

1 **Title:** Enhancing of accuracy assessment for forest above-ground biomass estimates obtained
2 from remote sensing via hypothesis testing and overfitting evaluation.

3

4 **Authors:** Valbuena R^{1,2*}, Hernando A³, Manzanera JA³, Görgens EB⁴, Almeida DRA⁵, Mauro
5 F³, García-Abril A³, Coomes DA¹.

6 1: University of Cambridge, Department of Plant Sciences, Forest Ecology and Conservation.
7 Downing Street, CB2 3EA Cambridge, UK

8 2: University of Eastern Finland, Faculty of Forest Sciences, P.O. Box 111, Joensuu, Finland

9 3: Universidad Politecnica de Madrid, College of Forestry and Natural Environment, Research
10 Group SILVANET, Ciudad Universitaria, 28040 Madrid, Spain.

11 4: Universidade Federal dos Vales do Jequitinhonha e Mucuri, Department of Forestry,
12 Campus JK, CEP 39100-000, Diamantina, Brazil.

13 5: University of São Paulo, Luiz de Queiroz College of Agriculture, Department of Forest
14 Sciences, Av. Pádua Dias, 11. CEP 13418-900 Piracicaba, Brazil.

15 * Corresponding author: rv314@cam.ac.uk

16

17

18

19

20

21 **Abstract**

22 The evaluation of accuracy is essential for assuring the reliability of ecological models. Usually,
23 the accuracy of above-ground biomass (*AGB*) predictions obtained from remote sensing is
24 assessed by the mean differences (*MD*), the root mean squared differences (*RMSD*), and the
25 coefficient of determination (R^2) between observed and predicted values. In this article we
26 propose a more thorough analysis of accuracy, including a hypothesis test to evaluate the
27 agreement between observed and predicted values, and an assessment of the degree of
28 overfitting to the sample employed for model training. Using the estimation of forest *AGB* from
29 LIDAR and spectral sensors as a case study, we compared alternative prediction and variable
30 selection methods using several statistical measures to evaluate their accuracy. We showed that
31 the hypothesis tests provide an objective method to infer the statistical significance of
32 agreement. We also observed that overfitting can be assessed by comparing the inflation in
33 residual sums of squares experienced when carrying out a cross-validation. Our results suggest
34 that this method may be more effective than analysing the deflation in R^2 . We proved that
35 overfitting needs to be specifically addressed since, in light of *MD*, *RMSD* and R^2 alone,
36 predictions may apparently seem reliable even in clearly unrealistic circumstances, for instance
37 when including too many predictor variables. Moreover, Theil's partial inequality coefficients,
38 which are employed to resolve the proportions of the total errors due to the unexplained
39 variance, the slope and the bias, may become useful to detect averaging effects common in
40 remote sensing predictions of *AGB*. We concluded that statistical measures of accuracy,
41 precision and agreement are necessary but insufficient for model evaluation. We therefore
42 advocate for incorporating evaluation measures specifically devoted to testing observed-versus-
43 predicted fit, and to assessing the degree of overfitting.

44 **Key words:** model assessment; overfitting; Theil's partial inequality coefficients; LIDAR.

45 **Introduction**

46 The evaluation of accuracy is an essential step indicating the reliability of a given prediction
47 method, thereby informing researchers about the level of confidence they should place in their
48 predictions and allowing them to compare alternatives (Tedeschi, 2006). Accuracy assessment
49 must be supported by rigorous statistical inference, with the ultimate target of evaluating the
50 ability to generalize from the sample data to the population of interest (Särndal et al., 1992;
51 Naeset, 2002; McRoberts et al., 2013; Asner & Mascaro, 2014; Chen et al., 2015; Mauro et
52 al., 2016). Several quantitative techniques can be used to verify if the predicted values differ
53 significantly from the observed, including squared sums of prediction errors (Wallach and
54 Goffinet, 1989), coefficient of determination (R^2) or other correlation-like measures (Willmott,
55 1981), a reliability index (Leggett & Williams, 1981), distribution hypothesis testing (Freese,
56 1960), and regression of predicted versus observed (Theil, 1958; Graybil, 1976; Reynolds &
57 Chung, 1986) or *vice versa* (Piñeiro et al 2008). The advantages and disadvantages of these
58 approaches have been evaluated (e.g., Fox, 1981; Willmott, 1982). Since each scientific
59 application has its own particularities, it is recognised that no single measure of model
60 performance is appropriate in all circumstances (Smith & Rose, 1995). This article explores
61 open questions on accuracy assessment in the context of predicting forest above-ground
62 biomass (*AGB*) from remote sensing sources. The accuracy assessment measures proposed here
63 can nonetheless be generalizable to many other contexts where predictions of ecological
64 variables from different sources of auxiliary information are sought.

65 *Common measures for accuracy assessment and aspects needing revision*

66 When assessing the performance of their methods, remote sensing researchers usually report:
67 (1) mean difference between observed and predicted values, which evaluates the degree of
68 under- or over-prediction of the dependent variable, *AGB* in this case; (2) the precision of the

69 prediction, often reporting the root mean squared differences (*RMSD*); and (3) the level of
70 agreement between observed and predicted values, commonly considered by indicating their
71 R^2 (e.g., Zhao et al., 2009; Erdody & Moskal 2010; McInerney et al., 2010; d'Oliveira et al
72 2012; Chen & Zhu, 2013; Straub et al., 2013; Asner & Mascaró, 2014; Valbuena et al., 2014).
73 There is, however, no strong consensus, and it is not uncommon to find studies reporting
74 alternative or complementary measures, for instance analysing the regression of predicted
75 versus observed (Bright et al., 2012; Wing et al. 2012) or alternatives to R^2 (Yebra & Chuvieco,
76 2009; García et al., 2010; Almeida et al., 2016). Some studies (e.g, d'Oliveira et al., 2012;
77 Estornell, et al. 2014) perform hypothesis tests comparing distributions, similar to those
78 suggested by Freese (1960). Moreover, the degree of overfitting is rarely accounted for
79 (Valbuena et al., 2013a; Latifi et al., 2015a; Almeida et al., 2016), despite of being a common
80 pitfall in predictive modelling (Weisberg, 1985; Hurvich & Tsai, 1989; Hawkins, 2004). In the
81 context of remote sensing prediction of forest *AGB*, we detected two key aspects of accuracy
82 lacking consensus (plus a third additional one, see Valbuena et al., 2018):

83 *Evaluating regression of observed versus predicted.* Piñeiro et al. (2008) argued that the correct
84 assessment is done by setting the predicted values as independent variable (in the x-axis) and
85 the observed values as dependent variable (in the y-axis), to properly evaluate their regression
86 coefficients (Reynolds & Chung, 1986). However, when evaluating remote sensing predictions
87 of forest attributes, many authors have presented predicted (in the y-axis) vs. observed (in the
88 x-axis) instead (e.g., McRoberts et al., 2002; Holmgren et al., 2008; Zhao et al., 2009;
89 McInerney et al., 2010; Chen & Zhu 2013; Valbuena et al., 2014). Furthermore, they usually
90 lack reporting the regression of observed against predicted (e.g., Naesset, 2002; García et al.
91 2010; Straub et al., 2013). Although some report the coefficients (e.g., Yebra & Chuvieco,
92 2009; Bright et al., 2012; Wing et al. 2012), they may still miss the hypothesis test suggested
93 by Piñeiro et al. (2008). There have therefore not been reports on the importance of carrying

94 out these hypothesis tests in the context of remote sensing predictions of *AGB*. Complementary
95 statistics may also be included in order to fully comprehend the source of prediction errors,
96 such as Theil's (1958) partial inequality coefficients (Smith & Rose, 1995). They disaggregate
97 the total error into model variance (unsystematic error), bias (systematic error), and slope
98 (averaging effects) (Paruelo et al., 1998). To our knowledge, these coefficients have not been
99 employed in the context of remote sensing estimates of forest characteristics before.

100 *The degree of overfitting to the sample.* Franco-Lopez et al. (2001) argued that statistical
101 measures to assess model overfitting should be included when reporting the accuracy
102 assessment of remote sensing estimates. Those measures of overfitting have been, however,
103 largely overlooked in remote sensing estimations of forest attributes (Latifi et al., 2015a).
104 Overfitting is usually prevented beforehand by avoiding over-parameterization with variable
105 selection methods (e.g., Naesset, 2002; Hudak et al., 2006; García et al., 2010; Wing et al.,
106 2012; Spriggs et al., 2015). These methods, however, have been suspected of being insufficient
107 to truly avoid model overfitting (Allen, 1974; Vanclay & Skovsgaard, 1997; Hurvich & Tsai,
108 1989; Rencher & Pun, 1993). As an alternative, some authors recommend preventing model
109 overfitting using replication methods such as cross-validation, and compare their results against
110 model residuals (Weisberg, 1985; Hawkins, 2004). These would also be particularly convenient
111 for non-parametric machine learning methods, whose flexibility makes them especially prone
112 to overfitting (Franco-Lopez et al., 2001; Hawkins, 2004), and which are of widespread use in
113 remote sensing predictions of forest attributes (McRoberts et al., 2002; Hudak et al., 2008;
114 Packalén & Maltamo, 2008; McInerney et al., 2010). However, overfitting is rarely addressed
115 in the context of remote sensing predictions of forest variables (Franco-Lopez et al., 2001;
116 Valbuena et al., 2013a; Latifi et al., 2015a; Almeida et al., 2016).

117 These alternative methods for testing the reliability of *AGB* predictions obtained by using
118 remotely sensed sources may also be employed to minimise errors in the estimation of

119 ecological variables in general. Results may therefore be relevant to other contexts too, for
120 example studies on ecosystem management responses to climate change or habitat suitability
121 for fauna, where the use of models to predict ecological attributes from auxiliary variables is
122 common.

123 *Objectives*

124 The objective of this research is to call into question the sufficiency of statistical measures
125 commonly used for accuracy assessment of predictions of ecological variables from auxiliary
126 information, and suggest the convenience of incorporating additional ones, with a focus on
127 remote sensing estimations of forest *AGB*. Our hypothesis is that the statistics usually reported
128 in *AGB* assessments may be insufficient for accepting the degree of agreement between predicted
129 and observed as reliable, and also that the fact that overfitting effects may remain undetected.
130 This article therefore aspires to present a thorough analysis of accuracy that applies to
131 ecological modelling in general, and to explain how to interpret the suggested statistical metrics
132 for readers unfamiliar to them in the given context.

