

title

Patrick Schratz^a, Jannes Muenchow^a, Eugenia Iturritxa¹, Alexander Brenning^a

^a*Department of Geography, GIScience group, Grietgasse 6, 07743, Jena, Germany*

Abstract

Keywords: hyperspectral imagery, forest health, machine-learning, variable importance, model comparison

1. Introduction

Data retrieved from remote sensing satellites is successfully used in forestry to monitor temporal changes across large areas (Martinez del Castillo et al., 2015; Sexton et al., 2015). The use of Synthetic Aperture Radar (SAR) techniques enables scientists to estimate Above-Ground Biomass (AGB) (Lu et al., 2016; Sinha et al., 2015). Forest health is commonly assessed using optical data from multi-/hyperspectral satellites by applying temporal change detections (Zhang et al., 2016). With the recent success story of machine-learning methods in the field of remote sensing, modeling techniques such as Random Forest (RF) are frequently used to model relationships of possible triggers to forest health (Belgiu & Drăguț, 2016; Lary et al., 2016; Michez et al., 2016).

With a robust model, predictions of the modelled response to large areas is possible, giving valuable information about the condition of this variable in unknown regions. To model forest health, usually few variables are extracted based on the spectral signatures of affected and unaffected trees (Lelong et al., 2010). However, spectral (vegetation-)indices have shown the potential to con-

*Corresponding author

Email address: `patrick.schratz@uni-jena.de` (Patrick Schratz)

tribute valuable information to increase predictive accuracy of forest pathogens (Jiang et al., 2014; Adamczyk & Osberger, 2015).

However, the amount of possible (vegetation-)indices that could be calculated is often limited due to a low spectral resolution of freely available data from optical multispectral sensors (e.g. Sentinel-2). Also, there is currently no free data available from hyperspectral sensors that could be used for such studies (after the decommission of the EO-1 Hyperion satellite). If the spatial resolution of the data is too coarse (e.g. $> 5m$), the value of a pixel usually contains information from multiple trees and possibly even bare-ground information. This makes the resulting information almost useless to be used for forest health monitoring on a tree level.

In this study we will use hyperspectral data with a spatial resolution of one meter and 126 spectral bands to model the health status of Monterey Pine (*Pinus radiata*) plantations in northern Spain. The trees in the study area suffer from infections of invasive pathogens such as *Diplodia sapinea*, *Fusarium circinatum*, *Armillaria mellea* or *Heterobasidion annosum* leading to a spread of cankers or defoliation before the tree dies (Mesanza et al., 2016; Iturrutxa et al., 2017). In-situ measurements of defoliation on a tree level are used as a proxy to model tree health. The fungus are assumed to infect the trees through open wounds, possibly caused by previous hail damage (Iturrutxa et al., 2014). The dieback of these trees, which are mainly used as timber, causes high economic damages (Ganley et al., 2009). Hyperspectral remote sensing data in combination with state-of-the-art machine-learning techniques is used to help monitoring the health status in this region by early detecting affected trees/plots.

To extract the most information from the available remote sensing data, we not only calculated the most common vegetation indices like *NDVI* to link against defoliation but all possible ones within the spectral region of the data (400 nm - 1000 nm) that were implemented in the *hsdar* package in R (Lehnert et al., 2018). Additionally, all possible combinations of Normalized Ratio Indices (NRI) were calculated from the data and supplied to a selection of machine-learning algorithms as predictors.

Specifically the following objectives are addressed:

- Comparison of multiple algorithms on their performance to model defoliation of *Pinus radiata* trees using highly-correlated indices
- Exploration of the most important indices of the best performing model
- Prediction of defoliation to *Pinus radiata* plots with an unknown defoliation level

2. Data and study area

2.1. In-situ data

The *Pinus radiata* plots of this study, named *Laukiz 1*, *Laukiz 2*, *Luiando* and *Oiartzun*, are located in the northern part of the Basque Country (Figure 1). *Oiartzun* has the most observations ($n = 529$) while *Laukiz 2* has the largest area size (1.44 ha). All plots besides *Luiando* are located nearby the coast (Figure 1). The data collection took place in September 2016.

