

Sentinel 2 and Machine Learning for biomass modelling

Davor Korman¹, Alen Berta¹, Zrinka Mesić¹, Ana Ostojić¹, Nela Jantol¹, Ivona Žiža¹, Vladimir Kušan¹

¹ Oikon Ltd Institute of Applied Ecology

INTRODUCTION

Remote sensing techniques have been used in forestry for last several decades, but due to the variety of available imagery data and data processing algorithms, the single best method for aboveground biomass (AGB) assessment is not existent. Furthermore, there is no quorum whether is it more appropriate to create AGB assessment models with data averaged on seasonal or monthly temporal scale in order to create the most accurate model. Besides, classical statistical methods often result in models with moderate or poor performance, while on the other hand Machine Learning techniques, coupled with the advancements in computer hardware and an increase in processing power, present new opportunities in remote sensing and AGB assessment modelling.

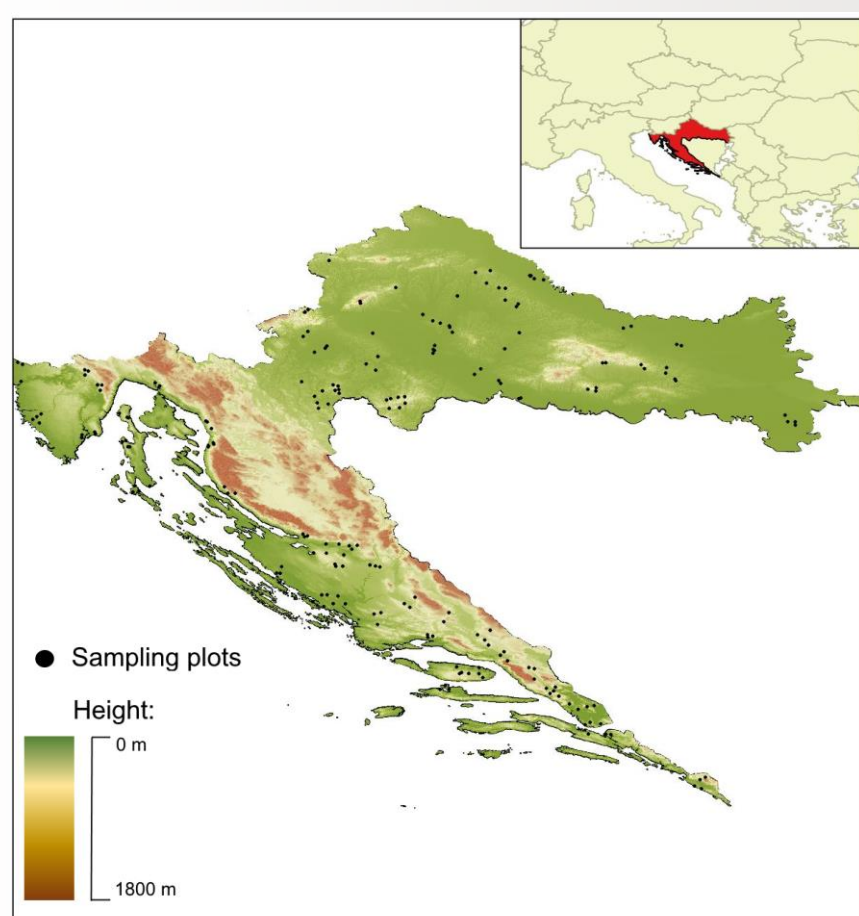
RESEARCH OBJECTIVE

In this research we have tested the effectiveness of using Sentinel 2 imagery for AGB assessment and designed a practical approach (algorithm) for extracting information from large amounts of satellite data.

Our goals were to:

- 1) create models for aboveground biomass assessment for thicket, maquis and continental forests of first age class using Sentinel 2 data
- 2) investigate and compare the suitability of monthly and seasonally averaged data for AGB assessment
- 3) explore the possibility of improving model accuracy by implementing various machine learning techniques.

STUDY AREA



METHODS – GROUND TRUTH DATA

Field data was gathered on 204 sample plots covering six different forest stands across Croatia:

Management class	Number of sample plots for field measurement AGB assessment
Common beech	24
Sessile oak	24
Common oak	20
Narrow-leaved ash	16
Continent total	84
Thicket	76
Maquis	44
TOTAL	204

Gathered data was analysed in several steps for the purpose of obtaining total AGB per sample plot:

- I. Measurement of tree height and diameter of each tree inside the sample plot
- II. Calculation of individual tree volumes using Schumacher-Hall volume equation

$$V = a \cdot db \cdot hc \cdot f$$
- III. Calculation of total volume per sample plot

