

DISS. ETH NO. ....

**Integration of Small Area Estimation Procedures of Timber  
Volume Resources in Large Scale Forest Inventories**

A dissertation submitted to the  
ETH ZURICH

for the degree of  
Doctor of Sciences

presented by

ANDREAS CHRISTIAN HILL  
MSc. of Sciences

born 15 January 1986  
citizen of Germany

accepted on the recommendation of  
Prof. Dr. H.R. Heinimann, examiner  
PD Dr. D. Mandallaz, supervisor, co-examiner  
Prof. Dr. J. Saborowski, co-examiner  
Dr. Johannes Breidenbach, co-examiner

2018

# Contents

<b>Acknowledgements</b>	<b>1</b>
<b>1 The R Package forestinventory: Design-Based Global and Small Area Estimations for Multi-Phase Forest Inventories</b>	<b>2</b>
1.1 Introduction . . . . .	4
1.2 Methods and structure of the package . . . . .	6
1.3 Two-phase estimators and their application . . . . .	10
1.4 Three-phase estimators and their application . . . . .	21
1.5 Calculation of confidence intervals . . . . .	27
1.6 Special cases and scenarios . . . . .	28
1.7 Analysis and visualization . . . . .	32
1.8 Future plans . . . . .	35
<b>2 Combining canopy height and tree species map information for large scale timber volume estimations under strong heterogeneity of auxiliary data and variable sample plot sizes</b>	<b>37</b>
2.1 Introduction . . . . .	39
2.2 Materials and Methods . . . . .	41
2.3 Results . . . . .	48
2.4 Discussion . . . . .	55
2.5 Conclusion . . . . .	58
<b>3 A double-sampling extension of the German National Forest Inventory for design-based small area estimation of timber volume resources on forest district levels</b>	<b>60</b>
<b>4 Accuracy Assessment of Timber Volume Maps using Forest Inventory Data and LiDAR Canopy Height Models</b>	<b>63</b>
4.1 Introduction . . . . .	65
4.2 Materials and Methods . . . . .	67
4.3 Results . . . . .	76
4.4 Discussion . . . . .	80
4.5 Conclusion . . . . .	82
<b>Bibliography</b>	<b>84</b>
<b>Curriculum Vitae</b>	<b>91</b>

# List of Figures

1.1	Artificial representation of a local density surface. The spatial distribution of a hypothetical density function for every point in a forested area is represented as a wavy piecewise constant green surface. Sample plots (white dots) inform the inventorist of the value of the density function at that point. Note that the plateaus of constant $Y(x)$ values here have the shape of squares whereas in reality they are likely to be formed by the intersection of circles around trees. . . . .	6
1.2	(a) Concept of multi-phase sampling. The square represents the forest area for which an inventory is being conducted. The points denote the sample locations $x$ . Filled points indicate available information. (b) Illustration of the small area estimation problem. . . . .	8
1.3	Structure of the multi-phase estimators in the <i>R</i> package <b>forestinventory</b> . . . . .	9
1.4	Concept of (a) exhaustive and (b) non-exhaustive calculation of explanatory variables including boundary adjustment at the support level. Auxiliary data are in both cases available over the entire inventory area marked by the large rectangle. A vector of explanatory variables $\mathbf{Z}(x)$ is calculated within the supports (small squares) at each sample location $x$ (points) that falls into the forest area (green underlying polygon). . . . .	12
2.1	Spatial distribution of the BWI3 cluster samples over Rhineland-Palatinate . . . . .	42
2.2	Separate ALS acquisitions in Rhineland-Palatinate over the years. The colors also indicate the quality of the data: <i>light</i> : low point densities ( $0.04/m^2$ ), <i>dark</i> : high point densities ( $>4/m^2$ ). Blue semitransparent layer: state and communal forest area. . . . .	44
2.3	Identification (a) and visualization (b) of potential support radii used for calculating the predictor variables on plot level based on ECDF of maximum limiting distances of all BWI3 sample locations in RLP. . . . .	47
2.4	Classification accuracy for the main tree species of a sample location <i>before</i> and <i>after</i> calibration: <i>top</i> ) overall accuracies. <i>bottom</i> ) user's accuracies. <i>ind</i> : plot individual support sizes. . . . .	49
2.5	10-fold RMSE <sub>cv</sub> [%] and adjusted $R^2$ realized under various support choices for the CHM and <i>treespecies</i> explanatory variables . . . . .	51
2.6	Effect on the adjusted $R^2$ when substituting the actual main tree species with the predicted main tree species of a sample plot. The <i>dotted</i> line tracks the the model with the highest adjusted $R^2$ under the use of the error-free <i>treespecies</i> variable. Semitransparent colours for the data points are used to visualize overlap. . . . .	52
2.7	Visualization of the timber volume prediction function ( <i>final regression model</i> ) on sample plot level for each main plot tree species and ALS acquisition year. For visualization purposes, the predictor variable <i>stddev</i> was set to its average value within the respective <i>treespecies</i> and <i>ALSpyear</i> categories. The terrestrially observed timber volume values are plotted in the background. . . . .	53
2.8	$R^2$ -values of the final regression model, submodel 1 and submodel 2 achieved <i>within</i> the ALS acquisition year strata. . . . .	55

4.1	Study site, including the distribution of the 67 field plots (regional forest inventory) that are part of the forest area derived from the TLM3D (Topographic Landscape Model) data (with approval of swisstopo JA100120/JD100042). . . . .	68
4.2	Terrestrial observed timber volume on plot level for all 67 sample plots of the study area (histogram bandwidth = 50 m <sup>3</sup> /ha; origin = 0 m <sup>3</sup> /ha). . . . .	69
4.3	Conceptual model illustrating the workflow of producing the timber volume map and assessing its accuracy. . . . .	70
4.4	Schematic design of the timber volume map. . . . .	75
4.5	Model-predicted timber volume on plot level against the observed timber volume of the 67 field plots; the distribution of the predictions and observations are also indicated on the respective axis. . . . .	76
4.6	Overall accuracies with 95% confidence intervals and corresponding kappa coefficients for a pre-defined number of classes, calculated for constant and locally-adapted class widths. . . . .	77
4.7	Volume map with model-predictions on a continuous scale for the entire study area covering 2000 hectares of forest with a spatial resolution of 25 meters; the subareas show the classified volume map using a constant class width of 200 m <sup>3</sup> /ha (upper) and locally-optimized class widths for five classes (lower). . . . .	79

# List of Tables

2.1	Descriptive statistics of the forest observed on NFI sample plots located iwithin communal and state forest area (n=5791). . . . .	43
2.2	Accuracy metrics for submodels of final OLS regression model. $p$ gives the number of parameters for each model. Interaction terms are indicated by ':'. . . . .	53
2.3	$R^2$ , RMSE and RMSE% of final regression model within ALS acquisition year strata ( $ALSyyear$ ). $Area_{ALSyyear}$ : Area covered by ALS acquisition given in $\text{km}^2$ . $n$ : number of validation data. . . . .	54
4.1	Summary statistics of the timber volume observed at the 67 terrestrial sample plots (given in $\text{m}^3/\text{ha}$ ). . . . .	69
4.2	Summary statistics of the timber volume observed at the 67 terrestrial sample plots (given in $\text{m}^3/\text{ha}$ ). . . . .	70
4.3	Summary statistics of the CHM metrics calculated at the 67 terrestrial sample plots (in meters). . . . .	72
4.4	Summary statistics of the regression model. . . . .	76
4.5	Constant class width of 200 $\text{m}^3/\text{ha}$ (five timber volume classes); OAA ( $\pm 95\%$ confidence interval) (%): 55.22 (42.58, 67.4); kappa: 0.36; square sum of class width: 200'000. . . . .	78
4.6	Optimized class width for five timber volume classes; OAA ( $\pm 95\%$ confidence interval) (%): 64.18 (51.53, 75.53); kappa: 0.54; square sum of class width: 176'600. Class widths smaller than 200 $\text{m}^3/\text{ha}$ are indicated by *. . . . .	78

# Acknowledgements

The thesis reveals the aggregated results of my work at the Section of Forest Growth Modelling and Computer Sciences of the Northwest German Research Institute and the Department of Forest Economics and Forest Management of the University of Göttingen. Along the way, many people made significant contribution to the presented work. Its completion would have been unlikely without the constructive help and the continuous support from my colleagues, family and friends, whom I thank herewith. In particular, I thank

- Prof. Dr. Jürgen Nagel - for the supervision of the thesis, critical revision of my workings and the freedom to explore new ideas.
- Prof. Dr. Bernhard Möhring - for acting as co-referee, the constructive and rational support of economic problems and the opportunity to keep going.
- Prof. Dr. Thomas Kneib - for being oral examiner in the disputation.
- Prof. Dr. Hermann Spellmann - for the decision support to write this thesis and the trustful and cooperative working environment.
- Christoph Fischer - for valuable help in countless tasks.
- Dr. Jan Hansen - for the numerous hours of help to connect the TreeGrOSS packages with R.
- My co-authors Alexander Lange, Elmar Spiegel and Sabine Rumpf - for the successful submissions.
- Andreas Hill - for many hours with constructive and interesting discussions and for proofreading.
- Dr. Matthias Schmidt and Dr. Egbert Schönfelder - for supporting statistical and methodological questions.
- Jan Butschkow, Tim Koddenberg, Michael Hill and Sebastian Ohrmann - for proofreading.
- The partners from the BioEconomy research project - for supporting practical questions regarding forestry, logistics and wood processing.
- My family - for miscellaneous supports.
- Laura - for giving me the courage to tackle the challenges and the strength to carry on.

# Chapter 1

## The *R* Package **forestinventory**: Design-Based Global and Small Area Estimations for Multi-Phase Forest Inventories

Andreas Hill<sup>1</sup>, Alexander Massey<sup>1</sup>

<sup>1</sup>ETH Zürich

Department of Environmental Systems Science, Universitätstrasse 16, 8092 Zürich, Switzerland

Submitted to:  
*Journal of Statistical Software* (in review).

- Alexander Massey is co-author of the R package *forestinventory* and supported writing of the manuscript.

# Abstract

Forest inventories provide reliable evidence-based information to assess the state and development of forests over time. They typically consist of a random sample of plot locations in the forest that are assessed individually by hand. Due to the high costs of these terrestrial campaigns, remote sensing information available in high quantity and low costs is frequently incorporated in the estimation process in order to reduce inventory costs or improve estimation precision. With respect to this objective, the application of multi-phase forest inventory methods (e.g., double- and triple-sampling regression estimators) has proved to be efficient. While these methods have been successfully applied in practice, the availability of open-source software has been rare if not non-existent. The *R* package **forestinventory** provides a comprehensive set of global and small area regression estimators for multi-phase forest inventories under simple and cluster sampling. The implemented methods have been demonstrated in various scientific studies ranging from small to large scale forest inventories, and can be used for post-stratification, regression and regression within strata. This article summarizes the mathematical theory of this family of design-based estimators and demonstrates their application in the *R* environment.

## 1.1 Introduction

In many countries, forest inventories have become an indispensable tool for evaluating the current state of forests as well as for tracking their development over time. They provide accurate quantitative information that can be used to define management actions and to adapt forest management strategies according to guidelines on national and international levels. As conducting a full census of all trees within any sizable forest area is clearly infeasible due to time and cost restrictions, forest inventories usually gather their information by means of statistical sampling methods. Typically this means that discrete sample locations (sample plots) are randomly chosen in the forest, making up the framework of a terrestrial inventory. This terrestrial sample data is then used to make estimates for the entire forested area and provide a measure of precision for those estimates in the form of confidence intervals. There is a broad range of literature describing the concepts and methods regarding the choice of different estimators under various sample designs (Gregoire & Valentine, 2007; Köhl et al., 2006; Schreuder et al., 1993; Mandallaz, 2008a).

Terrestrial inventories have the benefit of being very flexible in the sense that they can be used to produce high quality estimates for a wide-variety of different forest attributes. However, they have the downside of being very expensive. Improving the precision of the estimates by increasing the number of sample plots essentially means that travel costs will rise as trained inventors are sent to more and more plot locations. This is why the number of terrestrial samples is often limited. Although national inventories usually possess a sufficiently large terrestrial sample size to provide high estimation accuracies for larger areas, this is often not the case for smaller areas, such as forest management units. As a result, there has been an increasing need for alternative inventory methods that can maintain the same estimation precision at lower costs, or achieve higher estimation precision at identical costs (von Lüpke, 2013). A method which has become particularly attractive is called multi-phase sampling. The core concept is to enlarge the sample size in order to gain higher estimation precision without enlarging the terrestrial sample size. This is done by using predictions of the terrestrial target variable at additional sample locations where the terrestrial information has not been gathered. These predictions are produced by regression models that use explanatory variables derived from auxiliary data, commonly in the form of spatially exhaustive remote sensing data in the inventory area. Regression estimators using this concept can consider either one additional sample of plot locations (two-phase or double-sampling) or two additional samples available in different sample sizes (three-phase or triple-sampling) (Gregoire & Valentine,

2007; Saborowski et al., 2010; Mandallaz, 2013a,d; von Lüpke et al., 2012). (Gregoire & Valentine, 2007; Saborowski et al., 2010; Mandallaz, 2013a,d; von Lüpke et al., 2012). Their application to existing forest inventory systems have already showed their efficiency in terms of cost reduction and gain in estimation precision (Breidenbach & Astrup, 2012; von Lüpke & Saborowski, 2014; Mandallaz et al., 2013; Magnussen et al., 2014; Massey et al., 2014).

