

Modeling defoliation of Pinus Radiata trees using hyperspectral remote sensing data

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.,

*Corresponding author

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

2010). However, spectral (vegetation-)indices have shown the potential to contribute 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 fungi 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 (?). 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). In total 1750 observations are available (*Laukiz 1* = 479, *Laukiz 2* = 451, *Luiando* = 291, *Oiartzun* = 529). The data was surveyed 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 using a AISAEAGLE-II sensor. All preprocessing steps (geometric, radiometric, atmospheric) have been conducted by the Institut Cartografic i Geologic de Catalunya (ICGC). The first four bands were corrupted, leaving 122 bands with valid information. Additional metadata information is available in Table 1:

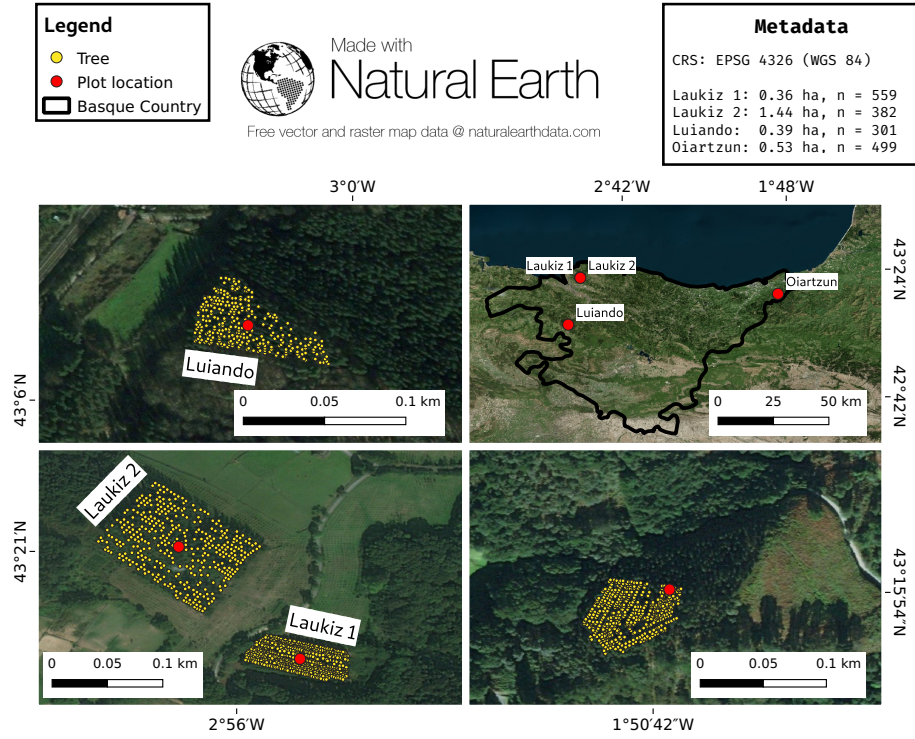


Figure 1: Information about the plot locations, the area of hyperspectral coverage and the number of trees per plot.

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), *kerlab* (?) (Sup-

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

port Vector Machine), *?)* and *glmnet* (Friedman et al., 2010) (Ridge Regression). We used the R package *mlr* for 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 parallelization (Bischl et al., 2016).

3.1. Derivation of indices

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.

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 ($\frac{125 \times 126}{2}$). Due to four corrupted bands and some other 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 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)

(?) 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. The probably most popular machine-learning algorithm, Random Forest, was not considered for this study: 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).

110 3.3.1. Performance estimation

The algorithms were benchmarked in two ways: (1) Using spatial cross-validation (CV) for each plot using on the k-means clustering approach of ?. 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 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

120 To tune the hyperparameters of the algorithms, we used Sequential-based Model Optimization (SMBO) via the R package *mlrMBO* (Bischl et al., 2017). This Bayesian approach first composes n randomly chosen hyperparameter settings 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 criterion, defined by the user, is reached (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

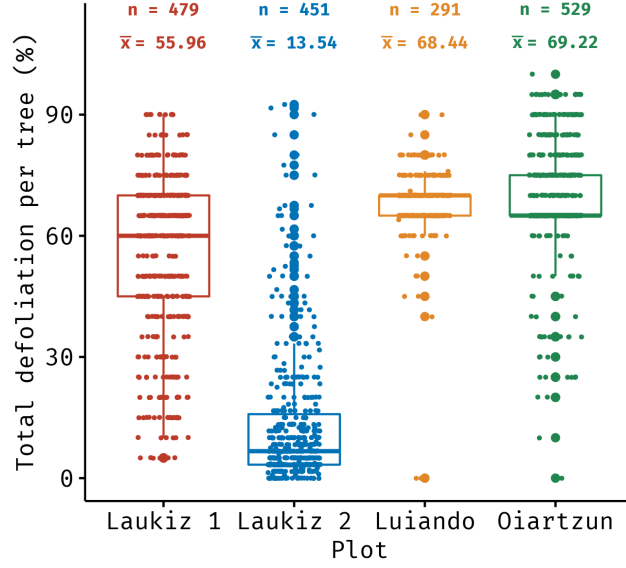


Figure 2: Descriptive statistics of the response variable *defoliation*.