2.2. Hyperspectral data

The airborne hyperspectral data was acquired during two flight campaigns on September 28th and October 5th 2016, both around 12 am. The images were taken by an AISAEAGLE-II sensor from the Institut Cartografic i Geologic de Catalunya (ICGC). All preprocessing steps (geometric, radiometric, atmospheric) have been conducted by ICGC. The first four bands were corrupted, leaving 122 bands with valid information.

Additional data characteristics are provided in Table 1:

3. Methods

For all analysis steps we used the open-source statistical programming language R (R Core Team, 2017). The algorithm implementations of the following packages have been used: *xgboost* (Chen & Guestrin, 2016), *kernelab* (Karatoglou et al., 2004) (Support Vector Machine), Vapnik (1998)) and *glmnet*

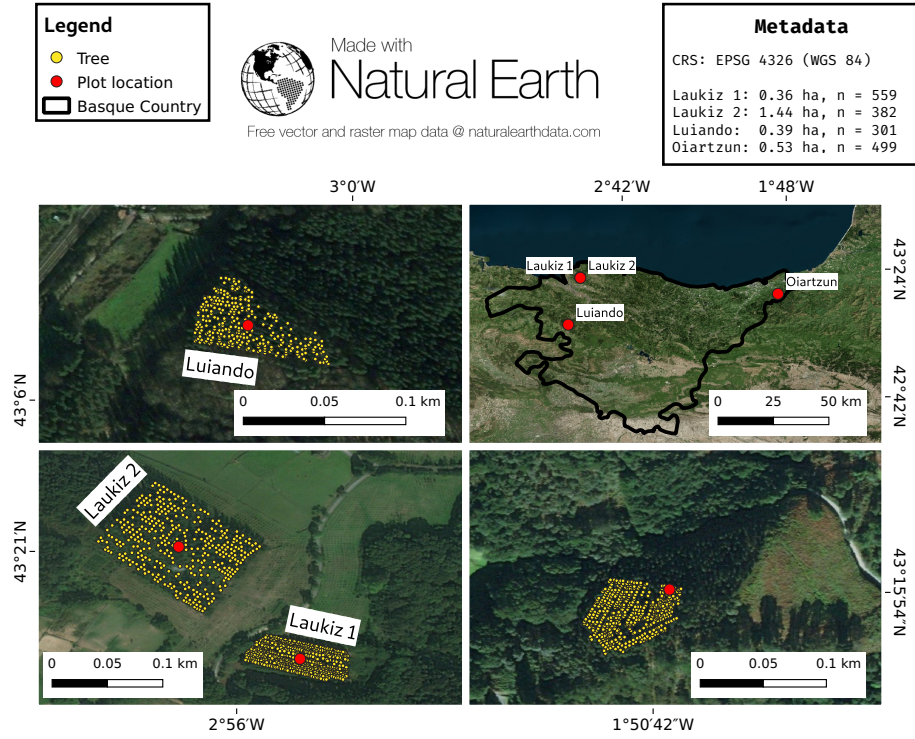


Figure 1: Information about the plot locations, the area of hyperspectral coverage and the number of trees per plot.

(Friedman et al., 2010) (Ridge Regression). We used the R package *mlr* for
75 all modeling related steps. It provides a standardized interface for a wide variety of statistical and machine-learning models in R simplifying essential modeling tasks such as hyperparameter tuning, model performance evaluation and

Table 1: Specifications of hyperspectral data.

Characteristic	Value
Geometric resolution	1 m
Radiometric resolution	12 bit
Spectral resolution	126 bands (404.08 nm - 996.31 nm)
Correction:	Radiometric, geometric, atmospheric

parallelization (Bischi et al., 2016).

3.1. Derivation of indices

80 All vegetation indices (90 total) suitable for the wavelength range of the hyperspectral data that were available in the R package *hsdar* have been calculated. Additionally, all possible NRI were calculated from the data using the formula:

$$NRI_{i,j} = \frac{b_i - b_j}{b_i + b_j} \quad (1)$$

where i and j are the respective band numbers.

85 To account for geometric offsets, we used a buffer of two meters around the centroid of the respective tree. The mean value of all pixels touched by the buffer was assigned as the final value for each index. Missing values were removed from the mean value calculation. In total, 7875 Normalized Ratio Indices NRI have been calculated ($\frac{125*126}{2}$). Due to four corrupted bands and some other
90 numerical problems, few indices returned NA values for some observations. These indices were removed from the dataset, leaving a total of 7471 variables without missing values.