133 **Material and Methods**

134 *Field and Remote Sensing Datasets*

135 The field datasets consisted of $n = 37$ plots surveyed during summer 2006 in the Scots pine
136 (*Pinus sylvestris* L.) dominated forests of Valsaín (Spain, approx. lat.: 41°04' N, lon.: 4°09' W;
137 1.3-1.5 km a.s.l.). These plots consisted of two concentric circles of radii 10 and 20 m.
138 Diameters at breast height (*dbh*, cm) were measured for every tree located within the inner sub-
139 plot, whereas at the outer sub-plot only those with $dbh > 10$ cm were measured (Valbuena et
140 al., 2013b). Differentially-corrected global navigation satellite systems (GNSS) were used to

141 obtain the positions of these plots with centimetre accuracy (Valbuena et al., 2012), enabling to
142 link the field and remote sensing information.

143 Locally-adjusted tree allometry specific for *P. sylvestris* was employed to obtain the above-
144 ground biomass (*agb*, kg) of each individual tree from the field measurements (Montero et al.,
145 2005):

$$146 \quad agb = 0.08439 \cdot dbh^{2.41194} \quad (1)$$

147 These tree-level *agb* estimations were aggregated to plot-level totals (*AGB*, Mg·ha⁻¹), after
148 referring each of them to per-hectare equivalents according to the differing size of the sub-plot
149 from which each tree was sampled (inner or outer). In this study, we used *AGB* as a response
150 variable to be predicted throughout the target forest by using the remote sensing predictor
151 variables.

152 The predictor variables were statistical metrics describing the distributions of signals received
153 at those same field plots from both active LIDAR and passive multispectral sensors. This
154 remotely sensed information was acquired on September 10, 2006, from a laser scanner ALS50-
155 II (Leica Geosystems, Switzerland) and a digital mapping camera system (Zeiss-Intergraph,
156 Germany). Simultaneously operating onboard a plane flying at a height of 1500 m, the LIDAR
157 dataset was obtained with an average scan density of 1.15 pulses·m⁻², whereas images had
158 spatial resolutions of 15 cm from panchromatic and 60 cm for multispectral. A back-projection
159 data fusion algorithm using information from on-flight GNSS and inertial navigation systems
160 assured a nearly perfect fit of all the sensor and field information (Valbuena, 2014). Back-
161 projecting consists in mathematically rendering the position of each LIDAR return onto the
162 camera at the time of exposure, retrieving back its radiometric information and effectively
163 colouring the LIDAR return with an accuracy close to pixel size (Valbuena et al., 2011). Returns
164 obtained from the LIDAR sensor, considered to represent the ground – by means of Axelsson’s

165 (2000) classification algorithms –, were interpolated into a digital terrain model, which was
166 used as a reference from which to calculate the heights above ground (h , m) for every single
167 LIDAR return. The radiometric information acquired from the digital camera system was
168 employed to calculate a value of normalised difference vegetation index ($NDVI$; Rouse et al.,
169 1974) correspondent to each LIDAR first return. Using FUSION software (USDA Forest
170 Service; McGaughey, 2012), the returns backscattered from each field plot were extracted, and
171 several metrics describing the distributions and relative proportions of h and $NDVI$ with each
172 plot were computed (Manzanera et al., 2016). All these metrics were employed as initial dataset
173 of predictors in all the predictive procedures.

174 *Modelling biomass from airborne remote sensing data*

175 Three prediction methods commonly employed for forest AGB predictions from remote sensing
176 were compared within the R statistical environment (version 3.3.1; R Development Core Team,
177 2016):

178 *Non-parametric modelling based on the most similar neighbour (MSN)* method to obtain AGB
179 predictions (Moeur & Stage, 1995) was applied using the “yaImpute” package of R (version
180 1.0-18; Crookston & Finley, 2007). MSN belongs to a type of non-parametric imputation
181 approaches known as nearest neighbour methods and commonly abbreviated as k -NN (Franco-
182 Lopez et al., 2001; McRoberts et al., 2002; McInerney et al., 2010), k being the number of
183 neighbours used in the algorithm. In the particular case of MSN, the feature space – where
184 distances to neighbours are measured – is modified according to canonical correlation
185 projectors (Hudak et al., 2008; Packalén & Maltamo, 2008). The nearest neighbour algorithm
186 was set for $k = 3$ and averaging by inverse distance weighting, also including a prior variable
187 selection based on variance-weighted canonical correlation analysis (CCA). The value of k was
188 kept low due to the small n available since, although a higher k may improve the precision of

189 the estimation, it can also have an averaging effect (i.e., bias extreme values toward the average)
190 (Eskelson et al., 2009; Almeida et al., 2016). The selection was done by recursively restricting
191 the number of predictors (p) from $p = 30$ to $p = 1$, on the grounds of the absolute values of
192 their coefficients in the canonical regression (Cohen et al., 2003; Manzanera et al., 2016). The
193 highest p was intentionally left unrealistically large, given the subsequent low n/p ratio, to test
194 the results that accuracy measures could provide in such an extreme case. An optimal p was
195 selected according to a combination of accuracy measures, which restricted the p on the basis
196 of a hypothesis test (Piñeiro et al., 2008) and avoiding model overfitting (Weisberg, 1985;
197 Hawkins, 2004), as explained below. This same approach for restricting p (see “restricted”
198 alternatives below) was also incorporated to optimize the *best-subset* and *step-wise* variable
199 selection procedures typically used in parametric modelling for remote sensing prediction of
200 *AGB*.

201 *Parametric modelling based on variable selection via step-wise regression* (Weisberg, 1985).
202 A linear model was fitted using a natural logarithm transform of the response variable, as it is
203 typically done in remote sensing predictions of forest attributes (e.g., Naesset, 2002; Hudak,
204 2005; Asner & Mascaro, 2014). Baskerville’s (1972) correction for bias in log-transformed
205 responses was applied taking into account the number of fitted parameters when calculating the
206 standard error of the estimate (Sprugel, 1983). Function “stepAIC” of R was used for applying
207 a backward selection of independent variables in linear regression models (Venables & Ripley,
208 2002). The final p was limited on the basis of the delta parameter (Δ ; Burnham & Anderson,
209 2002), which measured the relative increase in Sugiura’s (1978) corrected AIC (Akaike
210 Information Criterion) at each step (Valbuena et al., 2013b) (hereafter denominated “step-
211 wise”). The result was compared to an alternative incorporating the above-mentioned
212 restrictions – hypothesis test plus avoided overfitting –, which modified the p derived from the
213 step-wise procedure (hereafter denominated “step-wise restricted overfitting”).

214 *Parametric modelling based on variable selection via best-subset regression* (Miller, 1984;
 215 Hudak et al., 2006). This approach also consisted of a linear model with log-transformed
 216 responses and bias correction (Baskerville, 1972; Sprugel, 1983). In this other case, package
 217 “leaps” of R (Lumley & Miller, 2009) was employed for this modelling approach. This
 218 approach exhaustively searches for all variable combinations. The limiting criterion for p was
 219 set to be based on minimization of Mallows’ Cp (Mallows, 1973) (hereafter denominated “best-
 220 subset”). Its result was also compared to a version incorporating the novel restrictions –
 221 hypothesis test plus avoided overfitting – to the best subset procedure for predictor variable
 222 selection (hereafter denominated “best-subset restricted overfitting”).

223 *Statistical measures for accuracy assessment of AGB predictions*

224 Leave-one-out cross-validation was carried out to assess all the prediction methods considered.
 225 Thus, after removing one case (i) from the total n , the remaining were used to calculate a new
 226 *AGB* prediction of the response for that given case (pre_i^{cv}). Hereafter, the superscript/subscript
 227 *cv* is used to distinguish measures calculated after the cross-validation procedure, as opposed
 228 to the superscript/subscript *fit* which will denote non-cross-validated measures, for instance
 229 the predictions that yield model residuals (pre_i^{fit}). The result was evaluated with observed
 230 versus leave-one-out predicted plots, from which we evaluated:

231 (1) The *mean difference* (*MD*) between the predicted minus the observed values:

$$232 \quad MD = \sum_{i=1}^n (pre_i^{cv} - obs_i) / n, \quad (2)$$

233 which evaluates the degree of under- or over-prediction of the method employed. Eq. (2) is
 234 equivalent to the difference between the means of the observed and predicted (e.g.,
 235 McInerney et al., 2010; Wing et al. 2012). *MD* was expressed in *AGB* units, whereas relative
 236 mean difference (*MD%*) was calculated by dividing *MD* by the observed mean *AGB* (\overline{obs}).

237 (2) The *precision* of the prediction, considered as the mean of absolute differences (*MAD*):

$$238 \quad MAD = \sum_{i=1}^n |pre_i^{cv} - obs_i|/n, \quad (3)$$

239 and also the root mean squared differences (*RMSD*) of predicted values with respect to the
240 observed ones:

$$241 \quad RMSD = \sqrt{SS^{cv}/n}, \quad (4)$$

242 where SS^{cv} was the sum of the squared differences between the observed values and the
243 predicted values obtained by cross-validation (a.k.a. predicted sum of squares *PRESS*; Allen,
244 1974; Geisser & Eddy, 1979; Weisberg, 1985: 217; e.g., Valbuena et al., 2013a):

$$245 \quad SS^{cv} = \sum_{i=1}^n (pre_i^{cv} - obs_i)^2. \quad (5)$$

246 Both *MAD* and *RMSD* represent the error in *AGB* units, the latter being more prone to the
247 presence of outliers (e.g., García et al., 2010). Their respective relative counterparts, *MAD%*
248 and *RMSD%* (a.k.a the coefficient of variation of *RMSD*; e.g. Valbuena et al., 2014), were
249 also calculated by dividing them by \overline{obs} .

250 (3) A *hypothesis test* testing whether observed and predicted values follow the 1:1
251 correspondence line (Graybill, 1976; Leite & Oliveira, 2002), was assessed from the
252 intercept (α) and slope (β) of the linear regression model between the observed and predicted
253 (Piñeiro et al., 2008):

$$254 \quad obs_i = \alpha + \beta pre_i^{cv}, \quad (6)$$

255 which is proven by not rejecting the null hypotheses that $H_0: \alpha = 0$ and $H_0: \beta = 1$ for
256 $pre_i^{cv} - obs_i = \alpha + \beta pre_i^{cv}$ (Eq. 9 in Piñeiro et al., 2008). Hence, this is a means for
257 assessing the residual distribution analytically, instead of evaluating it visually from a
258 residuals versus predicted scatterplot (e.g., Mauro et al., 2016: Fig. 2).