130 it substantially reduces the tuning budget which is needed to 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).

3.4. Variable importance

To find indices that contributed most to model performance, we used the internal variable importance measure of the *xgboost* algorithm. The score is
 135 calculated by taking the contribution of each feature for each tree in the fitted model. The higher the score of a variable, the more important it is for the fitted model when making predictions (Chen & Guestrin, 2016).

4. Results

140 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.

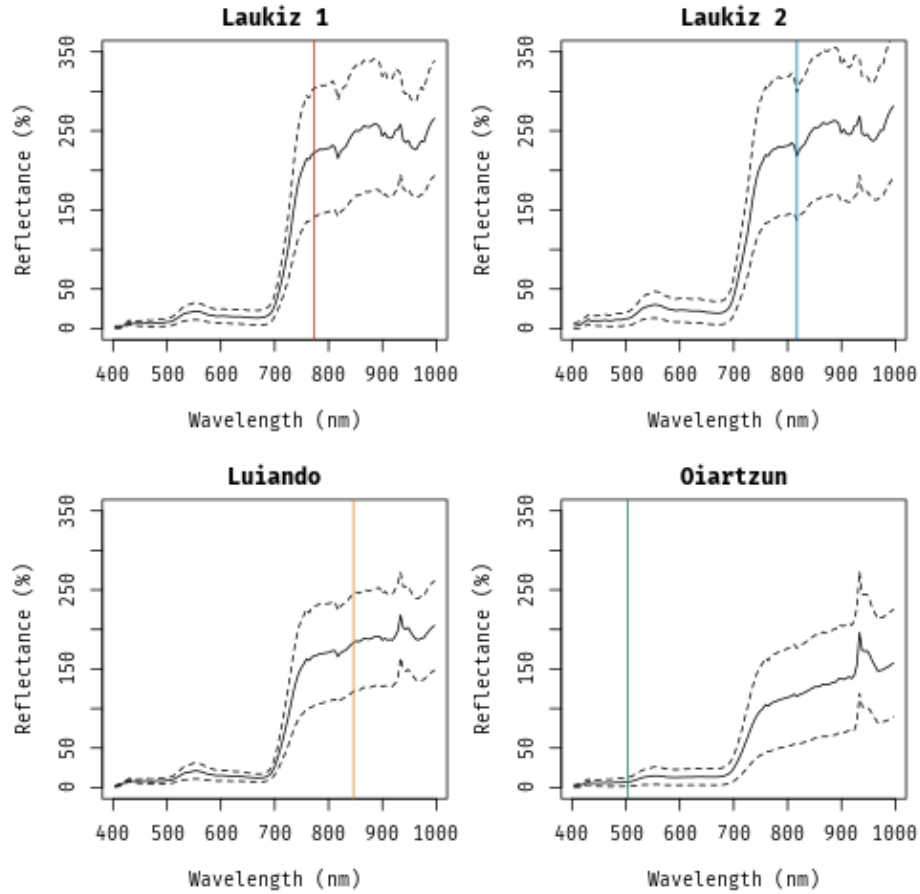


Figure 3: Spectral signatures (mean and standard deviation) of each plot. The colored lines show the most important band for each plot, respectively: Band 80 (773nm, red), band 89 (817nm, blue), band 95 (846nm, orange), band 23 (503nm, green).

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 a mean tree density of 61.59 m ??) while all other plots are more dense (34.64

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

RR	SVM	xgboost
59.10 (22.71)	36.23 (15.73)	33.26 (16.61)

(Laukiz 1), 33.01 (Luiando), 34.96 (Oiartzun)) (Figure 4).

4.2. Predictive performance

4.2.1. Algorithm benchmarking

The *xgboost* algorithm shows the lowest error (33.26 RMSE) when bench-
155 marking the learners on the complete dataset of all plots (Table 2). While the
SVM performance is only slightly worse (36.23 RMSE), RR shows a large drop
in performance compared to xgboost (59.10 RMSE).