Multi-stage and multi-phase estimation has already been implemented in commercial as well as open-source software, such as the survey sampling procedures in *SAS* (*SAS* Institute Inc., 2015) and the **survey** package in *R* (Lumley, 2016). However, both are targeted towards list-sampling as it is applied in official statistics. Available software providing multi-phase sampling methods better suited for forest inventories has been rare. Two exceptions are the *R* package **JoSAE** by Breidenbach (2015) and the **maSAE** package by Cullmann (2016). However, a more comprehensive software package covering a larger variety of sample designs and estimators for forest inventories has not yet been available, which is the motivation behind the *R* package **forestinventory**. The package provides global and small area estimators for two-phase and three-phase forest inventories under simple and cluster sampling, which have been developed under the infinite population approach by Daniel Mandallaz at ETH Zurich between 2008 and 2017. The implemented methods have been demonstrated by case studies in Switzerland (Massey et al., 2014; Massey & Mandallaz, 2015b; Mandallaz et al., 2013) and Germany (Hill et al., 2017). The implemented estimators cover 32 inventory scenarios and can be used for post-stratification, regression and regression within strata (Massey, 2015). The long-term objective of **forestinventory** is to make the broad range of estimators available to a large user community and to facilitate their application in science as well as operational forest management.

The objectives of this article are to a) establish the link between the mathematical description of the estimators and their implementation in our package, b) illustrate their application to real-world inventory data sets and c) address special cases and demonstrate how the package-functions handle them.

## 1.2 Methods and structure of the package

### 1.2.1 Forest inventory in the infinite population approach

Most forest inventories gather terrestrial information by sending field crews to randomly (or systematically) selected points in the forest area defined by coordinates. The crew then defines a sample plot by measuring individual trees located within one or multiple constructed inclusion circles around the sample point  $x$ , and aggregating their characteristics (e.g., timber volumes) to local plot densities (e.g., the timber density in  $\text{m}^3/\text{ha}$ ). The estimators implemented in **forestinventory** use the so called infinite population approach in order to bridge this inventory process to the mathematics behind the estimators. This approach assumes that the spatial distribution of the local density,  $Y(x)$ , over the forest area is determined by a fixed piecewise constant function, as visualized in Fig. 1.1. The population total of the target variable (e.g., the total timber volume of the forest) is mathematically equivalent to the integral of that density function, which is depicted in Fig. 1.1 as the volume under the density surface. From this perspective, the practical challenge is that the form of this function is unknown. Theoretically, we could get the total timber volume by observing the function value, i.e., the local density  $Y(x)$ , at each possible point  $x$  over the forest area and taking their sum. However, this is impossible because there is an infinite number of points in the forest area. Our strategy is thus to take a sample of points,  $s_2$ , from an infinite population of possible points and use their associated local densities  $Y(x)$  to estimate the integral  $Y = \frac{1}{\lambda(F)} \int_F Y(x) dx$  with  $\hat{Y} = \frac{1}{n_2} \sum_{x \in s_2} Y(x)$ . The total timber volume can then be obtained by multiplying  $\hat{Y}$  by the surface area of the forest,  $\lambda(F)$ . All estimators included in **forestinventory** rest upon this theoretical perspective. The key point in the infinite population approach is that a local density value  $Y(x)$  is associated with the sample point  $x$ , which constitutes the sample unit, and not with the sample plot area. This has some theoretical advantages over the finite population approach, where the sample units are the actual plot areas usually assumed to be either circular or rectangular. This is mainly due to the impossibility of a perfect tessellation over an amorphous forest area by any choice of plot shape. Hence, the population in the finite approach is, strictly speaking, not well defined with respect to the forest area. The consideration of an underlying infinite population of sample points will also play an important role when incorporating auxiliary information in the frame of two- and three-phase estimation methods.

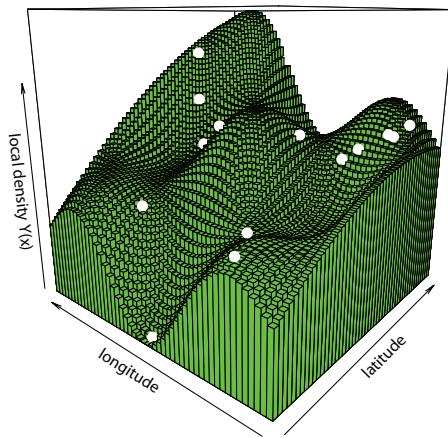


Figure 1.1: Artificial representation of a local density surface. The spatial distribution of a hypothetical density function for every point in a forested area is represented as a wavy piecewise constant green surface. Sample plots (white dots) inform the inventorist of the value of the density function at that point. Note that the plateaus of constant  $Y(x)$  values here have the shape of squares whereas in reality they are likely to be formed by the intersection of circles around trees.

### 1.2.2 Two-phase sampling

The two-phase or double-sampling estimators use inventory information from two nested samples which are commonly referred to as phases (Fig. 1.2a). The first phase  $s_1$  comprises  $n_1$  sample locations that each provide a set of explanatory variables described by the column vector  $\mathbf{Z}(x) = (z(x)_1, z(x)_2, \dots, z(x)_p)^\top$  at each point  $x \in s_1$ . These explanatory variables are derived from auxiliary information that is available in high quantity within the forest area  $F$ . The second phase  $s_2$  constitutes the terrestrial inventory conducted at  $n_2$  subsamples of the large phase  $s_1$  and provides the value of the target variable, i.e., the local density  $Y(x)$  (e.g., the timber volume per hectare). For every  $x \in s_1$ ,  $\mathbf{Z}(x)$  is transformed into a prediction  $\hat{Y}(x)$  of  $Y(x)$  using the choice of some model, which in **forestinventory** is always a linear model fit in  $s_2$  using ordinary least squares (OLS). The basic idea of this setup is to boost the sample size by providing a large sample of less precise but cheaper predictions of  $Y(x)$  in  $s_1$  and to correct any possible model bias, i.e.,  $\mathbb{E}(Y(x) - \hat{Y}(x))$ , using the subsample of terrestrial inventory units where the value of  $Y(x)$  is observed. In the design-based context, the two-phase estimator is typically unbiased regardless of the model used to produce the predictions. This property comes from the assumption that each phase's sample is selected via simple random sampling (see Section 1.2.5).

### 1.2.3 Three-phase sampling

Three-phase estimators extend the principle of two-phase sampling and use inventory information from three nested samples (phases) (Fig. 1.2a). The basic setup is that the explanatory variables calculated from the auxiliary information are available in two different frequencies. The phase  $s_0$  provides a large number  $n_0$  of auxiliary data, whereas the phase  $s_1$  provides additional auxiliary data that are only available at  $n_1$  subsamples of  $s_0$ . The terrestrial information is then collected at a further subsample  $s_2$  of  $s_1$ . The motivation for three-phase sampling is that the additional set of explanatory variables available at  $s_1$ , now denoted  $\mathbf{Z}^{(1)}(x)$ , adds considerable explanatory power to the set of variables available at all sample locations  $x \in s_0$ , denoted  $\mathbf{Z}^{(0)}(x)$ . From that it follows that we can define two nested regression models. The full set of predictor variables  $\mathbf{Z}^\top(x) = (\mathbf{Z}^{(0)\top}(x), \mathbf{Z}^{(1)\top}(x))$  can be used to calculate the predictions  $\hat{Y}(x)$  of  $Y(x)$  at all sample locations  $x \in s_1$ . The regression model applicable to the  $s_1$  phase is thus referred to as the full model. Less accurate predictions,  $\hat{Y}^{(0)}(x)$ , can be produced at all the sample locations  $x \in s_0$  using only the reduced set of explanatory variables  $\mathbf{Z}^{(0)}(x)$ . If there is a significant gain in model precision between the reduced and the full model and the sampling fraction between  $s_0$  and  $s_1$  is sufficiently large, the three-phase estimator normally leads to a further increase in estimation precision compared to the two-phase estimator.

### 1.2.4 Small area estimation

Small area estimation does not necessarily refer to small spatial areas but rather to areas that contain little or no terrestrial sample. To formulate this mathematically, we want to make an estimate for a subregion  $G$  of the entire inventory area  $F$  (Fig. 1.2b). As the sample size in the small area,  $n_{2,G}$ , is usually too small to provide sufficient estimation precision, multi-phase estimation can be efficient. However,  $n_{2,G}$  may also be too small to justify fitting a separate regression model just for that area because the estimates produce undesirably large confidence intervals. The idea is then to borrow strength from the entire terrestrial sample  $s_2$  of  $F$  to fit the model, and to apply this model to the small area. The potential bias of applying that model in  $G$  is then corrected for by using the empirical model residuals derived from that small area. If there are no terrestrial plots in  $G$  (i.e.,  $n_{2,G} = 0$ ), one cannot correct for a potential model bias in  $G$  and has to accept a potential bias in the estimator. These are called synthetic estimates and despite their potential bias, it is usually still possible to calculate their design-based variance.

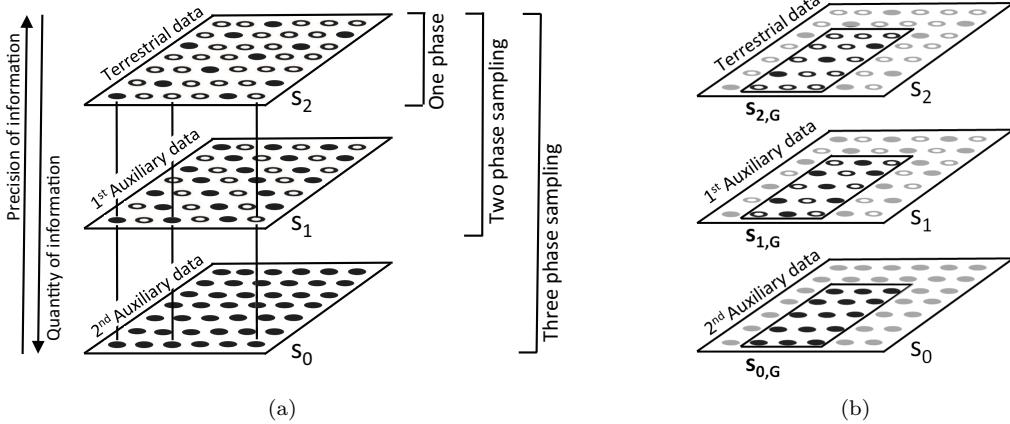


Figure 1.2: (a) Concept of multi-phase sampling. The square represents the forest area for which an inventory is being conducted. The points denote the sample locations  $x$ . Filled points indicate available information. (b) Illustration of the small area estimation problem.

### 1.2.5 Design-based vs. model-dependent approach

The subject of model selection gets a lot of attention in the field of forest inventory. This is why it is important to understand that the mathematical interpretation of how a model is used to produce estimates is fundamentally different between the design-based and model-dependent approach. In the model-dependent (also known as model-based) framework, the sample locations  $x$  are fixed and the observation  $Y(x)$  taken at location  $x$  is assumed to be a random variable as the forest is assumed to be the realization of a stochastic process. Although the model does not need to be fit from a probability sample, i.e., the sample locations could arbitrarily be chosen, the model should adequately describe the underlying stochastic process in order to efficiently ensure unbiased results. In practice this means that special attention must be made to ensure that the variable selection is appropriate to avoid overfitting, important variables are not omitted and all model assumptions are reasonably met through empirical verification. If a model is misspecified then estimation based on inference from that model may not be reliable. In the model-dependent framework one thus has to trust the model. In contrast, the design-based approach, on which all **forestinventory** estimators are based, rests upon the randomization of the sample locations  $x$ . While the sample locations  $x$  are independently and uniformly distributed in the forest, the forest itself and thus the values of the local density surface at any location  $x \in F$  are fixed and not the result of a stochastic process. A selected observation  $Y(x)$  still remains a random variable, but solely due to the random sample mechanism. A consequence of this approach is that the estimation properties of design-based regression estimators (e.g., unbiasedness) typically hold regardless of the model that is chosen. The philosophy of the design-based approach is thus to use prediction models to improve the efficiency of the estimators without having to rely on their correct specification, which makes them very attractive to be used in official statistics. They are therefore also referred to as model-assisted. It should be noted that the randomization of sample locations upon which design-based inference depends, is in practice often replaced by systematic grids to minimize travelling costs in the terrestrial survey. However, there is reasonable evidence that softening this assumption is acceptable for point and variance estimation as long as the grid does not interact with periodic features in the forest structure (Mandallaz, 2008a). The variance will in most cases be slightly overestimated and lead to wider, more conservative confidence intervals (Mandallaz, 2013a).