Response variable for image based modelling

METHODS – REMOTE SENSING DATA

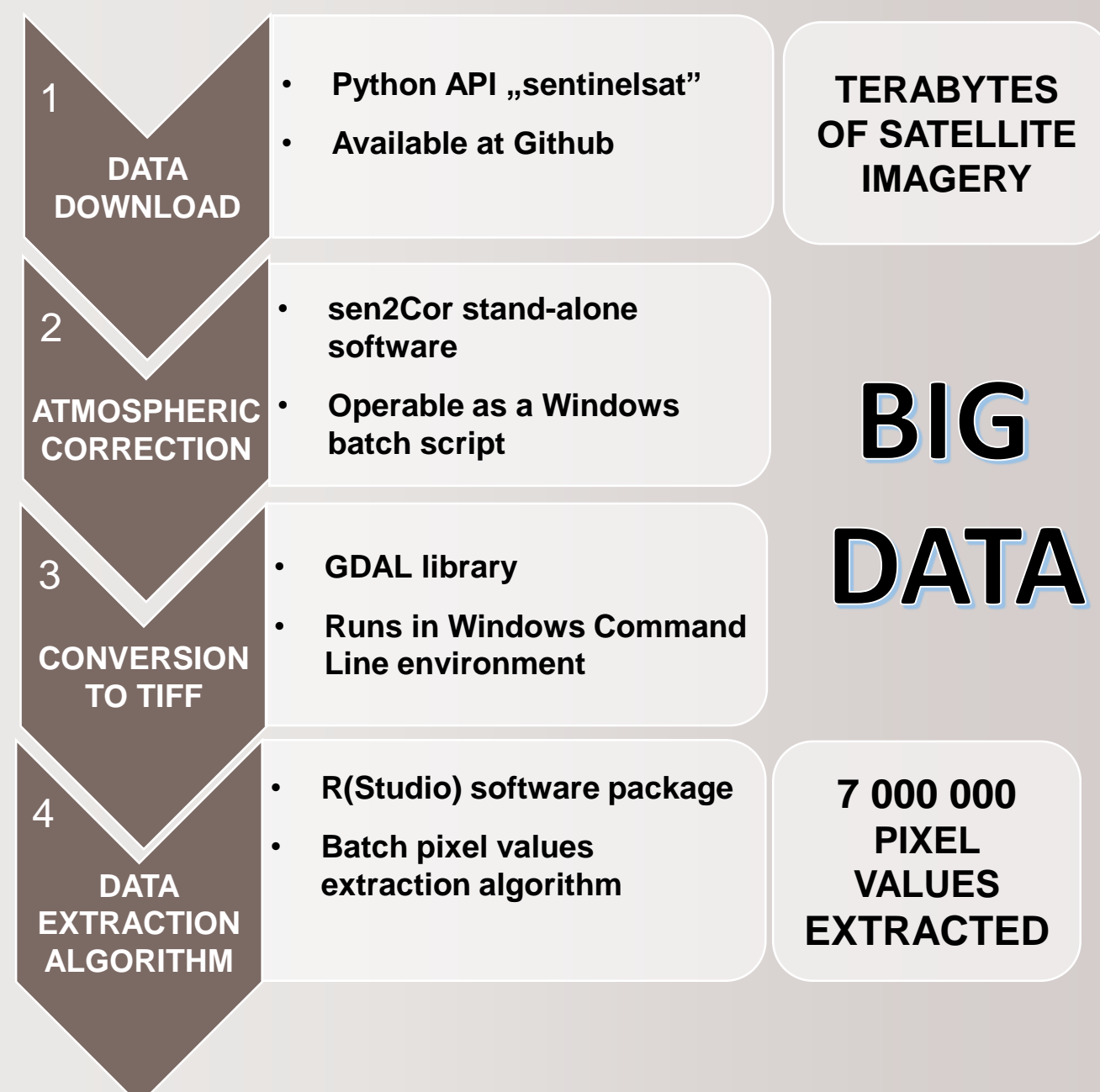
For the purpose of the research and for the comparison of seasonal and monthly based AGB models, all available Sentinel 2 imagery in 2016 growing season was downloaded