3.2. Exploratory analysis of plot characteristics

Plot characteristics like age, stand density and defoliation were analysed to
95 show differences among the plots. Additionally, the spectral signatures of each plot have been visualized.

3.3. Benchmarking of algorithms

Multiple algorithms were benchmarked on predictive performance to find the best performing one. Besides the well-known Support Vector Machines (SVM)
100 (Vapnik, 1998) we also used *xgboost* which is ensemble method relying on the idea of tree boosting that gained a lot of attention in recent years (Chen & Guestrin, 2016). We also added penalized L2 (Ridge) regression to the algorithm collection due to its ability to handle highly correlated covariates. We also want

to explain why we did not consider the probably most popular machine-learning
 105 algorithm, Random Forest: Due to the high number of variables model fitting
 times in the range hours for a single model fit were not practicable for this work.
 These high fitting times are caused by hyperparameter `mtry` which scales with
 the number of variables (Probst et al., 2018).

3.3.1. Performance estimation

110 The algorithms were benchmarked in two ways: (1) Using spatial cross-
 validation (CV) for each plot using on the k-means clustering approach of Bren-
 ning (2012). To reduce runtime we used a five-fold five-times repeated CV
 setup. (2) Using spatial CV on the plot level with each plot being the test set
 once. This results in four performance estimates, one for each fold. For (1) we
 115 only used the best performing algorithm from (2). The reason why the (2) was
 chosen for algorithm selection is that this model will also be used to spatially
 predict defoliation in other plots.

3.3.2. Hyperparameter tuning

To tune the hyperparameters of the algorithms, we used Sequential-based
 120 Model Optimization (SMBO) via the R package *mlrMBO* (Bischl et al., 2017).
 This Bayesian approach first composes n randomly chosen hyperparameter set-
 tings out of a user defined search space. After these n tries have been evaluated,
 a new hyperparameter setting to be evaluated next is proposed based on the
 setting that performed best. This strategy continues until a termination crite-
 125 rion, defined by the user, is reached (Hutter et al., 2011; Jones et al., 1998). In
 this work we used an initial design of 30 randomly composed hyperparameter
 settings and a termination criterion of 20 iterations, resulting a total budget of
 50 evaluated hyperparameter settings per fold. The advantage of this tuning
 approach is that it substantially reduces the tuning budget which is needed to
 130 find a setting close to the global minimum compared to methods that do not use
 information from previous runs such as *random search* or *grid search* (Bergstra
 & Bengio, 2012).

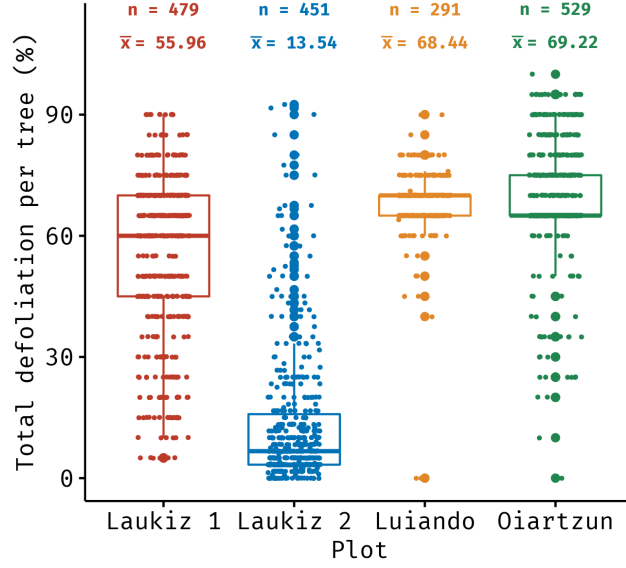


Figure 2: Descriptive statistics of the response variable *defoliation*.

3.4. Variable importance

To find indices that contributed most to model performance, we used permutation-based variable importance on the best performing algorithm.

4. Results

4.1. Exploratory data analysis

Oiartzun shows the highest defoliation ($\bar{x} = 69.22\%$) among the plots while *Laukiz 2* is the healthiest ($\bar{x} = 13.54\%$) (Figure 2). All plots besides *Luiando* show an evenly distributed level of defoliation across the entire plot.