### 1.2.6 Package structure

In the **forestinventory** package, estimators for two-phase and three-phase sampling are applied with the `twophase()` and `threephase()` functions. From these two overall function calls, various estimators for specific inventory scenarios under the chosen sampling design can be applied (Fig. 1.3). Choosing an estimator follows a tree-like structure which can serve the user as a guideline throughout this article as well as in future applications. The basic decision to make is whether an estimate and its variance should be computed for an entire inventory area (global estimators) or only for subregions of the entire inventory area (small area estimators). In the second case, the package offers three small area estimators that will in detail be described in the following sections. The estimators are available under exhaustive and non-exhaustive use of the auxiliary data. Additionally, the package can also calculate one-phase estimates solely based on terrestrial samples. All estimators are also available for cluster sampling, in which case a sample unit consists of multiple, spatially agglomerated samples. The following sections describe the mathematical details and the application of the multi-phase estimators implemented in the *R* package **forestinventory**. While Mandallaz (2008a, 2013c,b, 2015) provides an extensive derivation of all estimators, we will provide the mathematical formulas that are actually implemented in the package. We will also restrict discussion to simple sampling, while the formulas for cluster sampling are available in the technical reports (Mandallaz et al., 2016; Mandallaz, 2013c,b). A special case under cluster sampling is described in Section 1.6.

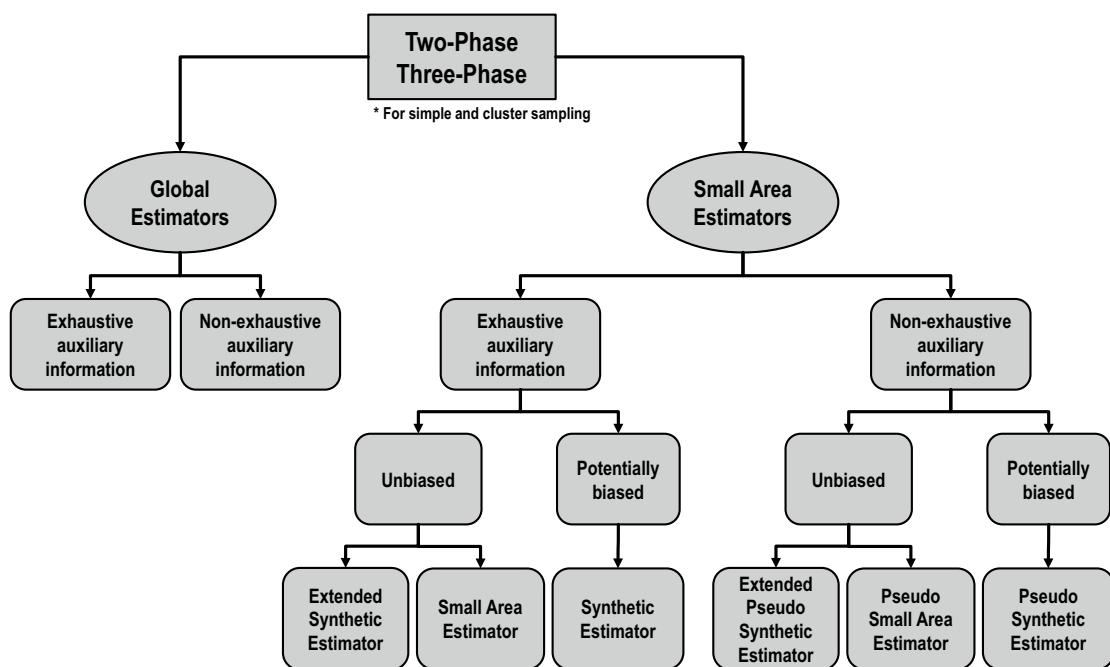


Figure 1.3: Structure of the multi-phase estimators in the *R* package **forestinventory**.

## 1.3 Two-phase estimators and their application

### 1.3.1 Global estimators

#### Mathematical background

The vector of regression coefficients of the OLS regression model is found by using the following solution to the sample-based normal equation:

$$\hat{\beta}_{s_2} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \left( \frac{1}{n_2} \sum_{x \in s_2} Y(x) \mathbf{Z}(x) \right) \quad (1.1)$$

The individual predictions can then be calculated as  $\hat{Y}(x) = \mathbf{Z}^\top(x) \hat{\beta}_{s_2}$  and the empirical model residuals, which are only available at all sample locations  $x \in s_2$ , are calculated as  $\hat{R}(x) = Y(x) - \hat{Y}(x)$ . Unless stated otherwise, **forestinventory** only uses internal models to calculate estimates. This means that the model fit, i.e.,  $\hat{\beta}_{s_2}$ , is derived from the current inventory data that are passed to the `twophase()` and `threephase()` functions. While virtually all inventorists fit their models using the current inventory data, sometimes there is reason to use formulas derived from external models where the sample used to train the model is assumed to be taken from an independent source (Massey & Mandallaz, 2015a). However, this usually occurs when using a model other than the OLS regression model and is beyond the scope of the package at this time.

The package provides the calculation of point estimates under exhaustive (EX) and non-exhaustive (NEX) use of the auxiliary information, which means to respectively apply  $\hat{\beta}_{s_2}$  to  $\bar{\mathbf{Z}}$ , i.e., the exact spatial mean of  $\mathbf{Z}(x)$ , or to  $\hat{\bar{\mathbf{Z}}}$ , i.e., an estimate of the spatial mean of  $\mathbf{Z}(x)$ :

$$\hat{Y}_{reg2p,EX} = \bar{\mathbf{Z}}^\top \hat{\beta}_{s_2} \quad (1.2a)$$

$$\hat{Y}_{reg2p,NEX} = \hat{\bar{\mathbf{Z}}}^\top \hat{\beta}_{s_2} \quad (1.2b)$$

Note that for internal linear models the mean of the empirical residuals  $\frac{1}{n_2} \sum_{x \in s_2} \hat{R}(x)$  is zero by construction (zero mean residual property) which is why it does not appear in the point estimate. More explanation about how to obtain the auxiliary means is given in the next subsection.

The **forestinventory** package implements two kinds of variances for each of these point estimates: the g-weight formulation that accounts for the fact that our model is in fact internal, and the external variance formulation that assumes a true external regression model and thus neglects the uncertainty in the regression coefficients (Mandallaz et al., 2016).

The g-weight formulation is

$$\hat{\mathbb{V}}(\hat{Y}_{reg2p,EX}) := \bar{\mathbf{Z}}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \bar{\mathbf{Z}} \quad (1.3a)$$

$$\hat{\mathbb{V}}(\hat{Y}_{reg2p,NEX}) := \hat{\bar{\mathbf{Z}}}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\bar{\mathbf{Z}}} + \hat{\beta}_{s_2}^\top \hat{\Sigma}_{\hat{\bar{\mathbf{Z}}}} \hat{\beta}_{s_2} \quad (1.3b)$$

where the g-weight variance-covariance matrix of  $\hat{\beta}_{s_2}$  is calculated as

$$\hat{\Sigma}_{\hat{\beta}_{s_2}} := \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \left( \frac{1}{n_2^2} \sum_{x \in s_2} \hat{R}^2(x) \mathbf{Z}(x) \mathbf{Z}^\top(x) \right) \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \quad (1.4)$$

and the uncertainty caused by using the  $s_1$  sample to estimate  $\bar{\mathbf{Z}}$  by  $\hat{\bar{\mathbf{Z}}}$  is accounted for by the

variance-covariance matrix of the auxiliary vector  $\hat{\mathbf{Z}}$

$$\hat{\Sigma}_{\hat{\mathbf{Z}}} = \frac{1}{n_1(n_1 - 1)} \sum_{x \in s_1} (\mathbf{Z}(x) - \hat{\mathbf{Z}})(\mathbf{Z}(x) - \hat{\mathbf{Z}})^\top \quad (1.5)$$

The external variance formulation for linear regression models is

$$\begin{aligned}\hat{\mathbb{V}}_{ext}(\hat{Y}_{reg2p,EX}) &= \frac{1}{n_2} \hat{\mathbb{V}}_{s_2}(\hat{R}(x)) \\ \hat{\mathbb{V}}_{ext}(\hat{Y}_{reg2p,NEX}) &= \frac{1}{n_1} \hat{\mathbb{V}}_{s_1}(\hat{Y}(x)) + \frac{1}{n_2} \hat{\mathbb{V}}_{s_2}(\hat{R}(x))\end{aligned}\quad (1.6a)$$

where  $\hat{\mathbb{V}}_{s_2}$  and  $\hat{\mathbb{V}}_{s_1}$  indicate taking the sample variance over  $s_2$  and  $s_1$  respectively.

Note that when applied to internal linear regression models, the external variance is asymptotically unbiased and usually slightly smaller than the g-weight variance, where the uncertainty of the regression coefficients is accounted for by the variance-covariance matrix (Eq. 1.4). The external variances are provided in the package **forestinventory** in case the user wants to compare linear models to another model type where no g-weight formulation is possible, as is the case with non-parametric models like kNN.

### Calculation of explanatory variables

We will now draw our attention to the calculation of the explanatory variables from the auxiliary data for both the non-exhaustive and exhaustive cases. Fig. 1.4b depicts how the non-exhaustive case often looks like in practice: a regular terrestrial grid  $s_2$  is given by a terrestrial inventory (the points surrounded by dotted circles) and densified to a larger sample  $s_1$  (the points). For every point  $x$ , each explanatory variable in the vector  $\mathbf{Z}(x) = (z(x)_1, z(x)_2, \dots, z(x)_p)^\top$  is calculated using a defined spatial extent of auxiliary information around that point called the support (the dark green square tiles). We emphasize that the value of the explanatory variables for  $\mathbf{Z}(x)$  are associated with the sample point whereas the support is the spatial extent of the auxiliary information used to calculate those values. So far this is in perfect agreement with the presented theory of the non-exhaustive estimator, except for using regular grids rather than randomly placed sample points. The **forestinventory** package calculates the empirical mean of  $\mathbf{Z}(x)$  automatically from the input data frame using  $\hat{\mathbf{Z}} = \frac{1}{n_1} \sum_{x \in s_1} \mathbf{Z}(x)$ .

The exhaustive case requires a closer look. In the infinite population approach,  $\mathbf{Z}(x)$  refers to the sample point  $x$  and not the area around it. Deriving the exact spatial mean,  $\bar{\mathbf{Z}} = \frac{1}{\lambda(F)} \int_F \mathbf{Z}(x) dx = (\frac{1}{\lambda(F)} \int_F z_1(x) dx, \dots, \frac{1}{\lambda(F)} \int_F z_p(x) dx)^\top$ , implies that we need to calculate the spatial mean of each component of  $\mathbf{Z}(x)$  using all possible points in  $F$ . This is much like the situation we had with calculating the mean of the local density surface for  $Y(x)$  in that we need to find the mean of  $\mathbf{Z}(x)$  over an infinite number of sample points (i.e.,  $n_1 = \infty$ ). Although it is practically infeasible to assess  $\mathbf{Z}(x)$  for every  $x$ , there are few cases where the exact mean can in fact be precisely calculated. The first case is when the explanatory variables are provided by polygon layers (e.g., map of development stages). In this case, one can calculate the exact mean as the area-weighted average of each categorical variable. The second case is when the exact mean can be calculated in one step, e.g., taking the mean of all height pixels of a raster canopy height model will perfectly equal the mean calculated by the use of an infinite number of supports (Mandallaz et al., 2013). However, for most types of explanatory variables we will try to get an approximation of  $\bar{\mathbf{Z}}$  that is only negligibly different.

One implementation to approximate the exact mean  $\bar{\mathbf{Z}}$  is shown in Fig. 1.4a, where the spatial arrangement of the supports (the dark green tiles) are tessellated to form a perfect partition over

the inventory area in order for all of the wall-to-wall auxiliary information to be exploited. It has to be noted that this setup would allow for a perfect calculation of the exact mean  $\bar{\mathbf{Z}}$  in the finite population approach, i.e., deriving  $\mathbf{Z}(x)$  for the finite population of supports that are considered the sampling units. While in the infinite population approach this implementation probably does not produce the true exact mean  $\bar{\mathbf{Z}}$ ,  $n_1$  is still expected to be reasonably large for the difference to be considered negligible as long as the size of the supports are not unreasonably large. However, the perfect tessellation implementation can also impose drawbacks. One is that a perfect tessellation by the supports strongly depends on the distance between the sample locations of  $s_1$  and the support size. Since in practice the support size should ideally be chosen to achieve a best possible explanatory power of the regression model (thus minimizing the residual variation) a perfect tessellation might often not be feasible. In the infinite population frame, the supports are allowed to overlap if this seems necessary to acquire a sufficiently large sample  $n_1$  to get a negligibly close approximation of  $\bar{\mathbf{Z}}$ . With this respect, the infinite population approach provides more flexibility than the finite approach.

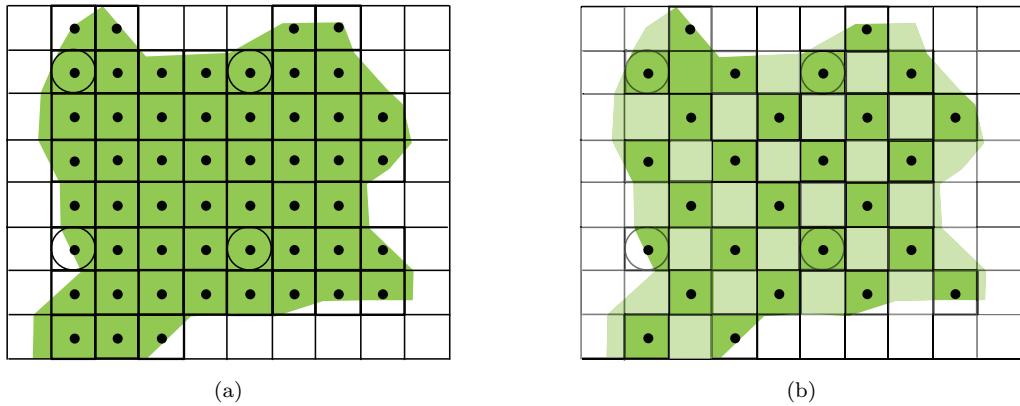


Figure 1.4: Concept of (a) exhaustive and (b) non-exhaustive calculation of explanatory variables including boundary adjustment at the support level. Auxiliary data are in both cases available over the entire inventory area marked by the large rectangle. A vector of explanatory variables  $\mathbf{Z}(x)$  is calculated within the supports (small squares) at each sample location  $x$  (points) that falls into the forest area (green underlying polygon).