(9 Sentinel 2 bands + 110 vegetation indices) * month | season

1100 independent variables

METHODS – DATA DOWNLOAD AND PROCESSING

DATA PROCESSING STEPS



TERABYTES OF SATELLITE IMAGERY

BIG DATA

7 000 000 PIXEL VALUES EXTRACTED

OPEN SOURCE SOFTWARE USED ONLY



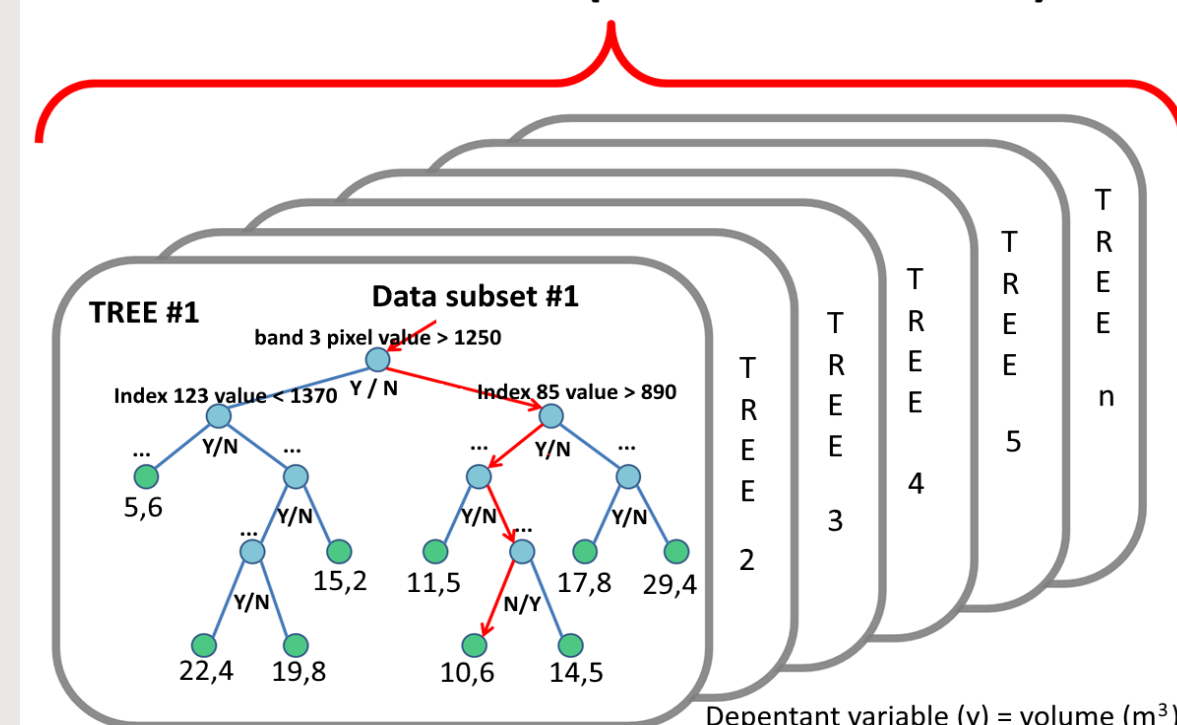
Code available at <https://github.com/DKorman>

METHODS – PRIMARY MODEL BUILDING PIPELINE

- I. Random forest for selecting best among ~1100 variables
- II. Stepwise regression for additional variable reduction
- III. Creation of multiple regression model:

$$\text{Volume} = a + b_1 \cdot \text{VegetationIndex1} + b_2 \cdot \text{VI2} \dots$$

Random forest (of decision trees)



High interpretability, but... (see results)

PRIMARY MODEL RESULTS & CONCLUSION

Seasonal vs. monthly models

Monthly averaged data has been selected as the most appropriate independent variable in all cases

Model accuracy

Forest stands	R ² adj.	Percent relative standard error
Common beech and sessile oak	0.65	41.75
Common oak and narrow-leaved ash	0.80	30.43
Thicket in Istria and Kvarner	0.64	33.12
Thicket in Dalmatia	0.55	36.87
Maquis	0.53	61.99

Sentinel 2 based AGB assessment resulted in models with moderate or even poor performance

IMPROVING ACCURACY WITH MACHINE LEARNING

Machine learning technique	Hyperparameters tuned	Role in the model creation pipelines
Bootstrap aggregating	Number of trees, tree depth, number of observation in nodes	Model creation and variable selection
Random forest	- -, number of features selected	Model creation and variable selection
Gradient boosting	Shrinkage, depth, number of observations in node, number of trees	Model creation and variable selection
Extreme gradient boosting	Depth, eta, number of trees, lambda, alpha	Model creation and variable selection
Genetic algorithm	Percentage of mutation, elitism, crossover	Model creation and variable selection
Support vector machines	Epsilon, cost, gama	Model creation
Princ. component analysis		Dimensionality reduction

RESULTS OF MACHINE LEARNING IMPLEMENTATION

Some of the pipelines tested (example on 2 forest stands)

ML PIPELINE USED	RMSE		ML PIPELINE USED	RMSE	
	sessile oak	thicket		sessile oak	thicket
Random forest - Linear model (PRIMARY MODEL)	22,04	5,25	Gradient boosting - LR	17,80	5,87
Support vector machines	21,32	6,013	Extreme gradient boosting (tree) - LR	18,73	5,87
Random forest - Random forest	11,61	3,78	Extreme g. boosting (linear) - LR	16,33	6,29
PCA - Random forest	18,32	5,73	Random forest - Extreme g. boosting	10,28	3,80
Genetic algorithm - LR	19,08	5,64	Random forest - Support vector machines	13,19	6,24
PCA -> LR	17,83	5,85	Random forest - gradient boosting	10,59	3,78

MACHINE LEARNING CONCLUSION

MACHINE LEARNING CAN BE AN EXCELLENT TOOL FOR BIOMASS MODELLING BUT