The high degree of defoliation of *Luiando* and *Oiartzun* is also visible in the spectral signatures of the plots (Figure 3). Both plots show lower mean reflectance values around the wavelength range 800 nm - 1000 nm compared to *Laukiz 1* and *Laukiz 2*. *Oiartzun* is almost completely missing the reflectance drop at around 815 nm that is visible for all other plots but instead shows a higher magnitude for the reflectance increase at around 920 nm. *Laukiz 2* shows

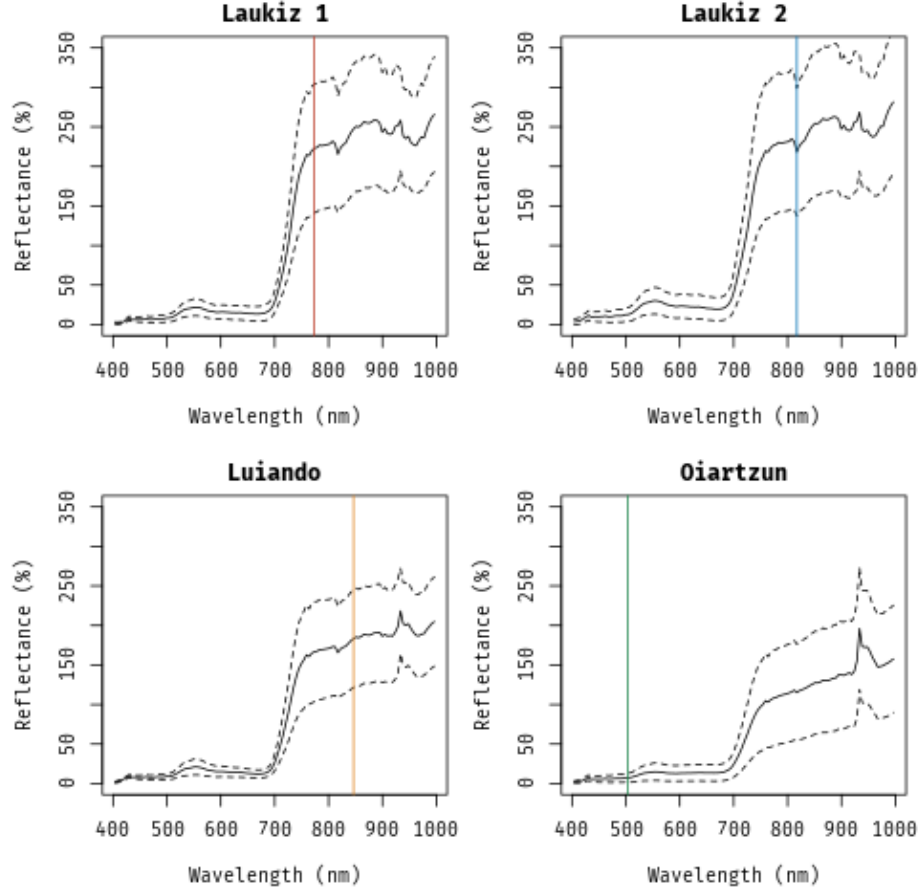


Figure 3: Spectral signatures (mean and standard deviation) of each plot. The colored lines show the most important band for each plot, respectively: Band 80 (773nm, red), band 89 (817nm, blue), band 95 (846nm, orange), band 23 (503nm, green).

a mean tree density of 61.59 m ??) while all other plots are more dense (34.64 (Laukiz 1), 33.01 (Luiando), 34.96 (Oiartzun)) (Figure 4).

4.2. Predictive performance

150 Ridge Regression (RR) shows the lowest error for three out of four plots (for Luiando *elasticnet* shows a slightly better performance) (Table 2). The magnitude of difference for RR compared to the other penalties for the plots in

which RR showed the best performance ranges between XX and XX percent. For the merged dataset, all penalties show a similar mean predictive performance
155 that outperform all single plot models besides the Laukiz 2 model.

When comparing the mean predictive performance of the plot level model against the performance of the super model at the plot level (when the respective plot served as the test set), the supermodel also outperforms the Laukiz 2 model (27.94 vs 30.37 RMSE) (Table 3).

160 The worst performance of the supermodel on the fold level is reported for Luiando (69.72 RMSE) while for the single plot models Oiartzun shows the highest error (106.65 RMSE).