### Boundary adjustment

An extension to the so-far published estimators by Mandallaz is the consideration of a boundary adjustment. In forest inventories, the sample is often restricted to those sample locations located within the forest area. In case a consistent forest definition can be applied to both the  $s_2$  and  $s_1$  sample (e.g., by a polygon forest mask layer), it might be desired to restrict the calculation of the explanatory variables to the forest area within the given support (see Fig. 1.4). This method was suggested in Mandallaz et al. (2013) and led to an improvement in estimation precision. In order to ensure an unbiased calculation of either  $\hat{\mathbf{Z}}$  or  $\bar{\mathbf{Z}}$ , the respective means have then to be calculated as the weighted mean (Eq. 1.7) where the weight  $w(x)$  is equal to the percentage of forested area within the support of sample location  $x$ .

$$\hat{\mathbf{Z}} = \frac{\sum_{x \in s_1} w(x)\mathbf{Z}(x)}{\sum_{x \in s_1} w(x)} \quad (1.7)$$

## Application

To demonstrate the use of the global two-phase estimators, we will use the **grisons** data set that comes with installing the package from the CRAN repository. The data set contains data from a simple (i.e., non-cluster) two-phase forest inventory conducted in 2007 that was used in Mandallaz et al. (2013) as a case study. The  $s_1$  sample is comprised of 306 sample locations arranged on a systematic grid containing auxiliary information in the form of airborne laserscanning (LiDAR) canopy height metrics (**mean**, **stddev**, **max**, **q75**). For a systematic subsample of 67 ( $s_2$  sample), terrestrial information of the timber volume per hectare (**tvol**) on the sample plot level is provided from a terrestrial survey. We can load **forestinventory** and examine the **grisons** data set in the *R* environment as follows:

```
R> library("forestinventory")
R> data("grisons", package = "forestinventory")
R> head(grisons)
```

	phase_id_2p	boundary_weights	mean	stddev	max	q75	smallarea	tvol
1	2		1.00	9.30	11.84	40.87	21.14	C 107.80
2	1		1.00	12.16	11.35	39.80	21.54	A NA
3	2		1.00	5.25	5.74	23.82	9.53	D 63.77
4	1		1.00	7.53	9.33	34.10	13.02	A NA
5	1		0.67	6.11	5.87	23.33	10.55	B NA
6	1		1.00	12.15	10.16	33.76	20.97	C NA
7	2		1.00	6.38	4.72	17.96	10.14	D 154.10
8	1		1.00	1.25	3.79	22.72	0.00	B NA
9	1		1.00	21.56	7.49	32.66	27.81	A NA
10	2		1.00	13.55	7.20	36.14	18.59	A 256.15

Estimates can be made using the **onephase()**, **twophase()** or **threephase()** functions. The data frame inputted to these functions must have the structure where each row corresponds to a unique sample location and the columns specify the attributes associated to that respective sample location. Attributes that are missing, e.g., because they are associated with sample locations that were not selected in the subsample for the subsequent phase, should be designated as **NA** and the phase membership is encoded as numeric.

For global two-phase estimation, we have to specify

- the regression model (**formula**) as specified in the **lm()**-function (R Core Team, 2017).
- the inputted **data.frame** containing the inventory information (**data**).
- the **list**-object **phase\_id** containing: the **phase.col** argument identifying the name of the column specifying membership to  $s_1$  or  $s_2$ , and the **terrgrid.id** argument specifying which numeric value indicates  $s_2$  membership in that column. Note that **forestinventory** implicitly assumes that all rows not indicated as  $s_2$  belong to the  $s_1$  phase.
- the name of the column containing the weights  $w(x)$  of the boundary adjustments (optional).

The non-exhaustive estimator with boundary weight adjustment can thus be applied as follows:

```
R> reg2p_nex <- twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   boundary_weights = "boundary_weights")
```

The **twophase()** function creates an S3 object of class "**twophase**" with subclass "**global**". A readable summary of the estimation results can be obtained by passing this object to the **summary()** function, which automatically interprets what type of estimator was used and returns pertinent information such as the regression model formula, the point estimate (**estimate**), the g-weight and external variance (**g\_variance** and **ext\_variance**) as well as the sample sizes and the model  $R^2$ :

```
R> summary(reg2p_nex)
Two-Phase global estimation

Call:
twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  boundary_weights = "boundary_weights")

Method used:
Non-exhaustive global estimator

Regression Model:
tvol ~ mean + stddev + max + q75

Estimation results:
estimate ext_variance g_variance n1 n2 r.squared
383.5354      279.954   271.5057 306 67 0.6428771

'boundary_weight'- option was used to calculate weighted means of auxiliary variables
```

For practical use, one should normally always prefer the g-weight variance over the external variance. This is because when we use internal models, the regression coefficients actually depend on the terrestrial sample realized by the sampling design. In contrast to the external variance, the g-weight variance accounts for this sampling variability which results in more reliable point and variance estimates and also enjoys better statistical calibration properties (g-weights). The external and g-weight variances are asymptotically equivalent but the external variance is really only included here in case the user wants to compare to another estimator where no g-weight variance exists.

The exhaustive estimator can be applied by additionally passing a vector containing the exact means of the explanatory variables, i.e.,  $\bar{\mathbf{Z}}$ , to the optional argument `exhaustive`. This vector must be calculated beforehand in such a way that any desired boundary adjustment has already been applied. Note that the vector input to `exhaustive` must be in the same order that the `lm()`-function processes a `formula` object including the intercept term whose exact mean will always be 1. Particular caution must be taken if categorical variables are present because the `lm()`-function, which is internally used to set up the design-matrix, automatically creates dummy variables with one of the categories used as a reference. Using our `grisons` example, the correct order can always be extracted by the following *R*-code:

```
R> colnames(lm(formula = tvol ~ mean + stddev + max + q75, data = grisons,
+   x = TRUE)$x)
```

The exhaustive estimator can be applied after defining the vector of exact means  $\bar{\mathbf{Z}}$  taken from Mandallaz et al. (2013), denoted as `true.means.Z`:

```
R> true.means.Z <- c(1, 11.39, 8.84, 32.68, 18.03)
R> reg2p_ex <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons,
+   phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   exhaustive = true.means.Z)
```

An alternative way to look at the estimation results without using the `summary()` is to query `reg2p_ex` directly:

```
R> reg2p_ex$estimation
estimate ext_variance g_variance n1 n2 r.squared
1 376.7426      202.5602   187.2787 Inf 67 0.6428771
```

Note that both variances of the exhaustive estimation are smaller than those of the non-exhaustive estimation. This is essentially because we eliminated one component of uncertainty by substituting the estimated means of the explanatory variables  $\hat{\bar{Z}}$  by their exact means  $\bar{Z}$ .

### 1.3.2 Small area estimators

#### Mathematical background

The **forestinventory** package provides three types of small area estimators each of which has an exhaustive and non-exhaustive form. We will use a different nomenclature for the non-exhaustive case in small area estimation since much of the existing literature shows preference for the label pseudo to indicate that the mean of the explanatory variables within the small area was based on a finite sample. The main idea for all these small area estimators is to calculate the regression coefficient vector  $\hat{\beta}_{s_2}$  and its variance-covariance matrix  $\hat{\Sigma}_{\hat{\beta}_{s_2}}$  on the entire  $s_2$  sample according to Eq. 1.1 and 1.4, and subsequently use that to make predictions for sample locations restricted to small area  $G$ .

We first introduce the small area estimator (SMALL), which uses exhaustively computed explanatory variables, and its non-exhaustive version, the pseudo small area estimator (PSMALL).

$$\hat{Y}_{G,SMALL,2p} = \bar{Z}_G^\top \hat{\beta}_{s_2} + \frac{1}{n_{2,G}} \hat{R}(x) \quad (1.8a)$$

$$\hat{Y}_{G,PSMALL,2p} = \hat{\bar{Z}}_G^\top \hat{\beta}_{s_2} + \frac{1}{n_{2,G}} \hat{R}(x) \quad (1.8b)$$

In the equations for the point estimates (Eq. 1.8a and 1.8b), we see that the globally derived regression coefficients are applied to the exhaustively or non-exhaustively calculated means of the explanatory variables ( $\bar{Z}_G$ ,  $\hat{\bar{Z}}_G$ ) which are now only based on the first-phase sample  $s_{1,G}$  located within small area  $G$ . A potential bias of the regression model predictions in the small area  $G$ , due to fitting the regression model with data also outside of  $G$ , is then corrected by adding the mean of the empirical model residuals in  $G$ . This is called the bias or residual correction term.

The package provides the g-weight variance for SMALL and PSMALL respectively (Eq. 1.9a, 1.9b) as well as the external variance (Eq. 1.10a, 1.10b). Again note that all components are restricted to those available at the sample locations in the small area ( $s_{1,G}$  and  $s_{2,G}$ ), with exception of the regression coefficient components  $\hat{\beta}_{s_2}$  and  $\hat{\Sigma}_{\hat{\beta}_{s_2}}$ .

$$\hat{\mathbb{V}}(\hat{Y}_{G,SMALL,2p}) := \bar{Z}_G^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \bar{Z}_G + \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.9a)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,PSMALL,2p}) := \hat{\bar{Z}}_G^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\bar{Z}}_G + \hat{\beta}_{s_2}^\top \hat{\Sigma}_{\hat{\bar{Z}}_G} \hat{\beta}_{s_2} + \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.9b)$$

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{G,SMALL,2p}) := \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.10a)$$

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{G,PSMALL,2p}) := \frac{1}{n_{1,G}} \hat{\mathbb{V}}_{s_{2,G}}(Y(x)) + \left(1 - \frac{n_{2,G}}{n_{1,G}}\right) \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.10b)$$

where  $\hat{\mathbb{V}}_{s_{2,G}}$  indicates taking the sample variance over  $s_{2,G}$ . If boundary adjustment is applied, the simple mean of the explanatory variable vector over the small area  $\hat{\bar{Z}}_G = \frac{1}{n_{1,G}} \sum_{x \in s_{1,G}} Z(x)$

is replaced by its weighted version  $\hat{\bar{\mathbf{Z}}}_G = \frac{\sum_{x \in s_{1,G}} w(x) \mathbf{Z}(x)}{\sum_{x \in s_{1,G}} w(x)}$ , and likewise for exhaustively used auxiliary information.

The synthetic estimator (SYNTH) and pseudo synthetic estimator (PSYNTH) are commonly applied when no terrestrial sample is available within the small area  $G$  (i.e.,  $n_{2,G} = 0$ ). In this case, the point estimates (Eq. 1.11a and 1.11b) are based only on the predictions generated by applying the globally derived regression model to the auxiliary vectors  $\bar{\mathbf{Z}}_G$  and  $\hat{\bar{\mathbf{Z}}}_G$  respectively. However, the bias correction using the observed residuals  $\hat{R}(x)$  is not applied as was the case in the small and pseudo small area estimator (Eq. 1.8a and 1.8b). Thus, the (pseudo) synthetic estimator has a potentially unobservable design-based bias. Also note that the residual variation can no longer be considered in the g-weight variance (Eq. 1.11c and 1.11d). Therefore, the synthetic estimators will usually have a smaller variance than estimators incorporating the regression model uncertainties, but at the cost of a potential bias. Due to the absence of available residuals in  $G$ , there is also no external variance form for the synthetic and pseudo synthetic estimator.

$$\hat{Y}_{G,SYNTH,2p} = \bar{\mathbf{Z}}_G^\top \hat{\beta}_{s_2} \quad (1.11a)$$

$$\hat{Y}_{G,PSYNTH,2p} = \hat{\bar{\mathbf{Z}}}_G^\top \hat{\beta}_{s_2} \quad (1.11b)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,SYNTH,2p}) = \bar{\mathbf{Z}}_G^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \bar{\mathbf{Z}}_G \quad (1.11c)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,PSYNTH,2p}) = \hat{\bar{\mathbf{Z}}}_G^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\bar{\mathbf{Z}}}_G + \hat{\beta}_{s_2}^\top \hat{\Sigma}_{\hat{\bar{\mathbf{Z}}}_G} \hat{\beta}_{s_2} \quad (1.11d)$$

where the variance-covariance matrix of the auxiliary vector  $\hat{\bar{\mathbf{Z}}}_G$  is estimated by

$$\hat{\Sigma}_{\hat{\bar{\mathbf{Z}}}_G} = \frac{1}{n_{1,G}(n_{1,G} - 1)} \sum_{x \in s_{1,G}} (\mathbf{Z}(x) - \hat{\bar{\mathbf{Z}}}_G)(\mathbf{Z}(x) - \hat{\bar{\mathbf{Z}}}_G)^\top \quad (1.12)$$

The synthetic estimators, SYNTH and PSYNTH, have attractively compact formulas but come with the downside of potential bias in their point estimates which can make the variances seem deceptively optimistic. The SMALL and PSMALL estimators overcome this issue by using a bias correction term, i.e.,  $\frac{1}{n_{2,G}} \sum_{x \in s_{2,G}} \hat{R}(x)$ . The motivation behind the extended synthetic and extended pseudo synthetic estimator (EXTSYNTH and EXTPSYNTH) is to use the same mathematically elegant formulas of the (pseudo) synthetic estimators while ensuring that the mean of the empirical prediction model residuals in the entire area  $F$  and the small area  $G$  are by construction both zero at the same time. This is accomplished by extending the vector of auxiliary information  $\mathbf{Z}(x)$  by a binary categorical indicator variable  $I_G(x)$  which takes the value 1 if the sample location  $x$  lies inside the target small area  $G$  and is otherwise set to 0. Recalling that linear models fitted using OLS have zero mean residual property by construction also if categorical variables are used, this leads to unbiased point estimates. The new extended auxiliary vector thus becomes  $\mathbf{Z}^\top(x) = (\mathbf{Z}^\top(x), I_G^\top(x))$  and can be used to replace its non-extended counterpart  $\mathbf{Z}^\top(x)$  wherever it is used in Eq. 1.11 and 1.12. Note that the package functions internally extend the data set by the indicator variable if the EXTSYNTH or EXTPSYNTH estimator is called.

