISCHE BIOMETRIE UND STATISTIK



#### Home

#### Main Menu

r	4	r	T	
		Ł	1	

- Staf
- How to find us
- Weblinks

#### Research

- Awards
- Projects
- Publication
- Software
- Employment Ad

#### Courses

- **Biometrical Consulting Service**
- Lectures
- Leselunci
- Seminar
- Student Research Project
- Veranstaltungen

#### Search

Search this site

Search

#### Marvin N. Wright, M.Sc.

Institut für Medizinische Biometrie und Statistik Universität zu Lübeck Universitätsklinikum Schleswig-Holstein, Campus Lübeck Ratzeburger Allee 160 23562 Lübeck fon: 0451-500 2782 fax: 0451-500 2999 wright@imbs.uni-luebeck.de



#### Publications

- Wright, M. N., Ziegler, A. & König, I. R. (2016). Do little interactions get lost in dark random forests? BMC Bioinformatics 17:145. http://dx.doi.org/10.1186/s12859-016-0995-8.
- Wright, M. N. & Ziegler, A. (2016). ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R. Journal of Statistical Software, in press. http://arxiv.org/abs/1508.04409.
- Schirmer, J. H., <u>Wright, M. N.</u>, Vonthein, R., Herrmann, K., Nölle, B., Both, M., Henes, F., Arlt, A., Gross, W. L., Schinke, S., Reinhold-Keller, E., Moosig, F. & Holle, J. U. (2016). Clinical presentation and long-term outcome of 144 patients with microscopic polyangiitis in a monocentric German cohort. *Rheumatology (Oxford)* 55:71-79. http://dx.doi.org/10.1093/rheumatology/kev286.
- Wright, M. N. & Ziegler, A. (2015). Multiple censored data in dentistry: A new statistical model for analyzing lesion size in randomized controlled trials. *Biometrical Journal* 57:384–394. http://dx.doi.org/10.1002/bimj.201400118.
- Paulick, C., Wright, M. N., Verleger, R. & Keller, K. (2014). Decomposition of 3-way arrays: A comparison of different PARAFAC algorithms. Chemometrics and Intelligent Laboratory Systems 137:97-109. http://dx.doi.org/10.1016/j.chemolab.2014.06.009.

#### Technical Reports

• Schmid, M., Wright, M. N. & Ziegler, A. (2015). On the Use of Harrell's C for Node Splitting in Random Survival Forests. Technical Report. http://arxiv.org/abs/1507.03092.



## >850,000 measurements of soil organic carbon

#### ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R

Marvin N. Wright Universität zu Lübeck Andreas Ziegler Universität zu Lübeck, University of KwaZulu-Natal

#### Abstract

We introduce the C++ application and R package **ranger**. The software is a fast implementation of random forests for high dimensional data. Ensembles of classification, regression and survival trees are supported. We describe the implementation, provide examples, validate the package with a reference implementation, and compare runtime and memory usage with other implementations. The new software proves to scale best with the number of features, samples, trees, and features tried for splitting. Finally, we show that **ranger** is the fastest and most memory efficient implementation of random forests to analyze data on the scale of a genome-wide association study.

Keywords: C++, classification, machine learning, R, random forests, Rcpp, recursive partitioning, survival analysis.

otto group	S.	by Tiand using da	i Chen · last ru ta from Otto G	n 11 months roup Produc	DOST   s ago · R r ct Classifi	VIODEI ON otebook · 4410 cation Challeng	Otto Da <sup>7 views</sup>	· ,	10 vote
Rep <mark>ort Code</mark>	Output (1)	Comments (5)	Log Ver	sions (10)	Forks (70)			Fork	Scri
D									
кероrt	Unde	rstand	ing XC Da	Boo atase	st M t	odel o	n Otto		
кероrt	Unde	rstand	ing XC Da Mic	GBOO atase	stM t	odel o	n Otto		
кероrt	<ul> <li>1 Introductio</li> <li>2 Preparation</li> <li>3 Model trair</li> <li>4 Model und</li> </ul>	n n of the data ing erstanding	ling XC Da Mid	GBOO atase chaël Benest	st M t	odel o	n Otto		
кероrt	<ul> <li>1 Introductio</li> <li>2 Preparation</li> <li>3 Model train</li> <li>4 Model und</li> <li>4.1 Featur</li> <li>4.2 Interpr</li> <li>4.3 Tree g</li> </ul>	n n of the data ing erstanding e importance etation raph	ling XC Da Mic	GBOO atase chaël Benest	st M t	odel o	n Otto		

of the model?

\_

#### XGBoost: A Scalable Tree Boosting System

Tianqi Chen University of Washington tqchen@cs.washington.edu Carlos Guestrin University of Washington guestrin@cs.washington.edu

#### ABSTRACT

Tree boosting is a highly effective and widely used machine learning method. In this paper, we describe a scalable endto-end tree boosting system called XGBoost, which is used widely by data scientists to achieve state-of-the-art results on many machine learning challenges. We propose a novel sparsity-aware algorithm for sparse data and weighted quantile sketch for approximate tree learning. More importantly, we provide insights on cache access patterns, data compression and sharding to build a scalable tree boosting system. By combining these insights, XGBoost scales beyond billions of examples using far fewer resources than existing systems.