259 (4) The *proportions of the total errors* which are due to the unexplained variance (U_{error}),
 260 the slope (U_{slope}), and the bias (U_{bias}), which were evaluated from Theil's (1958) partial
 261 inequality coefficients (Paruelo et al., 1998):

$$262 \quad U_{error} = \sum_{i=1}^n (est_i^{cv} - obs_i)^2 / SS^{cv} , \quad (7)$$

263 where $est_i^{cv} = \hat{\alpha} + \hat{\beta} \cdot pre_i^{cv}$ were the values estimated by the regression model (Eq. 6);

$$264 \quad U_{slope} = [(\beta - 1)^2 \sum_{i=1}^n (pre_i^{cv} - \overline{pre^{cv}})^2] / SS^{cv} ; \quad (8)$$

265 and

$$266 \quad U_{bias} = [n \cdot MD^2] / SS^{cv}. \quad (9)$$

267 We multiplied the values of Theil's (1958) partial inequality coefficients by 100, to make it
 268 straightforward to the reader that they express the percentage of the total error which is due
 269 to either an overall bias of the model (U_{bias}), the presence of trends in the residuals (U_{slope})
 270 or just the residual variance of the model (U_{error}).

271 (5) The *degree of overfitting* to the sample, which we assessed using a replication method
 272 comparing cross-validation results against model residuals (Allen, 1974; Snee, 1977;
 273 Vanclay & Skovsgaard, 1997; Geisser & Eddy, 1979; Hawkins, 2004). Most studies assume
 274 that overfitting is avoided if over-paramaterization of the prediction model is prevented by
 275 using condition number (κ ; Weisberg, 1985; e.g., Naesset, 2002), variance inflation factor
 276 (VIF; Fox & Monette, 1992; e.g., García et al., 2010), Mallows' (1973) Cp statistic (e.g.,
 277 Hudak et al., 2006), or information criterion indices: Akaike (1992) (AIC; e.g., Bright et al.,
 278 2012), Bayesian (BIC; e.g., Wing et al., 2012) or deviance (DIC; e.g., Spriggs et al, 2015).
 279 Many authors deem these insufficient, however, advocating for methods dealing with
 280 overfitting directly (Allen, 1974; Snee, 1977; Hurvich & Tsai, 1989; Rencher & Pun, 1993).

281 Moreover, Hawkins (2004) argued in favour of using replication methods for non-parametric
 282 machine learning approaches like MSN, which may lack the theoretical basis on which κ ,
 283 VIP, Cp or AIC are grounded. For this reason, we alternatively assessed overfitting directly
 284 from the sums of squares ratio (SSR) and R^2 ratio ($R2R$) (Ehrenberg, 1982; Weisberg, 1985:
 285 68-69, 217; Lipovetsky, 2013), both obtained by comparison of a same measure acquired by
 286 model fit against cross-validation.

287 The ratio between the square root of the sums of squares attained in the cross-validation
 288 (SS^{cv}) (Eq. 5) and that using the whole dataset (SS^{fit}) (Snee, 1977; e.g., Valbuena et al.,
 289 2013a) yielded the SSR :

$$290 \quad SSR = \sqrt{SS^{cv}} / \sqrt{SS^{fit}}, \quad (10)$$

291 where SS^{fit} was the sum of squares of the model residuals (j), i.e. the values fitted without
 292 cross-validation (Hawkins, 2004):

$$293 \quad SS^{fit} = \sum_{j=1}^n (pre_j^{fit} - obs_j)^2. \quad (11)$$

294 On the other hand, a similar measure was obtained using the R^2 of the regression of observed
 295 versus predicted values (Piñeiro et al., 2008). This was the ratio between the one obtained
 296 by cross-validation and that from model residuals: the R^2 ratio ($R2R$). Equation (5) derives:

$$297 \quad R_{cv}^2 = 1 - SS^{cv} / SS_{tot}, \quad (12)$$

298 where SS_{tot} was the sum of squared differences of each observation from the overall mean:

$$299 \quad SS_{tot} = \sum_{i=1}^n (obs_i - \overline{obs})^2. \quad (13)$$

300 Whereas from model residuals the coefficient of determination obtained is derived from Eq.
 301 (11) instead:

302 $R_{fit}^2 = 1 - SS^{fit}/SS_{tot}$. (14)

303 Then, the deflation observed by the cross-validation in the coefficient of determination can
 304 be then assessed as (Rencher & Pun, 1993; e.g., Latifi et al., 2015a):

305 $R2R = R_{fit}^2/R_{cv}^2 = \left(1 - \frac{SS^{fit}}{SS_{tot}}\right) / \left(1 - \frac{SS^{cv}}{SS_{tot}}\right)$, (15)

306 Comparing these two functions, Eqs. (10) and (15), it can be seen that *SSR* and *R2R* do, in
 307 essence, very similar tasks. While *R2R* is a ratio of decrease in explained variance
 308 experienced when cross-validating, *SSR* is a ratio of increase in unexplained variance
 309 (square-rooted, in this case). These two measures can therefore be employed to adjust the
 310 inflation of the unexplained variance (*SSR*) or deflation of explained variance (*R2R*) in the
 311 cross-validation to a desirable limit, for example 5% or 10% (Lipovetsky, 2013) (i.e., *SSR*
 312 or *R2R* would be lower than e.g. 1.05 or 1.10 – numerator and denominator in Eq. (15) have
 313 been swapped compared to Eq. (10), so that both *SSR* and *R2R* rise for increasing
 314 overfitting). It may be worthwhile to mention that although in the univariate case the cross-
 315 validation necessarily leads to an increase in the sums of squares and a decrease in the R^2
 316 (Ehrenberg, 1982; Weisberg, 1985), Lipovetsky (2013) showed that this property does not
 317 necessarily always hold in the multivariate case.

318 *Comparing alternatives*

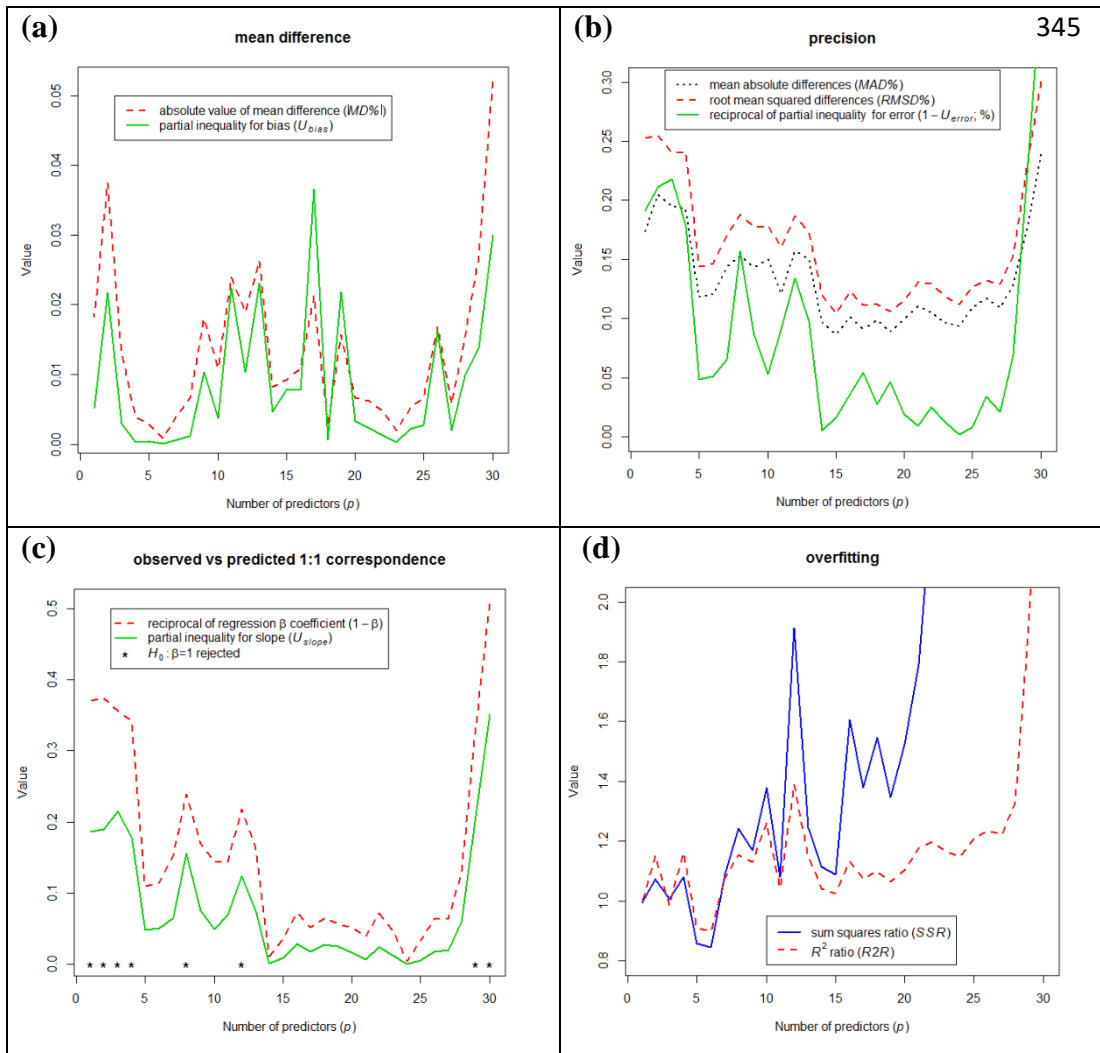
319 The relative merits of each of the proposed statistical measures – *MD*, *MD%*, *MAD*, *MAD%*,
 320 *RMSD*, *RMSD%*, α , β , U_{error} , U_{slope} , U_{bias} , *SSR* and *R2R* – were evaluated by analysing the
 321 results provided when applying different alternative prediction methods to the same dataset,
 322 and also by comparing their corresponding scatterplots of observed versus predicted values.
 323 Firstly, we compared results obtained while increasing the number of predictors in MSN. We
 324 purposely included unrealistically low n/p ratios, with the intention to realize which statistical

325 measures would flag up their unreliability. Additionally, we observed the correlations between
326 pairs of statistical measures to prove whether they are simply redundant or provide additional
327 information, using Spearman's rank correlation coefficient (ρ) because it could prove that two
328 methods would rank alternatives in a similar manner. Secondly, we compared automatic
329 variable selection procedures commonly employed in the assessment of remote sensing assisted
330 *AGB* estimations: step-wise and best subset. The additional statistical measures were
331 incorporated into these algorithms, showing that improvements in overfitting and avoiding
332 systematic errors may be achieved without excessively compromising the overall precision of
333 the estimates.