Not every equation needs to be re-written here, but to give an example of the notational change, the regression coefficient under extended model approach becomes

$$\hat{\theta}_{s_2} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \left( \frac{1}{n_2} \sum_{x \in s_2} Y(x) \mathbf{Z}(x) \right) \quad (1.13)$$

The point estimates and their g-weight variances can then be re-written as

$$\hat{Y}_{G,EXTSYNTH,2p} = \bar{\mathbf{Z}}_G^\top \hat{\boldsymbol{\theta}}_{s_2} \quad (1.14a)$$

$$\hat{Y}_{G,EXTPSYNTH,2p} = \hat{\bar{\mathbf{Z}}}_G^\top \hat{\boldsymbol{\theta}}_{s_2} \quad (1.14b)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,EXTSYNTH,2p}) = \bar{\mathbf{Z}}_G^\top \hat{\Sigma}_{\hat{\boldsymbol{\theta}}_{s_2}} \bar{\mathbf{Z}}_G \quad (1.14c)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,EXTPSYNTH,2p}) = \hat{\bar{\mathbf{Z}}}_G^\top \hat{\Sigma}_{\hat{\boldsymbol{\theta}}_{s_2}} \hat{\bar{\mathbf{Z}}}_G + \hat{\boldsymbol{\theta}}_{s_2}^\top \hat{\Sigma}_{\hat{\bar{\mathbf{Z}}}_G} \hat{\boldsymbol{\theta}}_{s_2} \quad (1.14d)$$

While the formulas look similar to the synthetic estimators, note that a decomposition of  $\hat{\boldsymbol{\theta}}_{s_2}$  reveals that the residual correction term is now included in the regression coefficient  $\hat{\boldsymbol{\theta}}_{s_2}$  (Mandallaz et al., 2016) and thus the estimates are asymptotically design-unbiased.

The package also provides the external variance for both the extended synthetic and extended pseudo synthetic estimator. Note that neither the extended model approach nor external variance estimates are possible in the absence of terrestrial samples and thus model residuals in  $G$ , which is precisely when one must rely on the (pseudo) synthetic estimates. The external variance forms of EXTSYNTH and EXTPSYNTH are

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{G,EXTSYNTH,2p}) = \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{\mathbb{R}}(x)) \quad (1.15a)$$

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{G,EXTPSYNTH,2p}) = \frac{1}{n_{1,G}} \hat{\mathbb{V}}_{s_{2,G}}(Y(x)) + \left(1 - \frac{n_{2,G}}{n_{1,G}}\right) \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{\mathbb{R}}(x)) \quad (1.15b)$$

where  $\hat{\mathbb{R}}(x)$  are the empirical residuals under the extended auxiliary vector.

To summarize, the synthetic estimators SYNTH and PSYNTH can be applied whether terrestrial inventory sample is found in the small area or not, but has a deceptively small g-weight variance due to its potential bias. When terrestrial sample is observed in the small area, we can produce (asymptotically) design-unbiased estimates and variances using either SMALL or PSMALL which remove this bias explicitly with a mean residual term, or more elegantly with EXTSYNTH or EXTPSYNTH which simply use the same synthetic formulas while including an indicator variable for the small area in the model formula to remove the bias by construction in OLS.

## Application

Small area estimates in the **forestinventory** package can be applied by specifying the optional argument **small\_area**. The input data set has to include an additional column of class **factor** that describes the small area membership of the sample location represented by that row. The argument **small\_area** requires a **list-object** that comprises

- the name of the column specifying the small area of each observation (**sa.col**).
- a vector specifying the small area(s) for which estimations are desired (**areas**).
- the argument **unbiased** that controls which of the three available estimators is applied.

In order to apply the pseudo small area estimator (PSMALL) with boundary adjustment, we set **unbiased=TRUE** as well as the optional argument **psmall=TRUE**:

```
R> psmall_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B"),
+   unbiased = TRUE), psmall = TRUE, boundary_weights = "boundary_weights")
R> summary(psmall_2p)
```

```
Two-phase small area estimation
```

Call:

```
twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  small_area = list(sa.col = "smallarea", areas = c("A", "B"),
    unbiased = TRUE), boundary_weights = "boundary_weights",
  psmall = TRUE)
```

Method used:

Pseudo small area estimator

Regression Model:

```
tvol ~ mean + stddev + max + q75
```

Estimation results:

	area estimate	ext_variance	g_variance	n1	n2	n1G	n2G	r.squared
A	393.9713	1009.034	1308.117	306	67	94	19	0.6428771
B	419.6416	1214.035	1259.472	306	67	81	17	0.6428771

'boundary\_weight'- option was used to calculate weighted means of auxiliary variables

The small area functions all return an S3 object of class "twophase" with subclass "smallarea". In addition to global estimation, the `estimation` object now comprises the estimates and variances for all small areas (column `area`). We can view the sample sizes by looking into the object itself

```
R> psmall_2p$samplesizes
```

```
$A  
  n1G n2G  n1 n2  
plots 94 19 306 67
```

```
$B  
  n1G n2G  n1 n2  
plots 81 17 306 67
```

The extended pseudo synthetic estimator (EXTPSYNTH) can be applied by setting `unbiased=TRUE` and leaving the optional argument `psmall` to its default value of FALSE:

```
R> extsynth_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B"),
+     unbiased = TRUE), boundary_weights = "boundary_weights")
R> extsynth_2p$estimation

area estimate ext_variance g_variance  n1  n2  n1G  n2G r.squared
1   A 391.9356    995.5602   1017.633 306 67  94  19 0.6526503
2   B 419.7231   1214.6053   1019.191 306 67  81  17 0.6428854
```

The **forestinventory** package automatically includes the indicator variable for the small area behind the scenes so there is no need for the user to implement it. Notice that the  $R^2$ 's (`r.squared`) under the EXTPSYNTH estimator vary between the small areas, while they are identical under the PSMALL estimator. This is because under the EXTPSYNTH estimator, the regression model is recalculated for each small area estimation after adding the indicator variable for the respective small area in the globally derived design matrix. In case of the PSMALL estimator, the regression model stays the same for each small area estimation. Although the results of both estimators should always be close to each other, we recommend applying both estimators and compare the

results afterwards in order to reveal unsuspected patterns in the data, particularly in the case of cluster sampling (see Section 1.6).

Setting the argument `unbiased=FALSE` applies the pseudo synthetic estimator to the selected small areas. Note that in the `grisons` data set, all small areas possess much more than the suggested minimum number of terrestrial observations (a rule of thumb is that  $n_{2,G} \geq 6$ ) required to produce reliable design-unbiased estimates. Hence, choosing to use PSYNTH is probably not desireable and is just applied here for demonstration purposes. In this case the residual correction will not be applied.

```
R> psynth_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B"),
+   unbiased = FALSE), boundary_weights = "boundary_weights")

R> psynth_2p$estimation

  area estimate ext_variance g_variance n1 n2 n1G n2G r.squared
1   A 421.8863          NA 546.8651 306 67  94  19 0.6428771
2   B 418.7399          NA 566.3361 306 67  81  17 0.6428771
```

We see here that the PSYNTH variances are almost only half the variances of the PSMALL and EXTPSYNTH estimator. However, PSMALL and EXTPSYNTH are design unbiased and their variances reflect the fact that they account for potential bias of the regression model predictions. The g-weight variance of PSYNTH completely neglects a potential bias and as a result risks severely overstating the estimation precision.

The exhaustive versions of the small area estimators (Eq. 1.8a, 1.9a, 1.10a, 1.11a, 1.11c) are specified via the optional argument `exhaustive`. Its application requires that we know the exact means of all explanatory variables within the small area(s) of interest. In contrast to the global estimators, the exact means have now to be delivered in the form of a `data.frame`, where each row corresponds to a small area, and each column specifies the exact mean of the respective explanatory variable. Note that likewise the case of global estimation, the order of the explanatory variables in the data frame has to match the order in which they appear in the design matrix defined by the `lm()`-function in *R*. In order to tell *R* which row describes which small area, the row names have to match the respective names of the small areas specified in the `areas` argument.

For the `grisons` data set, the exact means of the explanatory variables for the small areas used in Mandallaz et al. (2013) are thus defined by

```
R> colnames(lm(formula = tvol ~ mean + stddev + max + q75, data = grisons,
+   x = TRUE)$x)

R> true.means.Z.G <- data.frame(Intercept = rep(1, 4),
+   mean = c(12.85, 12.21, 9.33, 10.45),
+   stddev = c(9.31, 9.47, 7.90, 8.36),
+   max = c(34.92, 35.36, 28.81, 30.22),
+   q75 = c(19.77, 19.16, 15.40, 16.91))
R> rownames(true.means.Z.G) <- c("A", "B", "C", "D")

R> true.means.Z.G

  Intercept mean stddev max   q75
A         1 12.85  9.31 34.92 19.77
B         1 12.21  9.47 35.36 19.16
C         1  9.33  7.90 28.81 15.40
D         1 10.45  8.36 30.22 16.91
```

The extended synthetic estimator (EXTSYNTH) can then be applied by

```

R> extsynth_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+     data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+     small_area = list(sa.col ="smallarea", areas = c("A", "B"),
+     unbiased = TRUE), exhaustive = true.means.Z.G)

R> extsynth_2p$estimation
  area estimate ext_variance g_variance n1 n2 n1G n2G r.squared
1   A 372.6930    744.3658   696.5739 Inf 67 Inf 19 0.6526503
2   B 387.5116    693.8576   708.1105 Inf 67 Inf 17 0.6428854

```

Just as in the global case, we see that the variance has again been significantly decreased by substituting the exact auxiliary means and both first phase sample sizes are now infinity. Note that the function extracts the required exact means for small area "A" and "B" from the complete set of exact means defined in `true.means.Z.G`.