Laukiz 2 showed contrary results compared to all other plots when linking RMSE against coefficient of variation and mean point density (??). Comparing
165 RMSE against $CV/skewness$ shows a $\log_2(-x)$ relationship.

4.2.1. Variable importance

NRIs using bands in the wavelength range of 770 nm - 820 nm (band 80 - band 89), which belongs to the infrared region, appear most often among the ten highest coefficient estimates across all plots (??). Only one vegetation index

Table 2: Four-fold spatial CV performances of RR, SVM and xgboost using RMSE as the error measure. Mean and standard deviation are shown.

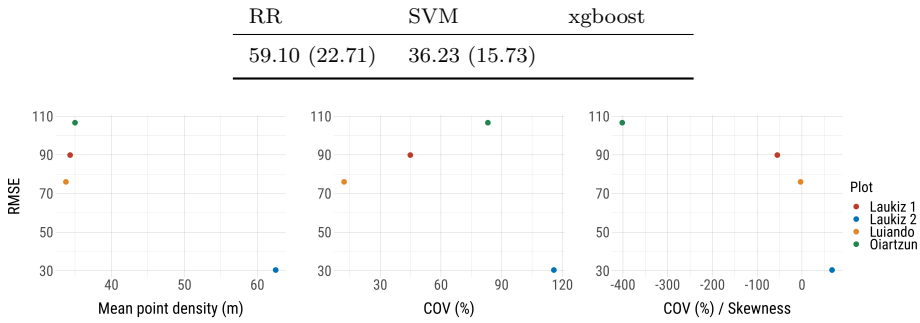


Figure 4: RMSE vs. mean point density, coefficient of variation and coefficient of variation / skewness.

Table 3: Predictive performance of xgboost using the merged dataset (supermodel) and observations on a plot level only (single plot) with RMSE as the error measure. The performance estimates for "All Plots" correspond to the fold for which the respective plot was serving as the test set. For "single plot", the values correspond to the mean value of the SpCV at the repetition level (10 folds, 20 repetitions), scored by models fitted on the data of the respective plot only.

Plot/Data	All Plots (Block CV)	Single plot (SpCV)
Laukiz 1	58.95	89.89
Laukiz 2	27.94	30.37
Luiando	69.72	76.02
Oiartzun	58.09	106.65

(Datt3) showed up among the most important predictors (Laukiz 1). Luiando and Oiartzun also preferred bands with longer (938.39 nm (band 114) - 996.31 nm (band 126)) and shorter wavelengths (480.30 nm (band 18) - 503.26 (band 23)). The first range again belongs to the infrared region while the second is within the region of the visible light, transitioning from blue to green.

5. Discussion

5.1. Index derivation

The exact number of contributing pixels to the final index value of an observations cannot be determined as it depends on the location of the tree within the pixel grid. If a tree is located at the border of a pixel, a buffer of e.g. three meters will include more pixels than if the point is located at the center of a pixel. Also, if a tree is located at the border of the plot, some directions of the buffer will not contain image values.

5.2. Plot characteristics

For Laukiz1, Luiando and Oiartzun RMSE seems to increase with a higher point density at a first glance. However, the point densities of these plots are very similar (33.7 m - 35.01 m) and should be interpreted as a group instead of

single values. With Laukiz2 being completely off from the other plots in terms of mean point density, no pattern can be extracted from this result. Linking RMSE vs coefficient of variation shows the same relationship as linking against mean
 190 point density. The interesting $\log_2(-x)$ relationship for RMSE vs. coefficient of variation / skewness should be interpreted with caution: The sample size of four plots is not representative to make general statements here. This finding should be verified with more observations in future studies.