334 **Results**

335 *Estimation with different number of predictors*

336 Let us first analyse the results observed when modifying the number of predictors p during the
337 variable selection procedure for MSN imputation. **Figure 1** shows the evolution of the statistical
338 measures for increasing p , grouped by the characteristics they describe: mean difference and
339 precision of predictions, their 1:1 correspondence with the observed values, and the degree of
340 overfitting. **Table 1** summarizes the numerical results attained for a relevant selection of these:
341 $p = 2, 3, 5, 8, 10, 15, 20$ and 30 . Their corresponding observed versus predicted plots are
342 shown in **Fig. 2**. Results obtained from the hypothesis tests applied to the fit of observed versus
343 predicted rejected the reliability of accepting the options using either $p = 1-4, 8, 12, 29$ or 30
344 (denoted with asterisks in **Figs. 1c**), whereas every other option passed the test successfully.



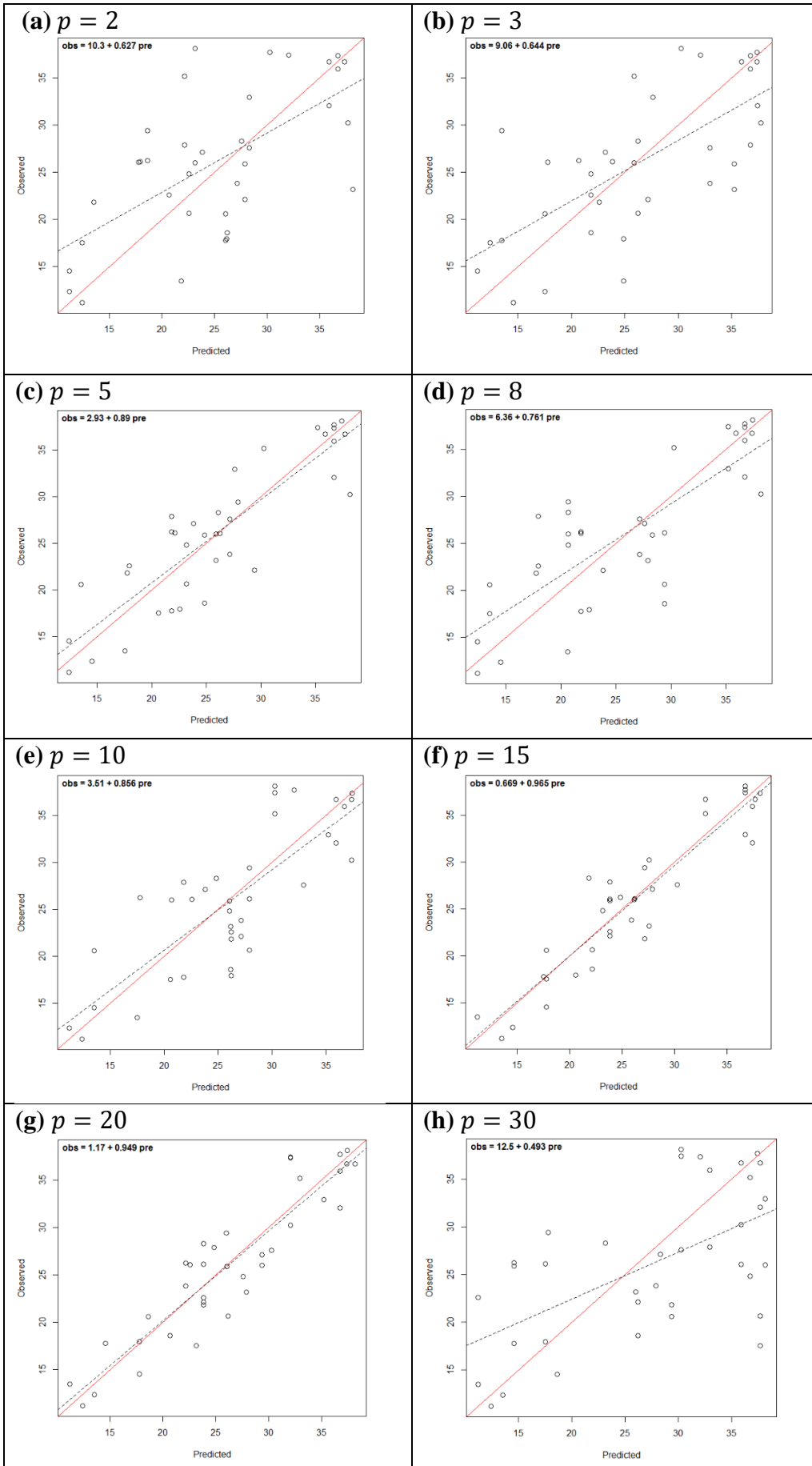
346 **Figure 1.** Statistical evaluation of MSN predictive method for increasing the number of
 347 predictors (p), grouped according to whether they define (a) the mean difference or (b)
 348 precision of predictions, (c) their 1:1 correspondence or (d) the degree of overfitting.

349

350 **Table 1.** Summary diagnosis of most similar neighbour (MSN) predictions for above-ground
 351 biomass (AGB , $Mg \cdot ha^{-1}$) using an increasing number of predictors (p).

		Number of predictors (p)							
		2	3	5	8	10	15	20	30
Prediction	MD	-98	.34	-.08	-.17	.28	.24	.17	1.35
bias	$MD\%$	-3.75	1.31	-.29	-.66	1.09	.93	.66	5.20
Prediction	MAD	5.33	5.08	3.09	4.00	3.93	2.26	2.60	6.23
precision	$MAD\%$	20.5	19.5	11.8	15.5	15.1	8.7	9.60	23.9
	$RMSD$	6.63	6.26	3.75	4.90	4.65	2.72	3.01	7.68
	$RMSD\%$	25.4	24.0	14.4	18.8	17.9	10.4	11.5	30.0
Hypothesis	α	10.3**	9.06**	2.93 ^{NS}	6.36*	3.51 ^{NS}	.66 ^{NS}	1.17 ^{NS}	12.5***
test	β	.63**	.64***	.89 ^{NS}	.76**	.86 ^{NS}	.97 ^{NS}	.95 ^{NS}	.49***
Partial	U_{error} (%)	78.8	78.2	95.1	84.3	94.7	98.3	98.0	61.9
inequality	U_{slope} (%)	18.9	21.5	4.84	15.6	4.9	.88	1.59	35.1
coefficients	U_{bias} (%)	2.16	.38	.01	.01	.02	.79	.40	3.00
Agreement	R_{cv}^2 (%)	38.7	45.8	76.3	64.2	63.7	87.1	84.3	33.1
Overfitting	SSR	1.07	1.01	.89	1.24	1.38	1.09	1.52	7.69
	$R2R$	1.15	.99	.91	1.15	1.26	1.02	1.11	2.97

MD : mean differences (Eq. 2). $MD\%$: relative MD . MAD : mean absolute differences (Eq. 3). $MAD\%$: relative MAD . $RMSD$: root mean squared differences (Eq. 4). $RMSD\%$: relative $RMSD$. α/β : intercept/slope of observed versus predicted regression (Eq. 6) (levels of significance for rejecting H_0 : *:.05; **:.01; ***:.001; ^{NS}: non-significant). $U_{error}/U_{slope}/U_{bias}$: Theil's (1958) partial inequality coefficients for error variance/slope/bias (Eqs. 7-9). R_{cv}^2 : cross-validated coefficient of determination (Eq. 12). SSR : sum of squares ratio (Eq. 10). $R2R$: R^2 ratio (Eq. 15). Relative figures and agreement/inequality coefficients have been multiplied by hundred to yield percentage units.



352 **Figure 2.** Observed versus predicted plots of most similar neighbour (MSN) imputation
353 models for above-ground biomass (*AGB*, Mg·ha⁻¹) using an increasing number of predictors
354 (*p*). The solid red line represents the 1:1 correspondence. Dashed line is the linear regression
355 fit between observed and predicted $obs_i = \alpha + \beta \cdot pre_i$.

356

357 Mean differences, i.e. over- or under-prediction, were negligible in most cases (**Table 1; Fig.**
358 **1a**), usually below $|MD| = 2\%$ (which in practice implies an approximate deviation of 0.5
359 Mg·ha⁻¹). Therefore, in almost every case the prediction methods would yield an unbiased
360 estimation of the mean *AGB* for the population. The absolute value of *MD* has been depicted in
361 **Fig. 1a** in order to express its magnitude regardless of whether it implies under- or over-
362 prediction. The results obtained by $|MD|$ were also corroborated by the low proportions of error
363 due to bias, as shown by its corresponding Theil's partial inequality coefficient (U_{bias}). These
364 two measures were very highly correlated $\rho_{(|MD|, U_{bias})} = 0.94$, and hence reiterative. The
365 largest over-predictions resulted from the MSN model with $p = 30$, which showed a $MD =$
366 5.20% with a proportion of the total error due to bias reaching $U_{bias} = 3.00\%$. Any other
367 alternative $p = 1-29$ could have been deemed as providing a reliable *AGB* prediction. However,
368 scatterplots in **Fig. 2a-b** show examples of some cases where the unreliability of predictions
369 could also have been perceived visually. Alternatively to visual assessment, lack of reliability
370 may also be automatically detected via significance of hypothesis tests (denoted by asterisks in
371 **Fig. 1c**).