## 1.4 Three-phase estimators and their application

### 1.4.1 Global estimators

#### Mathematical background

Solving the sample-based normal equations, the vector of regression coefficients  $\hat{\alpha}_{s_2}$  for the reduced model, i.e., using the reduced set of explanatory variables  $\mathbf{Z}^{(0)}(x)$  available at  $x \in s_0$ , and likewise the vector of regression coefficients  $\hat{\beta}_{s_2}$  for the full model, i.e., using the full set of explanatory variables  $\mathbf{Z}^\top(x) = (\mathbf{Z}^{(0)\top}(x), \mathbf{Z}^{(1)\top}(x))$  available only at a subset  $x \in s_1 \subset s_0$ , are derived as

$$\hat{\alpha}_{s_2} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}^{(0)}(x) \mathbf{Z}^{(0)\top}(x) \right)^{-1} \frac{1}{n_2} \sum_{x \in s_2} Y(x) \mathbf{Z}^{(0)}(x) \quad (1.16a)$$

$$\hat{\beta}_{s_2} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \frac{1}{n_2} \sum_{x \in s_2} Y(x) \mathbf{Z}(x) \quad (1.16b)$$

The package allows for the calculation of point estimates under exhaustive and non-exhaustive use of the auxiliary information in the  $s_0$  phase. Fitting the model using  $s_2$  (i.e., internally) ensures the zero mean residual property over  $s_2$ .

$$\begin{aligned} \hat{Y}_{reg3p,EX} &= \frac{1}{\lambda(F)} \int_F \mathbf{Z}^{(0)\top}(x) \hat{\alpha}_{s_2} + \frac{1}{n_1} \sum_{x \in s_1} (\mathbf{Z}^\top(x) \hat{\beta}_{s_2} - \mathbf{Z}^{(0)\top}(x) \hat{\alpha}_{s_2}) + \frac{1}{n_2} \sum_{x \in s_2} (Y(x) - \mathbf{Z}^\top(x) \hat{\beta}_{s_2}) \\ &= (\bar{\mathbf{Z}}_0^{(0)} - \hat{\bar{\mathbf{Z}}}_1^{(0)})^\top \hat{\alpha}_{s_2} + \hat{\bar{\mathbf{Z}}}_1^\top \hat{\beta}_{s_2} \end{aligned} \quad (1.17a)$$

$$\begin{aligned} \hat{Y}_{reg3p,NEX} &= \frac{1}{n_0} \sum_{x \in s_0} \mathbf{Z}^{(0)\top}(x) \hat{\alpha}_{s_2} + \frac{1}{n_1} \sum_{x \in s_1} (\mathbf{Z}^\top(x) \hat{\beta}_{s_2} - \mathbf{Z}^{(0)\top}(x) \hat{\alpha}_{s_2}) + \frac{1}{n_2} \sum_{x \in s_2} (Y(x) - \mathbf{Z}^\top(x) \hat{\beta}_{s_2}) \\ &= (\hat{\bar{\mathbf{Z}}}_0^{(0)} - \hat{\bar{\mathbf{Z}}}_1^{(0)})^\top \hat{\alpha}_{s_2} + \hat{\bar{\mathbf{Z}}}_1^\top \hat{\beta}_{s_2} \end{aligned} \quad (1.17b)$$

Intuitively, the three phase estimator is simply the mean of the predictions using the reduced model, corrected by the mean difference between the reduced model predictions and the more accurate full model predictions, corrected by the mean difference between the ground truth and the full model predictions. For the compact version of the formula in the non-exhaustive case, the estimated means of  $\mathbf{Z}^{(0)}(x)$  over both the  $s_0$  and  $s_1$  phase, as well as the estimated mean of  $\mathbf{Z}(x)$  over the  $s_1$  phase are calculated according to Eq. 1.18. If the exact mean over  $s_0$  is known, the estimated mean  $\hat{\bar{\mathbf{Z}}}_0^{(0)}$  can simply be replaced by the exact mean  $\bar{\mathbf{Z}}_0^{(0)}$ . Note that in case of applied boundary adjustment (Section 1.3), the simple mean is again replaced by the weighted mean analogous to Eq. 1.7.

$$\hat{\bar{\mathbf{Z}}}_0^{(0)} = \frac{1}{n_0} \sum_{x \in s_0} \mathbf{Z}^{(0)}(x), \quad \hat{\bar{\mathbf{Z}}}_1^{(0)} = \frac{1}{n_1} \sum_{x \in s_1} \mathbf{Z}^{(0)}(x), \quad \hat{\bar{\mathbf{Z}}}_1 = \frac{1}{n_1} \sum_{x \in s_1} \mathbf{Z}(x) \quad (1.18)$$

The package again provides the g-weight and external variances. The g-weight variance formulation is

$$\hat{\mathbb{V}}(\hat{Y}_{reg3p,EX}) = \frac{n_2}{n_1} \bar{\mathbf{Z}}^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \bar{\mathbf{Z}}^{(0)} + \left(1 - \frac{n_2}{n_1}\right) \hat{\bar{\mathbf{Z}}}_1^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\bar{\mathbf{Z}}}_1 \quad (1.19a)$$

$$\hat{\mathbb{V}}(\hat{Y}_{reg3p,NEX}) = \hat{\alpha}_{s_2}^\top \hat{\Sigma}_{\hat{\bar{\mathbf{Z}}}_0^{(0)}} \hat{\alpha}_{s_2} + \frac{n_2}{n_1} \hat{\bar{\mathbf{Z}}}_0^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \hat{\bar{\mathbf{Z}}}_0^{(0)} + \left(1 - \frac{n_2}{n_1}\right) \hat{\bar{\mathbf{Z}}}_1^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\bar{\mathbf{Z}}}_1 \quad (1.19b)$$

with the variance-covariance matrix of  $\hat{\mathbf{Z}}_0^{(0)}$  and the variance-covariance matrices of the regression coefficients  $\hat{\boldsymbol{\alpha}}_{s_2}$  and  $\hat{\boldsymbol{\beta}}_{s_2}$ :

$$\hat{\Sigma}_{\hat{\mathbf{Z}}_0^{(0)}} = \frac{1}{n_0(n_0 - 1)} \sum_{x \in s_0} (\mathbf{Z}^{(0)}(x) - \hat{\mathbf{Z}}_0^{(0)})(\mathbf{Z}^{(0)}(x) - \hat{\mathbf{Z}}_0^{(0)})^\top \quad (1.20a)$$

$$\hat{\Sigma}_{\hat{\boldsymbol{\alpha}}_{s_2}} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}^{(0)}(x) \mathbf{Z}^{(0)\top}(x) \right)^{-1} \left( \frac{1}{n_2^2} \sum_{x \in s_2} \hat{R}^{(0)2}(x) \mathbf{Z}^{(0)}(x) \mathbf{Z}^{(0)\top}(x) \right) \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}^{(0)}(x) \mathbf{Z}^{(0)\top}(x) \right)^{-1} \quad (1.20b)$$

$$\hat{\Sigma}_{\hat{\boldsymbol{\beta}}_{s_2}} = \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \left( \frac{1}{n_2^2} \sum_{x \in s_2} \hat{R}^2(x) \mathbf{Z}(x) \mathbf{Z}^\top(x) \right) \left( \frac{1}{n_2} \sum_{x \in s_2} \mathbf{Z}(x) \mathbf{Z}^\top(x) \right)^{-1} \quad (1.20c)$$

Note that  $\hat{R}(x) = Y(x) - \mathbf{Z}^\top(x) \hat{\boldsymbol{\beta}}_{s_2}$  denotes the empirical residuals of the full model, whereas  $\hat{R}^{(0)}(x) = Y(x) - \mathbf{Z}^{(0)\top} \hat{\boldsymbol{\alpha}}_{s_2}$  denotes the empirical residuals of the reduced model. The external variance form under linear regression models is defined as

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{reg3p,EX}) = \frac{1}{n_1} \hat{\mathbb{V}}_{s_2}(\hat{R}^{(0)}(x)) + (1 - \frac{n_2}{n_1}) \frac{1}{n_2} \hat{\mathbb{V}}_{s_2}(\hat{R}(x)) \quad (1.21a)$$

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{reg3p,NEX}) = \frac{1}{n_0} \hat{\mathbb{V}}_{s_0}(\hat{Y}^{(0)}(x)) + \frac{1}{n_1} \hat{\mathbb{V}}_{s_2}(\hat{R}^{(0)}(x)) + (1 - \frac{n_2}{n_1}) \frac{1}{n_2} \hat{\mathbb{V}}_{s_2}(\hat{R}(x)) \quad (1.21b)$$

where  $\hat{\mathbb{V}}_{s_0}$  indicates taking the sample variance over  $s_0$ .

## Application

In order to demonstrate the three-phase estimators in the package, we created an artificial three-phase scenario by recoding the phase indicators in the **grisons** data set (column **phase\_id\_3p**) according to the terminology used in this article (0 for  $s_0$ , 1 for  $s_1$ , 2 for  $s_2$ ). We now assume that the mean canopy height (**mean**) is available at all 306 sample locations  $x \in s_0$ , whereas we have the explanatory variables **stddev**, **max** and **q75** only at 128 subsamples  $s_1$  of  $s_0$ . At 40 further subsamples  $s_2$  we have the observations  $Y(x)$  from the field inventory. Based on this setup, we can now define the reduced and full regression model formulas to be used in the **threephase()** function (note that the models are nested):

```
R> formula.rm <- tvol ~ mean
R> formula.fm <- tvol ~ mean + stddev + max + q75
```

Compared to the **twophase()**-function, we now have to specify two regression models, i.e., the nested reduced (**formula.s0**) and full (**formula.s1**) regression model. In addition, we also have to specify the indication of the  $s_1$  phase (**s1.id**) in the argument **phase\_id** (note that **forestinventory** implicitly assumes that all other rows in the input data set belong to  $s_0$ ). The global three-phase estimation can thus be applied by

```
R> reg3p_nex <- threephase(formula.s0 = formula.rm, formula.s1 = formula.fm,
+   data = grisons, phase_id = list(phase.col = "phase_id_3p", s1.id = 1,
+   terrgrid.id = 2), boundary_weights = "boundary_weights")
R> summary(reg3p_nex)
Three-phase global estimation
```

```
Call:
threephase(formula.s0 = formula.rm, formula.s1 = formula.fm,
```

```

data = grisons, phase_id = list(phase.col = "phase_id_3p",
                                s1.id = 1, terrgrid.id = 2), boundary_weights = "boundary_weights")

Method used:
Non-exhaustive global estimator

Full Regression Model:
tvol ~ mean + stddev + max + q75

Reduced Regression Model:
tvol ~ mean

Estimation results:
estimate ext_variance g_variance n0 n1 n2 r.squared_reduced
372.0896    454.4064   451.3626 306 128 40          0.527363
r.squared_full
0.7166608

'boundary_weight'- option was used to calculate weighted means of auxiliary variables

```

The `summary()` of a `threephase()`-function now recalls both regression model formulas and also gives the  $R^2$  for both the reduced (`r.squared_reduced`) and the full (`r.squared_full`) models. We are told that including `stddev`, `max` and `q75` yields a 20 % improvement in  $R^2$ . When comparing to using only `mean` under a two-phase approach, we would see a considerable reduction in variance by the three-phase extension.

## 1.4.2 Small area estimators

### Mathematical background

The three two-phase small area estimators described in Section 1.3.2 can also be extended to the three-phase scenario. The general principle thereby stays the same, i.e., the regression coefficients of the reduced and full model and their variance-covariance matrices are calculated on the entire  $s_2$  sample according to Eq. 1.16a, 1.16b, 1.20b and 1.20c, and are subsequently used to make predictions for sample locations restricted to small area  $G$ .

The unbiased point estimates of the SMALL and PSMALL estimator are calculated by applying the globally derived reduced and full regression model coefficients to the small area means of the explanatory variables, and then corrected for a potential model bias in  $G$  by adding the small area mean of the full model residuals, i.e.,  $\hat{R}_G(x) = Y_G(x) - \bar{\mathbf{Z}}_G^\top(x)\hat{\boldsymbol{\beta}}_{s_2}$ , to the point estimate. The difference between the mean  $\hat{\bar{\mathbf{Z}}}_{1,G}^{(0)}$  and the more precise or exact mean  $\hat{\bar{\mathbf{Z}}}_{0,G}^{(0)}$  and  $\bar{\mathbf{Z}}_{0,G}^{(0)}$  is again considered as a correction term likewise in the global estimation (Eq. 1.17).

$$\hat{Y}_{G,SMALL,3p} = (\bar{\mathbf{Z}}_{0,G}^{(0)} - \hat{\bar{\mathbf{Z}}}_{1,G}^{(0)})^\top \hat{\boldsymbol{\alpha}}_{s_2} + \hat{\bar{\mathbf{Z}}}_{1,G}^\top \hat{\boldsymbol{\beta}}_{s_2} + \frac{1}{n_{2,G}} \hat{R}_G(x) \quad (1.22a)$$

$$\hat{Y}_{G,PSMALL,3p} = (\hat{\bar{\mathbf{Z}}}_{0,G}^{(0)} - \hat{\bar{\mathbf{Z}}}_{1,G}^{(0)})^\top \hat{\boldsymbol{\alpha}}_{s_2} + \hat{\bar{\mathbf{Z}}}_{1,G}^\top \hat{\boldsymbol{\beta}}_{s_2} + \frac{1}{n_{2,G}} \hat{R}_G(x) \quad (1.22b)$$

The g-weight variance is then calculated as

$$\hat{\mathbb{V}}(\hat{Y}_{G,SMALL,3p}) = \frac{n_2}{n_1} \bar{\mathbf{Z}}_{0,G}^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \bar{\mathbf{Z}}_{0,G}^{(0)} + (1 - \frac{n_2}{n_1}) \hat{\mathbf{Z}}_{1,G}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\mathbf{Z}}_{1,G} + \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.23a)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,PSMALL,3p}) = \hat{\alpha}_{s_2}^\top \hat{\Sigma}_{\hat{\mathbf{Z}}_{0,G}^{(0)}} \hat{\alpha}_{s_2} + \frac{n_2}{n_1} \hat{\mathbf{Z}}_{0,G}^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \hat{\mathbf{Z}}_{0,G}^{(0)} + (1 - \frac{n_2}{n_1}) \hat{\mathbf{Z}}_{1,G}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\mathbf{Z}}_{1,G} + \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.23b)$$

with the variance-covariance matrix

$$\hat{\Sigma}_{\hat{\mathbf{Z}}_{0,G}^{(0)}} = \frac{1}{n_{0,G}(n_{0,G}-1)} \sum_{x \in s_{0,G}} (\mathbf{Z}^{(0)}(x) - \hat{\mathbf{Z}}_{0,G}^{(0)}) (\mathbf{Z}^{(0)}(x) - \hat{\mathbf{Z}}_{0,G}^{(0)})^\top \quad (1.24)$$

The external variance is defined as:

$$\hat{\mathbb{V}}_{ext}(\hat{Y}_{G,SMALL,3p}) = \frac{1}{n_{1,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}^{(0)}(x)) + (1 - \frac{n_{2,G}}{n_{1,G}}) \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \quad (1.25a)$$

$$\begin{aligned} \hat{\mathbb{V}}_{ext}(\hat{Y}_{G,PSMALL,3p}) &= \frac{1}{n_{0,G}} \hat{\mathbb{V}}_{s_{2,G}}(Y(x)) + (1 - \frac{n_{1,G}}{n_{0,G}}) \frac{1}{n_{1,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}^{(0)}(x)) \\ &\quad + (1 - \frac{n_{2,G}}{n_{1,G}}) \frac{1}{n_{2,G}} \hat{\mathbb{V}}_{s_{2,G}}(\hat{R}(x)) \end{aligned} \quad (1.25b)$$

where  $\hat{R}^{(0)}(x) = Y(x) - \hat{Y}^{(0)}(x)$  with  $\hat{Y}^{(0)}(x) = \mathbf{Z}^{(0)\top}(x) \hat{\alpha}_{s_2}$ .