#### **CCS Concepts**

•Methodologies  $\rightarrow$  Machine learning; •Information systems  $\rightarrow$  Data mining;

#### Keywords

ISRIC

many applications. Tree boosting has been shown to give state-of-the-art results on many standard classification benchmarks [14]. LambdaMART [4], a variant of tree boosting for ranking, achieves state-of-the-art result for ranking problems. Besides being used as a stand-alone predictor, it is also incorporated into real-world production pipelines for ad click through rate prediction [13]. Finally, it is the de-facto choice of ensemble method and is used in challenges such as the Netflix prize [2].

In this paper, we describe XGBoost, a scalable machine learning system for tree boosting. The system is available as an open source package<sup>2</sup>. The impact of the system has been widely recognized in a number of machine learning and data mining challenges. Take the challenges hosted by the machine learning competition site Kaggle for example. Among the 29 challenge winning solutions <sup>3</sup> published at Kaggle's blog during 2015, 17 solutions used XGBoost. Among these solutions, eight solely used XGBoost to train the model, while most others combined XGBoost with neural nets in en-

**World Soil Information** 



## Journal of Statistical Software

November 2008, Volume 28, Issue 5.

http://www.jstatsoft.org/

#### Building Predictive Models in R Using the caret Package

Max Kuhn Pfizer Global R&D

#### Abstract

The caret package, short for classification and regression training, contains numerous tools for developing predictive models using the rich set of models available in R. The package focuses on simplifying model training and tuning across a wide variety of modeling techniques. It also includes methods for pre-processing training data, calculating variable importance, and model visualizations. An example from computational chemistry is used to illustrate the functionality on a real data set and to benchmark the benefits of parallel processing with several types of models.

Keywords: model building, tuning parameters, parallel processing, R, NetWorkSpaces.

## Results



## They would have been interested in this...



Vasili Dokuchaev

**The Russian School** 

Soil forming factors Soil forming processes Different Soils FACTORS OF SOIL FORMATION (1941)

A System of Quantitative Pedology Hans Jenny









Figure 6. Examples of fitted relationships for bulk density (above), pH (middle) and soil organic carbon (below). Plots show target variables and top three most important covariates as reported by the random forest model. DEPTH.f is the depth from soil surface, T09M0D3 is mean monthly temperature for September, TMDM0D3 is mean annual temperature, PRSMRG3 is total annual precipitation, M04M0D4 is mean monthly MODIS NIR band reflectance, P07MRG3 is mean monthly precipitation for July, T01M0D3 is mean monthly temperature for January, and



Figure 5. Fitted variable importance plots for target variables. Generated as an average between using the ranger and ×gboost packages, (for soil types results are based on the ranger model only). DEPTH.f is the depth from soil surface, T\*\*M0D3 and N\*\*M0D3 are mean monthly temperatures daytime and nighttime (red color), TWI, DEM, VBF and VDP are DEM-parameters (bisque color), M\*\*M0D4 are mean

## Variable importance soil types

DEMMRG5 VW4MOD1 VW3MOD1 VBFMRG5 TWIMRG5 VDPMRG5 VW5MOD1 ASSDAC3 T08MSD3 N07MSD3 N06MOD3 M11MOD4 N07MOD3 T09MSD3

WRB classes



**USDA classes** 





**Figure 2.** Example of soil variable-depth curves: original sampled soil profiles vs predicted values (SoilGrids) at seven standard depths (broken red line) and estimated soil organic carbon stock for depths 0–100 and 100–200 cm. Locations of points: mineral soil S1991CA055001 (-122.37°W, 38.25°N), and an organic soil profile S2012CA067002 (-121.62°W, 38.13°N).



Bulk density (FE) in kg / m3 (CV R-squared: 0.76) Counts Counts Counts predicted (SoilGrids250m) predicted (SoilGrids250m) measured measured Soil pH x 10 in KCI (CV R-squared: 0.77) CEC soil in cmolc/kg (CV R-squared: 0.64) Counts Counts Counts predicted (SoilGrids250m) predicted (SoilGrids250m) 

SOC in g/kg (CV R-squared: 0.64)

predicted (SoilGrids250m) measured

Soil pH x 10 in H2O (CV R-squared: 0.83)



Sand fraction in % (CV R-squared: 0.79)



Silt fraction in % (CV R-squared: 0.79)

measured



Clay fraction in % (CV R-squared: 0.73)







### Xeralfs















the state local

# Conclusions



## Conclusions

- → Traditional soil surveyors got it right! distribution of soil classes is mainly controlled by DEM morphometry (especially hydrological parameters).
- → Soil classification and polygon models of soils seem to make sense — in many parts of the world we see "soil groupings i.e. soil bodies"... but there are also transition zones and small individual patches... so it is really a hybrid that we need.
- → In the machine learning framework, much more time needs to be spent on preparing data

World Soil Information





BBC	Sign in	News	Sport	Weather	Shop	Earth	Travel	More -	Search	Q	L
NEW	′S										
Home Video	) World UK Busin	iess   Teo	h Scien	ce Magaz	zine Er	itertainmer	nt & Arts 🛛	Health	In Pictures	More 🕶	

Magazine

## Viewpoint: How creativity is helped by failure

() 14 November 2015 | Magazine



When it comes to creating a great work of art, practice makes perfect, writes Matthew Syed.

In today's Magazine

#### **Top Stories**

First Paris attacker identified

C 22 minutes ago

Democrats debate IS after Paris attacks

() 5 hours ago

Phones 'need bed mode' to help sleep

() 15 November 2015

#### Features & Analysis



Eau de Belgrade How a perfumer survived two wars and communism