372 With regards to the precision of predictions, results were also reasonably acceptable, ranging
373 $RMSD = 10.4-18.8\%$ for $p = 5-28$. Higher ($p = 29-30$) or lower ($p = 1-4$) number of
374 predictors reached larger $RMSD = 24.0-30.0\%$ (**Fig. 1b**). *RMSD* and *MAD* changed very
375 similarly for different *p*, *MAD* being systematically lower than *RMSD*, as it could be expected

376 from Eqs. 3-5. As a result, MSN imputations using $p = 1-4$ seemed apparently better when
 377 evaluated by their $MAD = 17.4-20.5\%$, as compared to observing their higher $RMSD =$
 378 $24.0-25.4\%$. In fact $\rho_{(MAD,RMSD)} = 0.99$, and hence there is no need to report both measures.
 379 Moreover, Theil's partial inequality for error (U_{error}) and the slope of the regression β also
 380 showed similar patterns as $RMSD$, being $\rho_{(RMSD,U_{error})} = -0.92$ and $\rho_{(RMSD,\beta)} = -0.85$.
 381 Significances in the test of lack of fit to the 1:1 correspondence were therefore closely
 382 associated to low precisions in the AGB prediction (**Fig. 2**). The use of β , however, provided
 383 the added value of incorporating a significance test that can be used as an objective threshold
 384 for rejecting excessively low precision in prediction error (denoted with asterisks in **Fig. 1c**).

 385 For assessing the degree of overfitting to the sample, the suggested statistical measures –
 386 SSR and $R2R$ – yielded diverging results for high values of p (**Fig. 1d**). Results in **Table 1** and
 387 **Fig. 1d** revealed that, for many of the alternatives, the ‘real’ (cross-validated) precision
 388 exceeded 10% of model residual variance (denoted by values of SSR or $R2R < 1.1$). Among
 389 all the alternatives considered, only those MSN imputations using $p = 1-7, 11$ and 15 obtained
 390 values of $SSR < 1.1$. Being 10% a fairly acceptable level of divergence, if such criterion is set
 391 in conjunction with the hypothesis tests for rejecting a given AGB estimation, then only the
 392 MSN predictions using $p = 5-7, 11$ and 15 would be acceptable options. On the other hand,
 393 $R2R$ was generally less sensitive to overfitting than SSR (**Table 1**). **Fig. 1d** shows that $R2R$
 394 was critically low at elevated values of p , which is in disagreement with what would have
 395 intuitively be assumed by the subsequent low n/p ratios, whereas SSR unveiled a dramatical
 396 increase in the overfit for most alternatives above $p = 7$. In fact, SSR correlated to the p itself
 397 – $\rho_{(p,SSR)} = 0.93$ –, while $R2R$ has a weaker relationship to the number of predictors used –
 398 $\rho_{(p,R2R)} = 0.62$ –, which indicates that comparing the deflation in R^2 may be not useful to
 399 avoid over-parameterization.

400 In light of our results, Theil's partial inequality coefficients can be useful for a detailed
401 evaluation of results. U_{bias} may detect systematic differences between observed and predicted
402 values. Additionally, large values for U_{slope} , such as those obtained for $p = 3$ or $p = 8-10$,
403 indicated a tendency for predicting towards the average *AGB* (**Fig. 2**) (i.e., over-predicting low
404 *AGB* areas and under-predicting large ones). Hence, even if the overall population mean may
405 be assumed unbiased in light of MD or U_{bias} , there is still a chance for the values shown at the
406 scale of the estimation units (the pixels in the remote sensing case) to be selectively under- or
407 over-predicted for certain values within the range of observed *AGB*. Our results showed that
408 this was indeed the case, since large values of $U_{slope} = 10.4-11.3\%$ were associated to
409 significant test results for either the α or β coefficient, or both (**Table 1**). On the other hand,
410 the alternatives for which the null hypotheses were not rejected by the tests (signified by non-
411 significances for the coefficients) obtained much lower values, such as $U_{slope} = 7.09\%$ for $p =$
412 5 and $U_{slope} = 0.51-4.22\%$ for $p = 15-25$. For instance, Theil's partial inequality coefficients
413 were particularly relevant for $p = 3$, (**Fig. 2b**), since its $U_{slope} = 21.5\%$ revealed and averaging
414 effect which remained concealed by its low $MD = 1.31\%$ (**Table 1**).

415 *Comparison of alternative modelling methods*

416 We also wanted to use the proposed measurements of accuracy to compare the results obtained
417 by the MSN imputation with two other modelling alternatives commonly employed in remote
418 sensing-assisted predictions of *AGB*: best-subset and step-wise regression (**Table 2**). Based on
419 the results detailed on the previous sub-section, we decided to incorporate two additional
420 constraints on variable selection (called 'restricted' in **Table 2**) on top of their original
421 limitation criteria (i.e., C_p for best-subset and Δ for step-wise). These were the hypothesis tests
422 and the degree of overfitting, i.e. a model would be declined if either of the null hypotheses H_0 :
423 $\alpha = 0$ or $H_0: \beta = 1$ were rejected, or $SSR > 1.1$. **Table 2** compares all these versions against

424 the previously-selected MSN imputation model for $p = 5$, which was selected as the optimal
425 MSN predictions under the same criteria. **Figure 3** shows the observed versus predicted plots
426 corresponding to each of these alternatives.

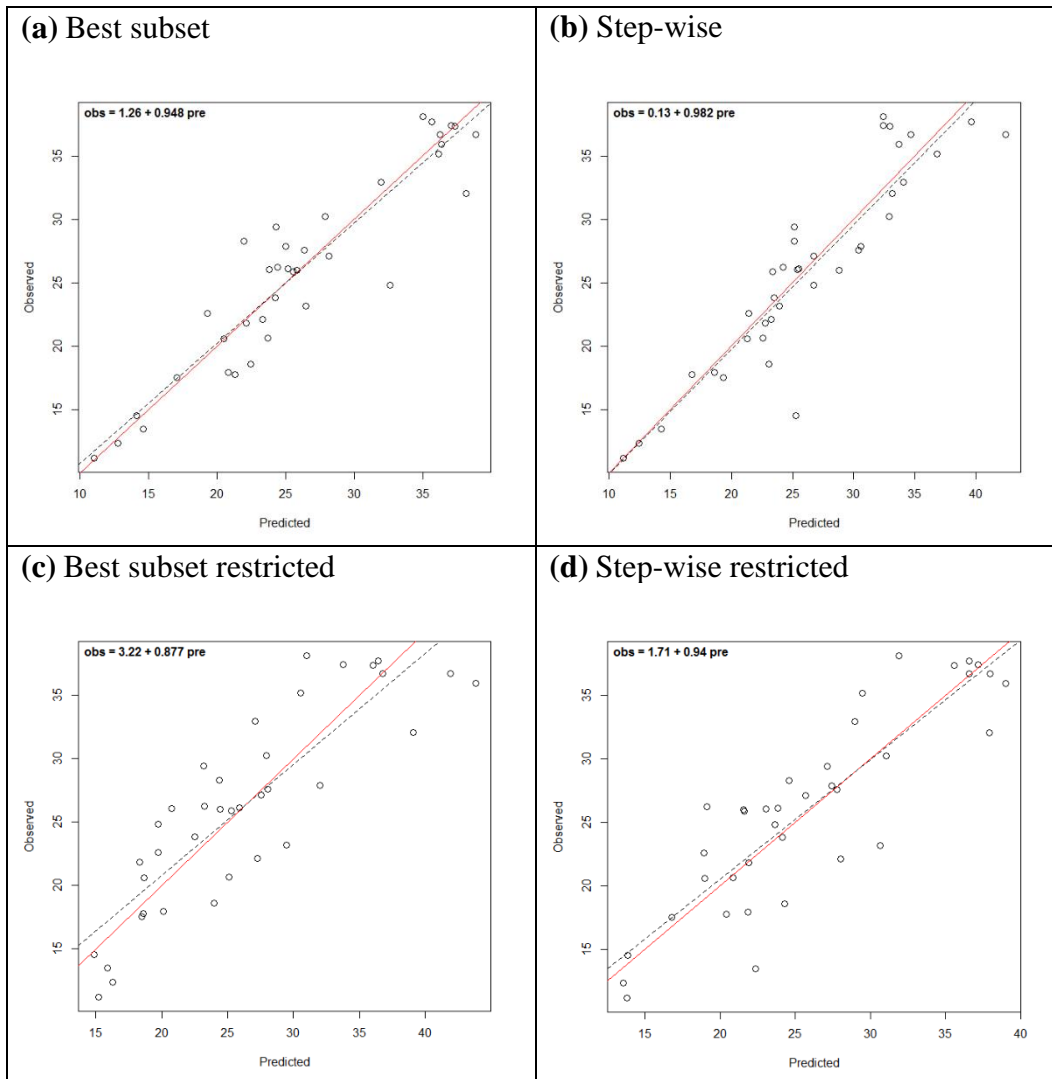
427 In a similar manner to the previous comparison of MSN imputation, all the alternatives resulted
428 in unbiased predictions of population mean ($|MD| = 0.10-1.32\%$ and $U_{bias} = 0.01-0.79\%$)
429 performing a reasonable error variance ($RMSD = 9.67-15.3\%$) and good agreement between
430 observed and predicted (see Valbuena et al., 2018). In this case they all passed the hypothesis
431 tests, as most of the overall errors were simply due to unsystematic sources affecting the error
432 variance of the model itself ($U_{errors} = 94.7-98.9\%$). None of the models therefore had to be
433 declined due to failing the hypothesis test on the correspondence between observed and
434 predicted. However, we detected an overfitting effect at both the best-subset model selected on
435 the grounds of Mallows' C_p ($SSR = 1.28$) and also at the step-wise regression selected via Δ 's
436 difference in Sugiura's corrected AIC ($SSR = 2.90$). It is noteworthy to point out that this
437 contingency could have simply remained overlooked if overfitting had been analysed according
438 to the deflation in R^2 , which was only $R2R = 1.04$ for best subset and $R2R = 1.17$ for the
439 step-wise regression model. Accordingly, we imposed the criterion of $SSR \leq 1.1$ to further
440 constrain the prediction dataset of these models. This resulted in unbiased models including just
441 $p = 2$ independent variables, which avoided overfitting ($SSR = 1.08$) while not excessively
442 compromising model precision ($RMSD = 15.3\%$ and $RMSD = 14.3\%$, respectively).

Table 2. Comparison of diagnoses for different prediction method and variable selection alternatives to obtain above-ground biomass (*AGB*, Mg·ha⁻¹) predictions.