The synthetic (SYNTH) and pseudo synthetic (PSYNTH) estimator can be applied if no terrestrial samples are available in the small area, i.e.,  $n_{2,G} = 0$ . Consequently, the residual correction and the residual variation term of the full model can no longer be applied and drops from the point estimate (Eq. 1.26a and 1.26b) and g-weight variance (Eq. 1.26c and 1.26d) formulas. The point estimates are again potentially biased since  $\frac{1}{n_{2,G}} \sum_{x \in s_{2,G}} \hat{R}(x) = 0$  for the full model residuals can not be ensured within small area  $G$ . Also the variance will be small but to the cost of ignoring the model uncertainties. Note that there is again no external variance formula for the synthetic and pseudo synthetic estimation.

$$\hat{Y}_{G,SYNTNH,3p} = (\bar{\mathbf{Z}}_{0,G}^{(0)} - \hat{\mathbf{Z}}_{1,G}^{(0)})^\top \hat{\alpha}_2 + \hat{\mathbf{Z}}_{1,G}^\top \hat{\beta}_{s_2} \quad (1.26a)$$

$$\hat{Y}_{G,PSYNTNH,3p} = (\hat{\mathbf{Z}}_{0,G}^{(0)} - \hat{\mathbf{Z}}_{1,G}^{(0)})^\top \hat{\alpha}_2 + \hat{\mathbf{Z}}_{1,G}^\top \hat{\beta}_{s_2} \quad (1.26b)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,SYNTNH,3p}) = \frac{n_2}{n_1} \bar{\mathbf{Z}}_{0,G}^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \bar{\mathbf{Z}}_{0,G}^{(0)} + (1 - \frac{n_2}{n_1}) \hat{\mathbf{Z}}_{1,G}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\mathbf{Z}}_{1,G} \quad (1.26c)$$

$$\hat{\mathbb{V}}(\hat{Y}_{G,PSYNTNH,3p}) = \hat{\alpha}_2^\top \hat{\Sigma}_{\hat{\mathbf{Z}}_{0,G}^{(0)}} \hat{\alpha}_2 + \frac{n_2}{n_1} \hat{\mathbf{Z}}_{0,G}^{(0)\top} \hat{\Sigma}_{\hat{\alpha}_{s_2}} \hat{\mathbf{Z}}_{0,G}^{(0)} + (1 - \frac{n_2}{n_1}) \hat{\mathbf{Z}}_{1,G}^\top \hat{\Sigma}_{\hat{\beta}_{s_2}} \hat{\mathbf{Z}}_{1,G} \quad (1.26d)$$

The extended synthetic (EXTSYNTH) and extended pseudo synthetic (EXTPSYNTH) estimator ensures that the residuals of the full model over both the entire inventory area  $F$  and the small area  $G$  are zero at the same time, i.e.,  $\frac{1}{n_2} \sum_{x \in s_2} \hat{R}(x) = \frac{1}{n_{2,G}} \sum_{x \in s_{2,G}} \hat{R}(x) = 0$ . This is again realized by extending the vector of explanatory variables by a binary categorical indicator variable  $I_G(x)$  which takes the value 1 if the observation lies inside the small area  $G$  and is otherwise set to 0. The extended auxiliary vector is thus defined as  $\mathbf{Z}^\top(x) = (\mathbf{Z}^{(0)\top}(x), \mathbf{Z}^{(1)\top}(x))$ , where  $\mathbf{Z}^{(0)\top}(x) = (\mathbf{Z}^{(0)\top}(x), I_G^\top(x))$ . In other words, when the extended option is chosen, **forestinventory** automatically adds the binary indicator variable for the desired small area for all observations in the input data frame (i.e.,  $s_0$ ). The regression coefficients, point estimates and variance estimates are calculated by replacing  $\mathbf{Z}$  with  $\mathbf{Z}$  (and likewise  $\mathbf{Z}^{(0)}$  with  $\mathbf{Z}^{(0)}$ ) into Eq. 1.16, 1.20, 1.25 and

1.26. Just as in the two-phase case, the resulting point estimates are now unbiased and have an associated g-weight variance that accounts for the variability of the regression coefficients resulting from the random sample  $s_2$ .

## Application

We will demonstrate the use of three-phase small area estimation in the package **forestinventory** by applying the EXTSYNTH and SYNTH estimator to the **grisons** data set. The setup is thus exactly the same as in the example for global three-phase estimation (Section 1.4.1). However, this time will use the exact auxiliary mean of the mean canopy height variable (**mean**) and assume that we do not know the exact means of the remaining explanatory variables **stddev**, **max** and **q75**. We thus first define the true means for each small area just as we did in the **twophase()** example (Section 1.3.2):

```
R> truemmeans.G <- data.frame(Intercept = rep(1, 4),
+     mean = c(12.85, 12.21, 9.33, 10.45))
R> rownames(truemmeans.G) <- c("A", "B", "C", "D")
```

Three-phase small area estimation in the package can in general be applied by additionally specifying the **small\_area** list argument. The exhaustive estimators can be called by optionally passing a **data.frame** containing the exact auxiliary means to the **exhaustive** argument. The EXTSYNTH estimator can be applied by setting the argument **unbiased** to TRUE (default):

```
R> extsynth_3p <- threephase(formula.rm, formula.fm, data = grisons,
+     phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+     small_area = list(sa.col = "smallarea", areas = c("A", "B"), unbiased = TRUE),
+     exhaustive = truemmeans.G, boundary_weights = "boundary_weights")

R> extsynth_3p$estimation

  area estimate ext_variance g_variance n0  n1  n2  n0G  n1G  n2G
1   A 382.6405    1642.055  1518.741 Inf 128 40 Inf  38  12
2   B 368.9013    1501.211  1530.576 Inf 128 40 Inf  34  11
  r.squared_reduced r.squared_full
1           0.5454824      0.7242913
2           0.5354637      0.7171512
```

The SYNTH estimator can be applied by changing the argument **unbiased** to FALSE, which causes the function to not apply a bias correction in the respective small area.

```
R> synth_3p <- threephase(formula.rm, formula.fm, data = grisons,
+     phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+     small_area = list(sa.col = "smallarea", areas = c("A", "B"), unbiased = FALSE),
+     exhaustive = truemmeans.G, boundary_weights = "boundary_weights")

R> synth_3p$estimation

  area estimate ext_variance g_variance n0  n1  n2  n0G  n1G  n2G
1   A 409.3390        NA  410.7529 Inf 128 40 Inf  38  12
2   B 375.4608        NA  461.8250 Inf 128 40 Inf  34  11
  r.squared_reduced r.squared_full
1           0.527363      0.7166608
2           0.527363      0.7166608
```

We see that the **threephase()**-function returns the sample sizes in the entire inventory area as well as within each small area. The value **Inf** for **n0G** indicates that the explanatory variables at  $s_0$  sample locations used in the reduced model were in our case derived exhaustively. If we compare

the two results, we see that the SYNTH estimation again yields a much smaller variance than the EXTSYNTH estimation, but at the cost of a potential bias.

We can also analyse how the exhaustive derivation of `mean` performed compared to the case where `mean` is non-exhaustively available but at a very large  $s_0$  phase with  $n_{0,G} \gg n_{1,G}$ . To do this, we additionally compute the EXTPSYNTH estimates. As we can see, the exhaustive derivation of `mean` only yielded a slightly smaller variance.

```
R> extsynth_3p <- threephase(formula.rm, formula.fm, data = grisons,
+   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B"), unbiased = TRUE),
+   boundary_weights = "boundary_weights")

R> extsynth_3p$estimation

  area estimate ext_variance g_variance n0  n1  n2  n0G  n1G  n2G
1   A 395.1882    1901.211  1858.204 306 128  40   94   38   12
2   B 389.8329    1846.995  1816.655 306 128  40   81   34   11
  r.squared_reduced r.squared_full
1           0.5454824      0.7242913
2           0.5354637      0.7171512
```

## 1.5 Calculation of confidence intervals

Converting the estimated variance into a 95% confidence interval (CI) allows for a more practical interpretation of a point estimate's precision. The correct interpretation of a CI is not that there is a 95% probability that it contains the true value. In the design-based context, the true value of the population parameter we are trying to estimate, albeit unknown, is fixed and the sample is randomly generated under the sample design. Theoretically, if we were to repeatedly conduct the inventory using the same estimation method, estimator and auxiliary information under newly drawn random samples and calculate the 95% CI from each sample, then 95% of the CIs are expected to contain the true population parameter. The confidence level  $1 - \alpha$  (e.g., 95%) is thus the expected frequency or proportion of possible confidence intervals to contain the unknown population parameter under resampling and is therefore often also referred to as coverage rate. The CI is also linked to hypothesis testing in that its associated point estimate is considered statistically different from any given value that lies outside the CI boundaries.

Based on the central limit theorem it can be assumed that under hypothetical repeated sampling the point estimates will asymptotically follow a normal distribution. However, on the recommendation of Mandallaz (2013a), better confidence intervals can be obtained using the Student's  $t$  distribution defined as

One-phase estimation:

$$CI_{1-\alpha}(\hat{Y}) = \left[ \hat{Y} - t_{n_2-1,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})}, \hat{Y} + t_{n_2-1,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})} \right] \quad (1.27)$$

Two-phase and three-phase global estimation:

$$CI_{1-\alpha}(\hat{Y}) = \left[ \hat{Y} - t_{n_2-p,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})}, \hat{Y} + t_{n_2-p,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})} \right] \quad (1.28)$$

Two-phase and three-phase small area estimation:

$$CI_{1-\alpha}(\hat{Y}) = \left[ \hat{Y} - t_{n_{2,G}-1,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})}, \hat{Y} + t_{n_{2,G}-1,1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{Y})} \right] \quad (1.29)$$

where  $\hat{Y}$  is the point estimate,  $\hat{V}(\hat{Y})$  is the estimated variance,  $1 - \alpha$  is the confidence level and  $p$  constitutes the number of parameters used in the (full) regression model. In case of cluster-sampling,  $n_{2,G}$  is the number of terrestrial clusters (a cluster constitutes the sample unit and comprises multiple sample plots). Note that in case of synthetic estimations (SYNTH, PSYNTH), the degrees of freedom are  $n_2 - p$  as is the case for global estimation. In **forestinventory**, the confidence intervals for all estimation methods and estimators can be computed by the S3 generic method **confint()**, which requires an estimation object created by either the **onephase()**, **twophase()** or **threephase()** function. For example, the 95% confidence interval for the small area estimates by the EXTPSYNTH estimator (Section 1.3.2) are calculated by:

```
R> confint(extpsynth_2p)

95% Confidence Intervals for twophase small area estimation

area estimate ci_lower_ext ci_upper_ext ci_lower_g ci_upper_g
1   A 391.9356    325.6463    458.2250    324.9155    458.9558
2   B 419.7231    345.8418    493.6043    352.0456    487.4006
```

## 1.6 Special cases and scenarios

### 1.6.1 Post-stratification

A special case of multi-phase regression estimation is post-stratification, which can further be divided into the cases of multi-phase sampling for stratification and multi-phase sampling for regression within strata. Both imply the use of one or more categorical variables in the regression model(s), leading to classical ANOVA and ANCOVA models.

To demonstrate post-stratification, we first create an artificial categorical variable development stage (`stage`) by clustering the mean canopy heights of the `grisons` data set into 3 height classes:

```
R> grisons$stage <- as.factor(kmeans(grisons$mean, centers = 3)$cluster)
```

Two-phase sampling for stratification is applied if the model only contains categorical variables, in this case the factor variable `stage`. Linear regression models only fitted with categorical variables produce ANOVA models, which when used in multi-phase regression estimators, is equivalent to post-stratification. For our example, this means that the model predictions are simply the means of the terrestrial response values within each development stage (within-strata means).

```
R> twophase(formula = tvol ~ stage, data = grisons,
+    phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+    boundary_weights = "boundary_weights")
```

Two-phase sampling for regression within strata implies the combination of continuous and categorical variables within the model (i.e., we have an ANCOVA model). If an interaction term is not present between categorical and continuous variables, the regression lines within the strata will have the same slope but different intercepts. If an interaction term is present, both the intercept and the slope are allowed to vary within the strata. Note that one can actually use the entire range of OLS regression techniques in the multi-phase estimators, including higher order terms and transformations of the explanatory variables, which makes them very flexible.

```
R> twophase(formula=tvol ~ mean + stddev + max + q75 + stage, data = grisons,
+    phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+    boundary_weights = "boundary_weights")
```

The variance of all design-based estimators included in `forestinventory` can be decreased by reducing the sum of squared residuals of the regression model. In case of post-stratification, this particularly implies minimizing the within strata residual square sum. Also, for post-stratification, the g-weight variance should be trusted over the external variance because it has the advantage that the strata weights are estimated from the large sample rather than the terrestrial sample  $s_2$ .