5.3. Variable importance

195 References

- Adamczyk, J., & Osberger, A. (2015). Red-edge vegetation indices for detecting and assessing disturbances in Norway spruce dominated mountain forests. *International Journal of Applied Earth Observation and Geoinformation*, 37, 90–99. doi:10/f64b6c.
- 200 Belgiu, M., & Drăguț, L. (2016). Random forest in remote sensing: A review of applications and future directions. *ISPRS Journal of Photogrammetry and Remote Sensing*, 114, 24–31. doi:10/f8ndk8.
- Bergstra, J., & Bengio, Y. (2012). Random Search for Hyper-parameter Optimization. *J. Mach. Learn. Res.*, 13, 281–305. 01590.
- 205 Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., & Jones, Z. M. (2016). mlr: Machine learning in R. *Journal of Machine Learning Research*, 17, 1–5.
- Bischl, B., Richter, J., Bossek, J., Horn, D., Thomas, J., & Lang, M. (2017). mlrMBO: A Modular Framework for Model-Based Optimization of Expensive
 210 Black-Box Functions. *ArXiv e-prints*, . arXiv:1703.03373.
- Brenning, A. (2012). Spatial cross-validation and bootstrap for the assessment of prediction rules in remote sensing: The R package sperrorest. In *2012 IEEE International Geoscience and Remote Sensing Symposium*. IEEE. doi:10.1109/igarss.2012.6352393 00052 R package version 2.1.0.

- 215 Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining KDD '16* (pp. 785–794). New York, NY, USA: ACM. doi:10.1145/2939672.2939785.
- Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization Paths for
220 Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software*, 33, 1–22. doi:10/bb3d.
- Ganley, R. J., Watt, M. S., Manning, L., & Iturritxa, E. (2009). A global climatic risk assessment of pitch canker disease. *Canadian Journal of Forest Research*, 39, 2246–2256. doi:10/bmj3nk.
- 225 Hutter, F., Hoos, H. H., & Leyton-Brown, K. (2011). Sequential Model-Based Optimization for General Algorithm Configuration. In *Lecture Notes in Computer Science* (pp. 507–523). Springer Berlin Heidelberg. doi:10.1007/978-3-642-25566-3_40 00686.
- Iturritxa, E., Mesanza, N., & Brenning, A. (2014). Spatial analysis of the risk
230 of major forest diseases in Monterey pine plantations. *Plant Pathology*, 64, 880–889. doi:10/gdq9pb. 00006.
- Iturritxa, E., Trask, T., Mesanza, N., Raposo, R., Elvira-Recueno, M., & Pat-
ten, C. L. (2017). Biocontrol of *Fusarium circinatum* Infection of Young *Pinus radiata* Trees. *Forests*, 8, 32. doi:10/f9t3d8.
- 235 Jiang, Y., Wang, T., de Bie, C. A. J. M., Skidmore, A. K., Liu, X., Song, S., Zhang, L., Wang, J., & Shao, X. (2014). Satellite-derived vegetation indices contribute significantly to the prediction of epiphyllous liverworts. *Ecological Indicators*, 38, 72–80. doi:10/f5q4b4.
- Jones, D. R., Schonlau, M., & Welch, W. J. (1998). Efficient global optimization
240 of expensive black-box functions. *Journal of Global Optimization*, 13, 455–492. doi:10/fg68nc.

- Karatzoglou, A., Smola, A., Hornik, K., & Zeileis, A. (2004). Kernlab – An S4 Package for Kernel Methods in R. *Journal of Statistical Software*, 11, 1–20. doi:10/gdq9pc. R package version 0.9-25.
- 245 Lary, D. J., Alavi, A. H., Gandomi, A. H., & Walker, A. L. (2016). Machine learning in geosciences and remote sensing. *Geoscience Frontiers*, 7, 3–10. doi:10/f79ddn.
- Lehnert, L. W., Meyer, H., & Bendix, J. (2018). *Hsdar: Manage, Analyse and Simulate Hyperspectral Data in R*. R package version 0.7.1.
- 250 Lelong, C. C. D., Roger, J.-M., Brégand, S., Dubertret, F., Lanore, M., Sitorus, N. A., Raharjo, D. A., & Caliman, J.-P. (2010). Evaluation of Oil-Palm Fungal Disease Infestation with Canopy Hyperspectral Reflectance Data. *Sensors*, 10, 734–747. doi:10/bb8wm6.
- Lu, D., Chen, Q., Wang, G., Liu, L., Li, G., & Moran, E. (2016). A survey
255 of remote sensing-based aboveground biomass estimation methods in forest ecosystems. *International Journal of Digital Earth*, 9, 63–105. doi:10/gdthzv.
- Martinez del Castillo, E., García-Martin, A., Longares Aladrén, L. A., & de Luis, M. (2015). Evaluation of forest cover change using remote sensing techniques and landscape metrics in Moncayo Natural Park (Spain). *Applied Geography*,
260 62, 247–255. doi:10/gdthzt.
- Mesanza, N., Iturrutxa, E., & Patten, C. L. (2016). Native rhizobacteria as bio-control agents of *Heterobasidion annosum* s.s. and *Armillaria mellea* infection of *Pinus radiata*. *Biological Control*, 101, 8–16. doi:10/f8xnp3.
- Michez, A., Piégay, H., Lisein, J., Claessens, H., & Lejeune, P. (2016). Classi-
265 fication of riparian forest species and health condition using multi-temporal and hyperspatial imagery from unmanned aerial system. *Environmental Monitoring and Assessment*, 188, 146. doi:10/f8q9wp.