		Best-subset	Best-subset restricted	Step-wise	Step-wise restricted	MSN restricted
Number of predictors (<i>p</i>)		8	2	23	2	5
Prediction bias	<i>MD</i>	-.06	-.02	.34	-.16	-.08
	<i>MD%</i>	-.24	-.10	1.32	-.63	-.29
Prediction precision	<i>MAD</i>	2.09	3.32	2.26	2.87	3.09
	<i>MAD%</i>	8.01	12.7	8.69	11.0	11.8
	<i>RMSD</i>	2.52	3.99	3.07	3.73	3.75
	<i>RMSD%</i>	9.67	15.3	11.8	14.3	14.4
Hypothesis test	α	.96 ^{NS}	3.21 ^{NS}	.13 ^{NS}	1.71 ^{NS}	2.93 ^{NS}
	β	.97 ^{NS}	.87 ^{NS}	.98 ^{NS}	.94 ^{NS}	.89 ^{NS}
Partial inequality coefficients	<i>U_{error}</i> (%)	98.9	94.7	98.6	98.5	95.1
	<i>U_{slope}</i> (%)	1.04	5.30	.17	1.30	4.84
	<i>U_{bias}</i> (%)	.06	.01	.01	.19	.01
Agreement	<i>R_{cv}²</i> (%)	88.9	73.4	83.6	75.7	76.3
Overfitting	<i>SSR</i>	1.28	1.08	2.90	1.08	.89
	<i>R2R</i>	1.04	1.05	1.17	1.04	.91

MD: mean differences (Eq. 2). *MD%*: relative *MD*. *MAD*: mean absolute differences (Eq. 3). *MAD%*: relative *MAD*. *RMSD*: root mean squared differences (Eq. 4). *RMSD%*: relative *RMSD*. α/β : intercept/slope of observed versus predicted regression (Eq. 6) (levels of significance for rejecting H₀: *.05; **:.01; ***:.001; ^{NS}: non-significant).

U_{error}/*U_{slope}*/*U_{bias}*: Theil's (1958) partial inequality coefficients for error variance/slope/bias (Eqs. 7-9). *R_{cv}²*: cross-validated coefficient of determination (Eq. 12). *SSR*: sum of squares ratio (Eq. 10). *R2R*: *R*² ratio (Eq. 15). Relative figures and agreement/inequality coefficients have been multiplied by hundred to yield percentage units.



443 **Figure 3.** Observed versus predicted plots of different modelling and variable selection

444 alternatives to obtain above-ground biomass (AGB , $Mg \cdot ha^{-1}$) predictions. The solid red line

445 represents the 1:1 correspondence. Dashed line is the linear regression fit between observed

446 and predicted $obs_i = \alpha + \beta \cdot pre_i$.

447 **Discussion**

448 *Importance of adding complementary analyses for assessing the accuracy of models*

449 The most important implication of the present results is that most of the alternatives contrasted

450 could have been reasonably judged as reliable when observing only statistical descriptors for

451 mean difference, precision and agreement. These three types of statistics are the ones most

452 commonly employed for assessing accuracy in this field (e.g., Zhao et al., 2009; Erdody &

453 Moskal 2010; McInerney et al., 2010; d'Oliveira et al 2012; Chen & Zhu, 2013; Straub et al.,
454 2013; Valbuena et al., 2014). In our analysis, by looking only at $MD\%$, $RMSD\%$ and R_{cv}^2 , and
455 also most scatterplots in **Figs. 2-3**, it could be rationally deduced that any option including a
456 MSN imputation with $p = 1-28$ would yield reliable accuracies, including the best-subset and
457 step-wise models as well. The suggested complementary analyses however, showed that many
458 more of the presented alternatives for *AGB* prediction should in fact be discarded.

459 Significances in the hypothesis tests suggested by Piñeiro et al. (2008) demonstrated that MSN
460 imputations using $p = 3$ or $p = 8$ gave an insufficient fit between observed and predicted
461 values. This diagnosis may have been difficult to make by merely observing the scatterplots
462 (**Figs. 2b,d**). Although testing the regression of observed versus predicted values is a well-
463 settled practice in ecological modelling (Graybill, 1976; Reynolds & Chung, 1986; Leite &
464 Oliveira, 2002; Piñeiro et al., 2008), to our knowledge, hypothesis tests have never before been
465 included in the evaluation of forest *AGB* using remote sensing, and they have seemingly been
466 simply overlooked. The results presented in this article suggest that there may be a need to
467 include them in future accuracy assessment procedures in this field as well. Furthermore, we
468 also wish to seek consensus and promote the arguments advanced by Piñeiro et al. (2008) in
469 favour using observed (on the y-axis) versus predicted (on the x-axis) – and not predicted versus
470 observed (e.g., McRoberts et al., 2002; Holmgren et al., 2008; Zhao et al., 2009; McInerney et
471 al., 2010; Chen & Zhu 2013; Valbuena et al., 2014) – for reporting the accuracy of remote
472 sensing-assisted *AGB* estimates. Piñeiro et al. (2008) showed that such distinction matters since
473 it may change the result and conclusions of model evaluation.

474 Regarding the overfitting tests based on cross-validation (Allen, 1974; Snee, 1977; Geisser &
475 Eddy, 1979; Weisberg, 1985; Hawkins, 2004), we wish to emphasize that *SSR* succeeded in
476 revealing both the best-subset and the step-wise models initially considered, and also any MSN

477 imputation using $p \geq 8$, as being unreliably overfitted to the sample and therefore hardly
478 generalizable. The described step-wise and best subset approaches to variable selection are very
479 frequently employed in remote sensing-assisted estimations of forest attributes (e.g., Naesset,
480 2002; Hudak et al., 2006; Wing et al., 2012; Straub et al., 2013; Estornell, et al. 2014). We
481 therefore suggest that accuracy assessment procedures for *AGB* predictions obtained from
482 remote sensing should be improved by using hypothesis testing and overfitting evaluation.

483 *Unveiling averaging effects: unbiased means, and yet over/under-predicting*

484 Even having an unbiased prediction method and a robust sampling design, the outcome is still
485 susceptible to under- and over-prediction within specific ranges of *AGB* values. Sometimes the
486 discrepancy between observed and predicted is due to an averaging effect, which in practice
487 translates into an underestimation of large *AGB* values and an overestimation at areas of lesser
488 *AGB*, which in turn may remain concealed if only observing the bias of the population mean.
489 The averaging effect is a typical and intrinsic weakness of nearest neighbours methods (Franco-
490 Lopez et al., 2001; McInerney et al., 2010). It is caused by the lack of available neighbours
491 beyond the limits of the observed *AGB* range, hence tending to shift the predictions towards the
492 average for values located in the borderline of that range. This effect therefore becomes more
493 accentuated as the n/p ratio decreases (McRoberts et al., 2002). Our results indicate this
494 shortcoming may be detected with the assistance of hypothesis tests suggested by Piñeiro et al.
495 (2008) and Theil's (1958) partial inequality coefficients (Smith & Rose, 1995; Paruelo et al.,
496 1998). Taking our results and as a rule of thumb, we would suggest that the proportions of error
497 due to causes other than the residual variance must not exceed the thresholds $U_{slope} \leq 10\%$
498 and $U_{bias} \leq 1\%$, and in general the model error itself should be no lesser than $U_{error} \geq 90\%$.
499 Under-prediction in areas of large *AGB* is a common problem in remote sensing assessments
500 (e.g., Bright et al., 2012; Asner & Mascaro, 2014), and these areas are of very high importance

501 for the purposes of the inventory. To our knowledge, however, these coefficients have not been
502 employed in the context of remote sensing estimates of forest characteristics before, and only
503 García et al. (2010) resolved the *RMSD* into systematic and unsystematic portions. U_{slope} could
504 still be useful for identifying these averaging effects, as it was revealed in our results for MSN
505 imputations using $p = 3$ or $p = 8$ (**Table 2**), where averaging effects were indeed undergoing
506 (**Fig. 2b,d**). This flaw was also detected by significant results in the hypothesis tests. Therefore,
507 averaging effects may be detected by either large values of U_{slope} , or via interpretation of α or
508 β coefficients. When statistical significance proves $\alpha \neq 0$ but cannot reject $\beta = 1$, it is an
509 indication for a source of systematic under- or over-prediction along the full *AGB* range. If $\alpha =$
510 0 cannot be rejected but $\beta \neq 1$ significantly, the under-prediction is concentrated in values of
511 large *AGB* only, for instance due to saturation of the remote sensor. A combination of $\alpha = 0$
512 and $\beta \neq 1$ may as well indicate an over-prediction for small *AGB* values. If both null
513 hypotheses are rejected and we accept $\alpha \neq 0$ and $\beta \neq 1$, then we are detecting an averaging
514 effect whenever $\beta < 1$, as was the case in many of the results presented in this study.

515 *Overfitting to the field sample training the prediction method*

516 We also detected potential problems of overfitting in some of the alternatives proposed. Such
517 contingency would in practice have a harmful effect when applying the resulting fit to the
518 predictor variables to obtain *AGB* maps. It is noteworthy that the added value of remote sensing,
519 compared to traditional design-based sampling using field plots only, is on the capacity to
520 provide *AGB* predictions throughout large inaccessible forest areas (Naesset, 2002; McRoberts
521 et al., 2013; Asner & Mascaro, 2014; Chen et al., 2015; Mauro et al., 2016). This advantage is
522 therefore lost if overfitting to the sample renders *AGB* predictions unreliable at the pixel scale,
523 even if the population mean estimate is unbiased. We therefore suggest the inclusion of

524 overfitting measures in addition to those already widespread: mean difference, precision and
525 agreement.

526 The true degree of overfitting will always remain elusive unless an external validation using an
527 independent field *AGB* dataset is carried out (Allen, 1974; Snee, 1977; Geisser & Eddy, 1979;
528 Hawkins, 2004). However, even in the event of having the opportunity to acquire a large enough
529 number of plots from the field, modellers would face trade-offs between the advantages
530 separating a subset for validation of the main dataset and the gain in incorporating them to the
531 model for increasing its degrees of freedom, strengthening the certainty of the relationships
532 found, and the power of their statistical inference (Cohen et al., 2003). As an alternative, the
533 cross-validation approach seems to provide a good indicative proxy for assessing overfitting
534 (Weisberg, 1985; Rencher & Pun, 1993; Vanclay & Skovsgaard, 1997; Hawkins, 2004). *SSR*
535 succeeded in identifying risk of overfitting for some of the alternatives in **Tables 1 and 2** that
536 could have otherwise remained undetected. For this reason, we suggest that *SSR* may provide
537 a useful indication that a given predictive method may undergo overfitting effects. For predictor
538 variable selection purposes, a desirable limit for model rejection may be chosen, as for instance
539 we suggested to limit $SSR \leq 1.1$. It is worth emphasizing that such limit should also be
540 employed in combination with the suggested hypothesis test, since otherwise *MSN* imputations
541 using $p = 2-3$ would have been deemed reliable if judged on the basis of *SSR* only (**Table 1**).
542 Surprisingly, decreasing p did not univocally lead to a decrease in *SSR* and *R2R*, and hence it
543 may be as detrimental to have either too few or too many predictors. The key question is
544 possibly to include in the model only non-collinear predictors which truly add separate portions
545 of explained variance in the observed *AGB* (Ehrenberg, 1982; Weisberg, 1985).