4.2.2. Single models vs. super model

When comparing the mean predictive performance of models fitted at the
160 plot level against the performance of the model that was fitted using all data, the
plot-level models show a better performance in all cases (Table 3). The highest
difference between both models types occurs for *Laukiz2* with a difference of
34.51 RMSE.

Table 3: Predictive performance of *xgboost* using all observations (All Observations) and observations from specific plots only (Single Plot Observations) with RMSE as the error measure. The performance estimates for "All Observations" correspond to the fold for which the respective plot was serving as the test set. Column "Single Plot Observations", shows the mean performances at the repetition level of a spatial CV (5 folds, 5 repetitions), scored by using data of the respective plot only.

Plot/Data	All Observations (Block CV)	Single Plot Observations (SpCV)
Laukiz 1	22.03	19.18
Laukiz 2	51.75	17.24
Luiando	13.20	8.30
Oiartzun	32.97	14.40

4.2.3. RMSE vs. plot characteristics

165 An increase of the error rate was observed with an increases of descriptive plot
measures such as mean point density and the coefficient of variation (based on
the response *defoliation*) (Figure 4).

4.3. Variable importance

The seven most important features of the super model in this study were veg-
170 etation indices with *EVI* (Huete et al., 1997) being the most important one
(Figure 5).

$$EVI = 2.5 * \frac{R_{800} - R_{670}}{R_{800} - (6 * R_{670}) - (7.5 * R_{475}) + 1} \quad (2)$$

where R = Reflectance at the respective wavelength.

Vegetation index *GDVI* appears three times among the first seven most impor-
tant features (Figure 5) with different . This is because it was computed four
175 times, with n ranging from 1 - 4 (Wu et al., 2008):

$$GDVI = \frac{R_{800}^n - R_{680}^n}{R_{800}^n + R_{680}^n} \quad (3)$$

The seven most important features (*EVI*, *GDVI4*, *D1*, *GDVI3*, *GDVI2*, *mNDVI*
and *mSR*) show a substantial difference compared to all following variables
(Figure 5).

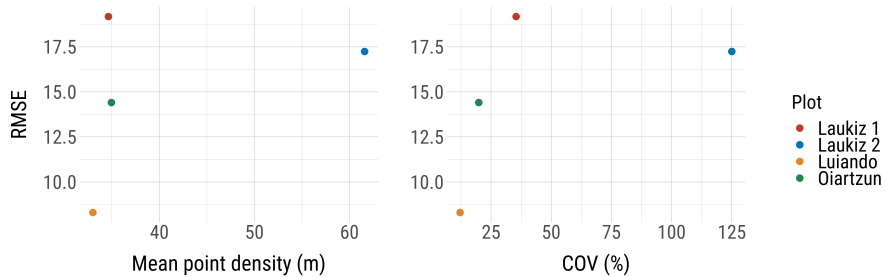


Figure 4: RMSE vs. mean point density and coefficient of variation (defoliation).

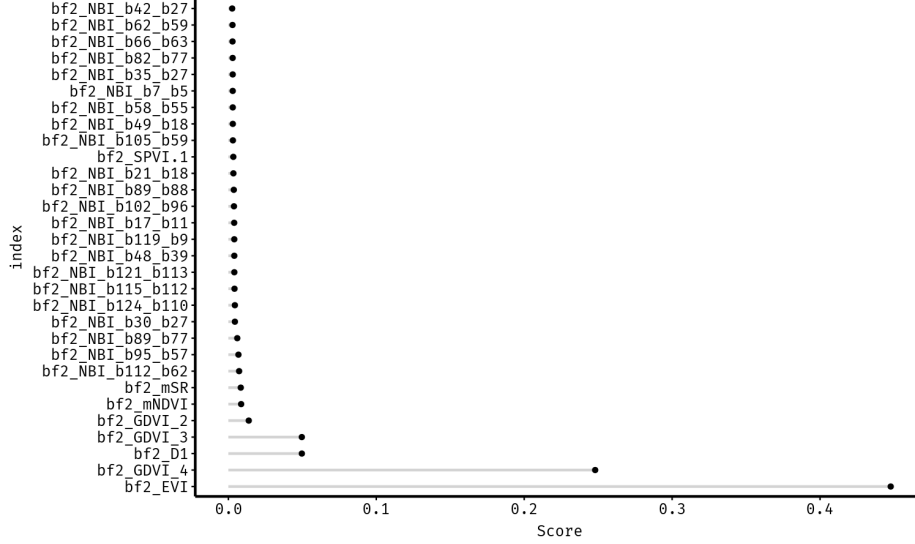


Figure 5: The 30 most important variables as estimated by the internal variable importance measure of the *xgboost* algorithm. The higher the score, the more important the feature. "bf2" means that a buffer of 2 meter was used to extract the variable information to the tree observation. "NRI" means that a normalized ratio index with the subsequent bands was calculated. Features without "NRI" prefix are vegetation indices, e.g. "bf2_EVI".