- 270 Probst, P., Wright, M., & Boulesteix, A.-L. (2018). Hyperparameters and Tuning Strategies for Random Forest. *ArXiv e-prints*, . arXiv:1804.03515.00000.
- R Core Team (2017). *R: A Language and Environment for Statistical Computing*. Vienna, Austria. 00000 R version 3.3.3.
- 275 Sexton, J. O., Noojipady, P., Anand, A., Song, X.-P., McMahon, S., Huang, C., Feng, M., Channan, S., & Townshend, J. R. (2015). A model for the propagation of uncertainty from continuous estimates of tree cover to categorical forest cover and change. *Remote Sensing of Environment*, 156, 418–425. doi:10/f6v7zc.
- 280 Sinha, S., Jeganathan, C., Sharma, L. K., & Nathawat, M. S. (2015). A review of radar remote sensing for biomass estimation. *International Journal of Environmental Science and Technology*, 12, 1779–1792. doi:10/gdthzw.
- Vapnik, V. (1998). The Support Vector Method of Function Estimation. In *Nonlinear Modeling* (pp. 55–85). Springer US. doi:10.1007/978-1-4615-5703-6_3.
- 285 Zhang, K., Thapa, B., Ross, M., & Gann, D. (2016). Remote sensing of seasonal changes and disturbances in mangrove forest: A case study from South Florida. *Ecosphere*, (p. e01366). doi:10.1002/ecs2.1366@10.1002/(ISSN)2150-8925.ExtremeColdSpells.

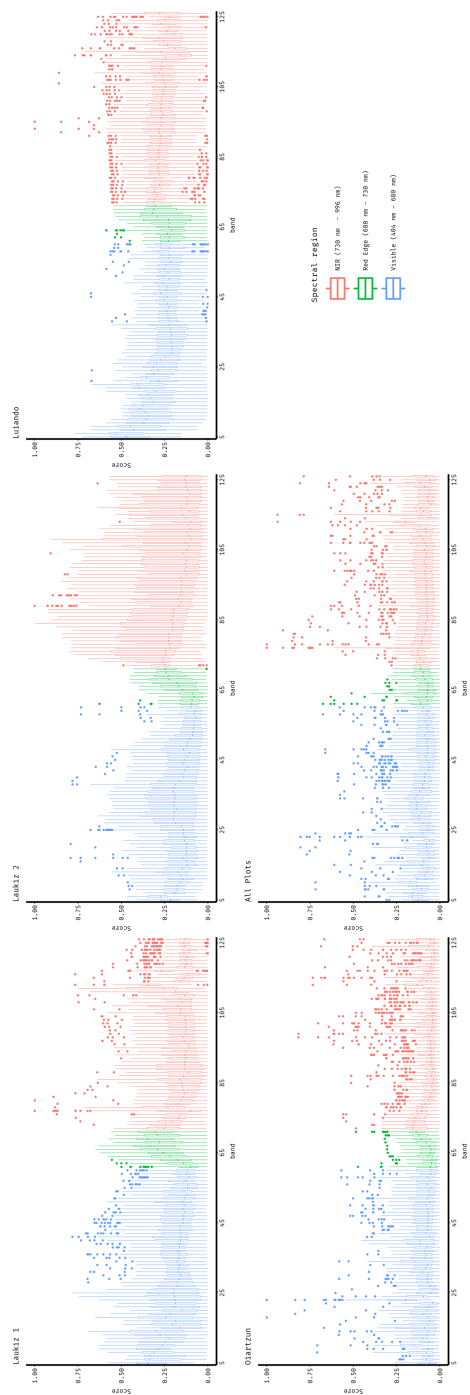


Figure 5: test