546 Regarding the choice of either *SSR* or *R2R* for assessing overfitting, our results showed
547 unexpected differences which may in practice be critical. **Fig. 1d** demonstrated that the values

548 obtained by *SSR* or *R2R* diverged from $p \geq 8$. As a result, *R2R* was too low at high values of
549 p , which in practice would imply insufficiently low n/p ratios, and therefore the reliability of
550 *R2R* as a measure of overfitting is questionable. We therefore suggest that evaluating the
551 inflation in the sums of squares of errors (Weisberg, 1985; e.g., Valbuena et al., 2013; Almeida
552 et al., 2016) may be a more sensible approach to assessing overfitting than analysing the
553 deflation in R^2 (Rencher & Pun, 1993; e.g., Latifi et al., 2015a).

554 Most studies assume that avoiding over-parametrized models via κ , VIP, Cp or AIC is sufficient
555 to avoid overfitting (e.g., Naeset, 2002; Hudak et al., 2006; Erdody & Moskal, 2010; García
556 et al., 2010; Bright et al., 2012; Wing et al., 2012; Latifi et al., 2015a; Spriggs et al., 2015).
557 Many of these indices, however, have been suspected in some occasions of being insufficient
558 to avoid model overfitting (Hurvich & Tsai, 1989; Rencher & Pun, 1993; Vanclay &
559 Skovsgaard, 1997). In the present research we also detected the need for incorporating further
560 restrictions to Cp and AIC (**Table 2**). As an alternative, Weisberg (1985) and Hawkins (2004)
561 recommended using cross-validation to prevent overfitting. Our results suggest that, while
562 model precision was not excessively compromised, the assessment of overfitting presented an
563 opportunity for increased reliability of remote sensing predictions of *AGB* (Franco-Lopez et al.,
564 2001; Valbuena et al., 2013b; Latifi et al., 2015a). We therefore suggest that in addition to the
565 use of κ , VIP, Cp or AIC, a specific measure devoted to evaluate the degree of overfitting, such
566 as *SSR*, should become a general requirement

567 **Conclusions**

568 Given the results presented in these comparisons, we wish to put forward a suggestion to
569 perform a more thorough analysis of accuracy in ecological models, which in particular we
570 wish to address to authors carrying out remote sensing-assisted predictions of *AGB*. We may
571 draw four main conclusions from the discussion of our results (plus an additional one, see

572 Valbuena et al., 2018). (1) By simply looking at the most common measures of accuracy
573 assessment – mean difference, precision and agreement –, there is a risk of interpreting as
574 reliable *AGB* predictions which are in fact unreliable. *MD*, *RMSD* and R^2 are useful statistics
575 for accuracy assessment, but perhaps not sufficient for truly evaluating the convenience of a
576 given prediction alternative. (2) Piñeiro et al.’s (2008) hypothesis tests were clearly useful in
577 providing objective means for inferring the statistical significance of the agreement between
578 observed and predicted values, which would otherwise be difficult to grasp just by visual
579 diagnosis of scatterplots. (3) Theil’s partial inequality coefficients can be useful for diagnosis
580 of the causes leading to disagreement, detecting averaging effects or other types of under- or
581 over-predictions occurring at specific ranges of *AGB*. (4) We also observed that overfitting
582 effects may remain concealed unless specifically addressed. When comparing the evaluation of
583 inflation in sums of squares versus deflation of R^2 , our results suggested the former to be a
584 more advantageous approach. We therefore recommend researchers to incorporate the
585 presented statistical measures for (2), (3) and (4) in their own accuracy assessment protocols.
586 This recommendation may, of course, be extended to other fields of applied ecological
587 modelling as well.

588 **Acknowledgments**

589 This work was partially supported by the Spanish Directorate General for Scientific and
590 Technical Research under Grant CGL2013-46387-C2-2-R. We also thank the Valsain Forest
591 Centre, of the National Park Body (Spain), for their valuable help. Dr Valbuena and Prof.
592 Coomes work is supported by an EU Horizon 2020 Marie Skłodowska-Curie Action entitled
593 “Classification of forest structural types with LIDAR remote sensing applied to study tree size-
594 density scaling theories” (LORENZLIDAR-658180). Danilo Almeida acknowledges support
595 from São Paulo Research Foundation (FAPESP) (grant 2016/05219-9). The authors are grateful
596 for the comments received from the editor and reviewers.

597 **References**

598 Allen DM (1974) The relationship between variable selection and data augmentation and a
599 method for prediction. *Technometrics* 16: 125-127

600 Almeida DRA, Nelson BW, Schietti J, Gorgens EB, Resende AF, Stark SC, Valbuena R. (2016)
601 Contrasting fire damage and fire susceptibility between seasonally flooded forest and upland
602 forest in the Central Amazon using portable profiling LiDAR. *Remote Sensing of Environment*
603 184: 153-160

604 Akaike H (1992) Information theory and an extension of the maximum likelihood principle. In:
605 Kotz S, Johnson NL (Eds) *Breakthroughs in statistics* 1. Springer, London, pp 610–624

606 Asner GP, Mascaro J (2014). Mapping tropical forest carbon: Calibrating plot estimates to a
607 simple LiDAR metric. *Remote Sensing of Environment* 140: 614-624

608 Axelsson P (2000) DEM generation from laser scanner data using adaptive TIN models.
609 *International Archives of Photogrammetry and Remote Sensing* 33, Part B4: 110-117

610 Baskerville G (1972) Use of logarithmic regression in the estimation of plant biomass.
611 *Canadian Journal of Forest Research* 2: 49-53

612 Bright BC, Hicke JA, Hudak AT (2012) Estimating aboveground carbon stocks of a forest
613 affected by mountain pine beetle in Idaho using lidar and multispectral imagery. *Remote*
614 *Sensing of Environment* 124: 270–281

615 Burnham KP, Anderson DR (2002) *Model selection and multimodel inference: a practical*
616 *information-theoretic approach* (2nd Edition). Secaucus, NJ, USA: Springer

617 Chen Y, Zhu X (2013) An integrated GIS tool for automatic forest inventory estimates of *Pinus*
618 *radiata* from LiDAR data. *GIScience & Remote Sensing* 50: 667-689

619 Chen Q, Vaglio Laurin G, Valentini R (2015) Uncertainty of remotely sensed aboveground
620 biomass over an African tropical forest: Propagating errors from trees to plots to pixels. *Remote*
621 *Sensing of Environment* 160: 134-143

622 Cohen WB, Maersperger TK, Gower ST, Turner DP (2003) An improved strategy for
623 regression of biophysical variables and Landsat ETM+. *Remote Sensing of Environment* 84:
624 561-571

625 Crookston NL, Finley AO (2007) yaImpute: An R package for k NN imputation. *Journal of*
626 *Statistical Software* 23 (10): 1-16

627 d'Oliveira MVN, Reutebuch SE, McGaughey RJ, Andersen HE (2012) Estimating forest
628 biomass and identifying low-intensity logging areas using airborne scanning lidar in Antimary
629 State Forest, Acre State, Western Brazilian Amazon. *Remote Sensing of Environment* 124: 479-
630 491

631 Ehrenberg ASC (1982). How good is best? *Journal of the Royal Statistical Society, Series A*
632 145: 364-366

633 Erdody TL, Moskal LM (2010) Fusion of LiDAR and imagery for estimating forest canopy
634 fuels. *Remote Sensing of Environment* 114, 725–737

635 Estornell J, Velázquez-Martí B, López-Cortés I, Salazar D, Fernández-Sarría A (2014)
636 Estimation of wood volume and height of olive tree plantations using airborne discrete-return
637 LiDAR data. *GIScience & Remote Sensing* 51, 17-29.

638 Fox DG (1981) Judging air quality model performance. *Bulletin of the American Meteorological*
639 *Society* 62: 599–609

640 Fox J, Monette G (1992). Generalized collinearity diagnostics. *Journal of the American*
641 *Statistical Association* 87: 178-183.

642 Franco-Lopez H, Ek AR, Bauer ME (2001) Estimation and mapping of forest stand density,
643 volume, and cover type using the k-nearest neighbors method. *Remote Sensing of Environment*
644 77(3): 251-274

645 Freese F (1960) Testing accuracy. *Forest Science* 6: 139-145.

646 Eskelson, B.N.I., Temesgen, H., Lemay, V., Barrett, T.M., Crookston, N.L. & Hudak, A.T.
647 2009. The roles of nearest neighbor methods in imputing missing data in forest inventory and
648 monitoring databases. *Scandinavian Journal of Forest Research* 24: 235–246.

649 García M, Riaño D, Chuvieco E, Danson FM (2010) Estimating biomass carbon stocks for a
650 Mediterranean forest in central Spain using LiDAR height and intensity data. *Remote Sensing*
651 *of Environment* 115: 1369-1379

652 Geisser S, Eddy W (1979) A Predictive Approach to Model Selection. *Journal of the American*
653 *Statistical Association* 74 (365): 153-160.

654 Graybill FA (1976) *Theory and Application of the Linear Model*. Duxbury Press: Belmont, CA.

655 Hawkins DM (2004) The problem of overfitting. *Journal of chemical information and*
656 *computer sciences* 44.1: 1-12.

657 Holmgren J, Persson A, Soderman U (2008). Species identification of individual trees by
658 combining high resolution LIDAR data with multi-spectral images. *International Journal of*
659 *Remote Sensing* 29: 1537–1552

660 Hudak AT, Crookston NL, Evans JS, Falkowski MJ, Smith AMS, Gessler PE, Paul E, Morgan
661 P (2006) Regression modeling and mapping of coniferous forest basal area and tree density

662 from discrete-return lidar and multispectral satellite data. *Canadian Journal of Remote Sensing*
663 32: 126-138

664 Hudak AT, Crookston NL, Evans JS, Hall DE, Falkowski MJ (2008) Nearest neighbor
665 imputation of species-level, plot-scale forest structure attributes from LiDAR data. *Remote*
666 *Sensing of Environment* 112 (5): 2232–2245

667 Hurvich CM, Tsai CL (1989) Regression and time series model selection in small samples.
668 *Biometrika* 76: 297–307

669 Latifi H, Heurich M, Hartig F, Müller J, Krzystek P, Jehl H, Dech S (2015a) Estimating over-
670 and understorey canopy density of temperate mixed stands by airborne LiDAR data. *Forestry*
671 89 (1): 69-81.