Table 4: Formulas of the five most important vegetation indices of the super model. R = Reflectance at wavelength, D = First derivation of reflectance value at wavelength.

Acronym	Name	Formula	Reference
EVI	Enhanced vegetation index	$2.5 * \frac{R_{800} - R_{670}}{R_{800} - (6 * R_{670}) - (7.5 * R_{475}) + 1}$	Huete et al. (1997)
GDVI	Generalized DVI*	$\frac{R_{800}^n - R_{680}^n}{R_{800}^n + R_{680}^n}$	Wu et al. (2008)
D1	Derivative Index	$\frac{D_{730}}{D_{706}}$	Zarco-Tejada et al. (2003)
mNDVI	Normalized DVI*	$\frac{R_{800} - R_{680}}{(R_{800} + R_{680} - 2 * R_{445})}$	Sims & Gamon (2002)
mSR	Simple Ratio Index	$\frac{R_{800} - R_{445}}{R_{680} - R_{445}}$	Sims & Gamon (2002)

* Difference Vegetation Index

5. Discussion

180 5.1. Index derivation

The exact number of contributing pixels to the final index value of an observation 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
185 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
190 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 point density. The interesting $\log_2(-x)$ relationship for RMSE vs. coefficient
195 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

References

- 200 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.
- Belgiu, M., & Drăguț, L. (2016). Random forest in remote sensing: A review
205 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.
- 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 Black-Box Functions. *ArXiv e-prints*, . [arXiv:1703.03373](https://arxiv.org/abs/1703.03373).
- 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. 05097.
- Huete, A. R., Liu, H. Q., Batchily, K., & van Leeuwen, W. (1997). A comparison of vegetation indices over a global set of TM images for EOS-MODIS. *Remote Sensing of Environment*, *59*, 440–451. doi:10/bgtpgv.
- 225 Iturritxa, E., Mesanza, N., & Brenning, A. (2014). Spatial analysis of the risk of major forest diseases in Monterey pine plantations. *Plant Pathology*, *64*, 880–889. doi:10/gdq9pb.
- Iturritxa, E., Trask, T., Mesanza, N., Raposo, R., Elvira-Recuenco, M., & Patten, C. L. (2017). Biocontrol of *Fusarium circinatum* Infection of Young *Pinus radiata* Trees. *Forests*, *8*, 32. doi:10/f9t3d8.
- 230 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.

- 235 Jones, D. R., Schonlau, M., & Welch, W. J. (1998). Efficient global optimization
of expensive black-box functions. *Journal of Global Optimization*, 13, 455–
492. doi:10/fg68nc.
- 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.
240 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.
- 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
245 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
of remote sensing-based aboveground biomass estimation methods in forest
ecosystems. *International Journal of Digital Earth*, 9, 63–105. doi:10/gdthzv.
- 250 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*,
62, 247–255. doi:10/gdthzt.
- Mesanza, N., Iturritya, E., & Patten, C. L. (2016). Native rhizobacteria as bio-
255 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-
fication of riparian forest species and health condition using multi-temporal
and hyperspatial imagery from unmanned aerial system. *Environmental Mon-
260 itoring and Assessment*, 188, 146. doi:10/f8q9wp.

- 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. 88058 R version 3.4.4.
- 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.
- Sims, D. A., & Gamon, J. A. (2002). Relationships between leaf pigment content and spectral reflectance across a wide range of species, leaf structures and developmental stages. *Remote Sensing of Environment*, 81, 337–354. doi:10/fb9nnj.
- 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.
- Wu, C., Niu, Z., Tang, Q., & Huang, W. (2008). Estimating chlorophyll content from hyperspectral vegetation indices: Modeling and validation. *Agricultural and Forest Meteorology*, 148, 1230–1241. doi:10/dhcp6r.
- Zarco-Tejada, P. J., Pushnik, J. C., Dobrowski, S., & Ustin, S. L. (2003). Steady-state chlorophyll a fluorescence detection from canopy derivative reflectance and double-peak red-edge effects. *Remote Sensing of Environment*, 84, 283–294. doi:10/c8gjtt.
- 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.