672 Latifi H, Fassnacht FE, Müller J, Tharani A, Dech S, Heurich M (2015b) Forest inventories by
673 LiDAR data: A comparison of single tree segmentation and metric-based methods for
674 inventories of a heterogeneous temperate forest. *International Journal of Applied Earth*
675 *Observation and Geoinformation* 42, 162-174

676 Leggett RW, Williams LR (1981) A reliability index for models. *Ecological Modelling* 13:
677 303–312.

678 Leite HG, Oliveira FHT (2002) Statistical procedure to test identity between analytical
679 methods. *Communications in Soil Science and Plant Analysis* 33 (7-8): 1105-1118

680 Lipovetsky S (2013) How good is best? Multivariate case of Ehrenberg-Weisberg analysis of
681 residual errors in competing regressions. *Journal of Modern Applied Statistical Methods* 12 (2):
682 14

683 Lumley T, Miller A (2009) *leaps: regression subset selection*. R package version 2.9.
684 <https://CRAN.R-project.org/package=leaps>

685 Mallows CL (1973) Some Comments on Cp. *Technometrics* 15 (4): 661–675

686 Manzanera JA, Garcia-Abril A, Pascual C, Tejera R, Martín Fernández S, Tokola T, Valbuena
687 R (2016) Fusion of airborne LiDAR and multispectral sensors reveals synergic capabilities in
688 forest structure characterization. *GIScience & Remote Sensing* 53: 723-738

689 Mauro F, Molina I, Garcia-Abril A, Valbuena R, Ayuga-Téllez E (2016) Remote sensing
690 estimates and measures of uncertainty for forest variables at different aggregation levels.
691 *Environmetrics* 27(4): 225-238

692 McGaughey RJ (2012). *FUSION/LDV: Software for LIDAR data analysis and visualization*.
693 Version 3.10. USDA Forest Service. Seattle, Washington, USA.

694 McInerney DO, Suárez J, Valbuena R, Nieuwenhuis M (2010). Forest canopy height retrieval
695 using Lidar data, medium-resolution satellite imagery and kNN estimation in Aberfoyle,
696 Scotland. *Forestry* 83(2): 195-206

697 McRoberts RE, Nelson MD, Wendt DG (2002) Stratified estimation of forest area using
698 satellite imagery, inventory data, and the k-nearest neighbors technique. *Remote Sensing of*
699 *Environment* 82(2-3): 457-468

700 McRoberts RE, Naesset E, Gobakken T (2013) Accuracy and precision for remote sensing
701 applications of nonlinear model-based inference. *IEEE Journal of Selected Topics in Applied*
702 *Earth Observations and Remote Sensing* 6 (1): 27-34

703 Miller A (1984) Selection of subsets of regression variables. *Journal of the Royal Statistical*
704 *Society, Series A* 147: 389-425

705 Moeur M, & Stage AR (1995) Most similar neighbor: an improved sampling inference
706 procedure for natural resource planning. *Forest Science* 41 (2): 337-359.

707 Montero G, Ruiz-Peinado R, Muñoz M. (2005) *Producción de biomasa y fijación de CO2 por*
708 *los bosques españoles*. Monografías Instituto Nacional de Investigación y Tecnología Agraria
709 y Alimentaria, Serie Forestal, Madrid, Spain (in Spanish).

710 Naesset E (2002) Predicting forest stand characteristics with airborne scanning laser using a
711 practical two-stage procedure and field data. *Remote Sensing of Environment* 80(1): 88-99

712 Packalén P, Maltamo M (2008) Estimation of species-specific diameter distributions using
713 airborne laser scanning and aerial photographs. *Canadian Journal of Forest Research* 38(7):
714 1750-1760

715 Paruelo JM, Jobbágy EG, Sala OE, Lauenroth WK, Burke I (1998) Functional and structural
716 convergence of temperate grassland and shrubland ecosystems. *Ecological Applications* 8(1):
717 194–206

718 Piñeiro G, Perelman S, Guerschman JP, Paruelo JM (2008) How to evaluate models: observed
719 vs. predicted or predicted vs. observed? *Ecological Modelling* 216(3): 316–322

720 R Development Core Team (2016) *R: a language and environment for statistical computing*.

721 Rencher AC, Pun CP (1993) Inflation of R^2 in best subset regression. *Technometrics* 22: 49-
722 53

723 Reynolds MR, Chung J (1986) Regression methodology for estimating model prediction
724 error. *Canadian Journal of Forest Research*, 16 (5): 931-938

725 Rouse JW, Haas RH, Scheel JA, Deering DW (1974) Monitoring vegetation systems in the
726 great plains with ERTS. *Proceedings, 3rd Earth Resource Technology Satellite (ERTS)*
727 *Symposium 1*: 48-62

728 Särndal CE, Swensson B, Wretman J (1992) *Model assisted survey sampling*. Springer-Verlag,
729 Inc. New York

730 Smith EP, Rose KR (1995) Model goodness of fit using regression and related techniques.
731 *Ecological Modelling 77*: 49-64

732 Snee R (1977) Validation of Regression Models: Methods and Examples. *Technometrics 19*(4):
733 415-428.

734 Spriggs R, Vanderwel M, Jones T, Caspersen J, Coomes D (2015) A simple area-based model
735 for predicting airborne LiDAR first returns from stem diameter distributions: an example study
736 in an uneven aged, mixed temperate forest. *Canadian Journal of Forest Research 45* (10): 1338-
737 1350

738 Sprugel DG. (1983) Correcting for bias in log-transformed allometric equations. *Ecology 64*:
739 209-210

740 Straub C, Tian J, Seitz R, Reinartz P (2013) Assessment of Cartosat-1 and WorldView-2 stereo
741 imagery in combination with a LiDAR-DTM for timber volume estimation in a highly
742 structured forest in Germany. *Forestry 86* (4): 463-473

743 Sugiura N (1978) Further analysts of the data by Akaike's information criterion and the finite
744 corrections. *Communications in Statistics - Theory and Methods 7* (1): 13-26.

745 Tedeschi LO (2006) Assessment of the adequacy of mathematical models. *Agricultural Systems*
746 89 (2-3): 225-247

- 747 Theil H (1958) *Economic Forecasts and Policy*. Amsterdam: North Holland
- 748 Valbuena R (2014) Integrating airborne laser scanning with data from global navigation satellite
749 systems and optical sensors. In: Maltamo M, Næsset E, Vauhkonen J (Eds.) *Forestry*
750 *applications of airborne laser scanning. Concepts and case studies*. Managing Forest
751 Ecosystems Series 27. Springer, Dordrecht. pp. 63-88
- 752 Valbuena R, Mauro F, Arjonilla F, Manzanera JA (2011) Comparing airborne laser scanning-
753 imagery fusion methods based on geometric accuracy in forested areas. *Remote Sensing of*
754 *Environment* 115: 1942–1954
- 755 Valbuena R, Mauro F, Rodriguez-Solano R, Manzanera JA (2012) Partial least squares for
756 discriminating variance components in global navigation satellite systems accuracy obtained
757 under Scots pine canopies. *Forest Science* 582: 139-153
- 758 Valbuena R, Maltamo M, Martín-Fernández S, Packalen P, Pascual C, Nabuurs GJ (2013a)
759 Patterns of covariance between airborne laser scanning metrics and Lorenz curve descriptors of
760 tree size inequality. *Canadian Journal of Remote Sensing* 39 (S1), S18–S31
- 761 Valbuena R, Packalen P, Mehtätalo L, Garcia-Abril A, Maltamo M (2013b) Characterizing
762 forest structural types and shelterwood dynamics from Lorenz-based indicators predicted by
763 airborne laser scanning. *Canadian Journal of Forest Research* 43: 1063–1074.
- 764 Valbuena R, Vauhkonen J, Packalén P, Pitkanen J, Maltamo M (2014) Comparison of airborne
765 laser scanning methods for estimating forest structure indicators based on Lorenz curves. *ISPRS*
766 *Journal of Photogrammetry and Remote Sensing* 95: 23-33
- 767 Valbuena R., Hernando A., Manzanera J.A., Görgens E.B., Almeida D.R.A., Mauro F., García-
768 Abril A. and Coomes D.A. (2018) Evaluating observed versus predicted: R-squared, index of
769 agreement or maximal information coefficient? (forthcoming).

770 Vanclay JK, Skovsgaard JP (1997) Evaluating forest growth models. *Ecological Modelling*
771 98: 1-12

772 Venables WN, Ripley B D (2002) *Modern Applied Statistics with S* (4th Ed.). Springer, New
773 York

774 Wallach D, Goffinet B (1989) Mean squared error of prediction as a criterion for evaluating
775 and comparing system models. *Ecological Modelling* 44: 299–306

776 Weisberg S (1985) *Applied linear regression* (2nd Ed.) John wiley & Sons, New York

777 White JD, Coops NC, Scott NA (2000) Estimates of New Zealand forest and scrub biomass
778 from the 3-PG model. *Ecological Modelling* 131: 175–190

779 Willmott CJ (1981) On the validation of models. *Physical Geography* 2: 184–194

780 Willmott CJ (1982) Some comments on the evaluation of model performance. *Bulletin of the*
781 *American Meteorological Society* 63 (11): 1309–1313

782 Wing BM, Ritchie MW, Boston K, Cohen WB, Gitelman A, Olsen MJ (2012) Prediction of
783 understory vegetation cover with airborne lidar in an interior ponderosa pine forest. *Remote*
784 *Sensing of Environment* 124: 730-741

785 Yebra M, Chuvieco E (2009) Linking ecological information and radiative transfer models to
786 estimate fuel moisture content in the Mediterranean region of Spain: Solving the ill-posed
787 inverse problem. *Remote Sensing of Environment* 113 (11): 2403-2411

788 Zhao K, Popescu S, Nelson R (2009) Lidar remote sensing of forest biomass: a scale-invariant
789 estimation approach using airborne lasers. *Remote Sensing of Environment* 113 (1): 182-196