### 1.6.2 Small area estimation under cluster sampling

As mentioned in Section 1.2.6, cluster sampling is a special case of sample designs where the sample consists of more than one spatially agglomerated sample points. One randomly places the sample location  $x$  in the inventory area as in the simple sampling design, but then  $M - 1$  additional sample locations  $x_2, \dots, x_M$  are created close to the cluster origin  $x$  by adding a fixed set of spatial vectors  $e_2, \dots, e_M$  to  $x$ . The idea of cluster sampling is to increase the amount of information without increasing the travel costs of the terrestrial campaign. However, the information gathered at all sub-locations of a cluster is then averaged on the cluster level, and this average value then references exactly one point, i.e., the cluster origin  $x$ . Without going into too much mathematical detail, the estimators under simple sampling are thus extended in a way that all parameters (local density, mean vector of explanatory variables, mean model residuals) have to be calculated as the weighted cluster means with  $M(x)$  being the cluster weights. Whereas the geometric form and the number of sample locations per cluster  $M$  is fixed (i.e., defined by the inventorist), the actual

number of points  $M(x)$  falling into the forest area  $F$  at sample location  $x$  is random because the cluster origin  $x$  is random. The **forstinventory** package identifies clusters via a unique cluster ID that is assigned to a column in the input data set. Its column name is passed to the argument `cluster` in the `twophase()` and `threephase()` function calls.

For small area applications, the scenario might occur where the points of a cluster at sample locations  $x$  spread over more than one small area, i.e., only a subset  $M_G(x) < M(x)$  is included in the small area of interest. In this case, the zero mean residual property within the small area,  $\frac{\sum_{x \in s_{2,G}} M(x) \hat{R}_c(x)}{\sum_{x \in s_{2,G}} M(x)} = 0$ , is no longer guaranteed when using the extended and pseudo extended synthetic estimator (see EXTSYNTH and EXTPSYNTH in Sections 1.3.2 and 1.4.2). In this case, it is adviseable to use the (pseudo) small area estimator (SMALL or PSMALL) where the zero mean residual property is still ensured.

In order to keep track of such cases, **forstinventory** tells the user to do so by returning a warning message:

```
R> extsynth.clust <- twophase(formula = basal ~ stade + couver + melange,
+   data = zberg, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   cluster = "cluster", small_area = list(sa.col = "ismallold", areas = c("1"),
+   unbiased = TRUE))
Warning message:
At least one terrestrial cluster not entirely included within small area 1.
Zero mean residual assumption for small area maybe violated.
Check mean_Rc_x_hat_G and consider alternative estimator 'psmall'

R> psmall.clust <- twophase(formula = basal ~ stade + couver + melange,
+   data = zberg, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   cluster = "cluster", small_area = list(sa.col = "ismallold", areas = c("1"),
+   unbiased = TRUE), psmall = TRUE)

R> extsynth.clust$estimation
  area estimate ext_variance g_variance n1 n2 n1G n2G r.squared
1    1 25.54748     14.03806   14.16853 298 73  29    8  0.205741

R> psmall.clust$estimation
  area estimate ext_variance g_variance n1 n2 n1G n2G r.squared
1    1 23.98581     16.30509   15.69473 298 73  29    8  0.1873795
```

Comparing the EXTPSYNTH and PSMALL estimates, we see that in this particular case the point estimates are close and more important, the external as well as the g-weight variances only differ marginally. This can be taken as evidence that the violation of the zero mean residual property can here be expected to have negligible consequences.

### 1.6.3 Violation of nesting in sample design

As explained in Section 1.2, a basic prerequisite for the application of multi-phase estimators is that the sample phases ( $s_0, s_1, s_2$ ) are nested in each other. The correct nesting thereby concerns the spatial arrangement of the sample phases (Fig. 1.2a) as well as the availability of terrestrial and auxiliary information per phase and sample location. For the latter, **forstinventory** runs validity checks in the background, provides warning and error messages and, if possible, applies first-aid adjustments to the inventory data set to prevent the calculations from failing. We will demonstrate possible nesting violations by applying the global three-phase estimator to the `grisons` and `zberg` data sets.

## Violation 1

Based on the nesting rule,  $s_2 \in s_1 \in s_0$ , each  $s_2$  and  $s_1$  sample location must have all explanatory variables available that are used in the full (and thus reduced) regression model. If e.g., an  $s_2$  and/or  $s_1$  point misses a variable which is used in the full and reduced model (in this case `mean`), the function will delete this sample point from the dataset and produce the following messages:

```
R> grisons[which(grisons$phase_id_3p == 2)[1], "mean"] <- NA
R> threephase(formula.s0 = tvol ~ mean,
+   formula.s1 = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   boundary_weights = "boundary_weights")
Warning messages:
1: In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  Sample design not nested: for 1 terrestrial plots at least one auxiliary
  parameter of the first phase (s1) is missing
2: In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  Sample design not nested: for 1 terrestrial plots at least one auxiliary
  parameter of the zero phase (s0) is missing
3: In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  1 rows deleted due to missingness in the set of auxiliary parameters for the
  zero phase (s0) (1 terrestrial plots affected by deletion)
```

## Violation 2

However, if an  $s_2$  and/or  $s_1$  point is missing a variable which is only used in the full regression model (in this example `q75`), the function will recode the phase indicator of that point to  $s_0$ , since the point still provides the required information for the reduced model. If this concerns an  $s_2$  sample location, the associated value of the response variable can no longer be used.

```
R> grisons[which(grisons$phase_id_3p == 2)[1], "q75"] <- NA
R> threephase(formula.s0 = tvol ~ mean,
+   formula.s1 = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col="phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   boundary_weights = "boundary_weights")
Warning messages:
1: In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  Sample design not nested: for 1 terrestrial plots at least one auxiliary
  parameter of the first phase (s1) is missing
2: In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  Changed the phase_id for 1 rows to the zero phase (s0) due to missingness in
  the set of auxiliary parameters for the first phase (s1) (1 terrestrial
  information no longer usable by this change)
```

## Violation 3

If an  $s_0$  point misses at least one of the explanatory variables used in the reduced model, the sample locations are deleted from the data set.

```
R> grisons[which(grisons$phase_id_3p == 0)[1], "mean"] <- NA
R> threephase(formula.s0 = tvol ~ mean,
+   formula.s1 = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   boundary_weights = "boundary_weights")
```

```

Warning message:
In threephase(formula.s0 = tvol ~ mean, formula.s1 = tvol ~ mean + :
  1 rows deleted due to missingness in the set of auxiliary parameters for the
  zero phase (s0) (0 terrestrial plots affected by deletion)

```

Note that all the automatic data adjustments (deletion, recoding) have to be accepted with caution. Recapitulating, the unbiasedness of estimators in the design-based framework is based on the uniform and independent randomization of the sample locations. This means that every possible location within the forest area  $F$ , as well as pairs of locations, have inclusion and joint inclusion probabilities greater than zero. Whereas this is already violated in practice by the use of regular grids, one can still expect that these grids do not exclude specific forest structures. If any information should be missing at the sample locations, one should clarify the reason for this and make sure that the information can reasonably be assumed to be completely missing at random.

#### **Violation 4**

If a categorical variable is used in the regression model(s) and the terrestrial sample  $s_2$  is considerably small compared to the  $s_1$  phase, it might occur that a category is only present in the  $s_1 \setminus s_2$  sample, and thus missing in the  $s_2$  sample. Therefore, an internal regression model cannot be calculated and the function stops with the following error message:

```

R> zberg <- zberg[-which(zberg.n$phase_id_2p == 2 & zberg.n$stade == "300"), ]
R> twophase(formula = basal ~ stade + couver + melange, data = zberg,
+   phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   cluster = "cluster")
Error in check.mandatoryInputs(formula, data, phase_id) :
  Level '300' of factor variable 'stade' existing in s1(s0)-but not in s2 sample.
  Calculation of coefficient not feasible.

```

## 1.7 Analysis and visualization

### 1.7.1 Analysis

We often want to compare the results and performances of different estimation methods and estimators for a given global or small area inventory, which can be easily accomplished in **forestinventory** using the `estTable()` function. This function restructures the results from the `onephase()`, `twophase()` and `threephase()` objects and merges them into one single data set that provides the basis for further analysis. For demonstration purposes, we will first recalculate the one-phase estimator as well as the two-phase and three-phase EXTPSYNTH and PSYNTH estimator for the `grisons` data set:

```
R> op <- onephase(formula = tvol~1, data = grisons,
+   phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D")))

R> extsynth_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
+   unbiased = TRUE), boundary_weights = "boundary_weights")

R> psynth_2p <- twophase(formula = tvol ~ mean + stddev + max + q75,
+   data = grisons, phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
+   unbiased = FALSE), boundary_weights = "boundary_weights")

R> extsynth_3p <- threephase(formula.s0 = tvol ~ mean,
+   formula.s1 = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
+   unbiased = TRUE), boundary_weights = "boundary_weights")

R> psynth_3p <- threephase(formula.s0 = tvol ~ mean,
+   formula.s1 = tvol ~ mean + stddev + max + q75, data = grisons,
+   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
+   small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
+   unbiased = FALSE), boundary_weights = "boundary_weights")
```

We can then aggregate all estimation objects in a `list` and pass it to the `estTable()`-function:

```
R> grisons.sae.table <- estTable(est.list = list(op, extsynth_2p, psynth_2p,
+   extsynth_3p, psynth_3p), sae = TRUE, vartypes = c("variance", "g_variance",
+   "ext_variance"))
```

The function merges the estimation results and returns a `list` object with the subclasses "`esttable`" "`smallarea`". The `vartypes` argument can be used to restrict the `estTable()` output to certain types of variances. If one prefers the `data.frame` format for further analysis, this can easily be done using `as.data.frame(grisons.sae.table)`. Note however that **forestinventory** provides several S3 generic methods specifically for the class "`esttable`".

The structure of an `esttable` object is very similar to the objects created by the small area estimation functions of the package. However, the point estimates and variances from all estimation objects passed to `estTable()` have been stored in one single column (`estimate` and `variance`) and can be distinguished by the variables `method`, `estimator` and `vartype` which specify the estimation method (one, two or three-phase), the estimator and the type of variance that was applied (`g_` for *g-weight* and `ext_` for external variance). By default, the confidence intervals are also added.

```
R> str(grisons.sae.table)
```

```
List of 20
$ area           : chr [1:28] "A" "A" "A" "A" ...
$ domain         : Factor w/ 2 levels "global","smallarea": 1 2 2 2 2 2 2 1 2 ...
$ method          : Factor w/ 3 levels "onephase","twophase",...: 1 3 3 3 2 2 2 ...
$ estimator       : Factor w/ 3 levels "onephase","psynth extended",...: 1 2 2 3 ...
$ vartype         : Factor w/ 3 levels "ext_variance",...: 3 1 2 2 1 2 2 3 1 2 ...
$ estimate        : num [1:28] 410 395 395 422 392 ...
$ variance        : num [1:28] 1987 1901 1858 726 996 ...
$ std             : num [1:28] 44.6 43.6 43.1 26.9 31.6 ...
$ error            : num [1:28] 10.86 11.03 10.91 6.39 8.05 ...
$ n2              : num [1:28] 19 40 40 40 67 67 67 17 40 40 ...
$ n2G             : num [1:28] NA 12 12 12 19 19 19 NA 11 11 ...
$ n1              : num [1:28] NA 128 128 128 306 306 306 NA 128 128 ...
$ n1G             : num [1:28] NA 38 38 38 94 94 94 NA 34 34 ...
$ n0              : int [1:28] NA 306 306 306 NA NA NA NA 306 306 ...
$ n0G             : int [1:28] NA 94 94 94 NA NA NA NA 81 81 ...
$ r.squared        : num [1:28] NA NA NA NA 0.653 ...
$ r.squared_reduced: num [1:28] NA 0.545 0.545 0.527 NA ...
$ r.squared_full   : num [1:28] NA 0.724 0.724 0.717 NA ...
$ ci_lower         : num [1:28] 317 299 300 367 326 ...
$ ci_upper         : num [1:28] 504 491 490 476 458 ...
- attr(*, "row.names")= int [1:28] 1 2 3 4 5 6 7 8 9 10 ...
- attr(*, "class")= chr [1:3] "list" "esttable" "smallarea"
```

Note that `estTable()` also returns the estimation error (`error`) that is defined as the standard error devideed by the point estimate:

$$error[\%] = \frac{\sqrt{\hat{V}(\hat{Y})}}{\hat{Y}} \cdot 100 \quad (1.30)$$

As multi-phase estimation techniques are primary intended to increase estimation precision, the function `mphase.gain()` can be applied to quantify the potential benefit of a multi-phase global or small area estimate compared to its respective one-phase estimate. The function takes an `esttable` object as input and returns a summary of which multi-phase method and estimator performed best using the precision from the one-phase estimator as a baseline. If the `esttable` object contains more than one multi-phase estimation object, `mphase.gain()` identifies the one with the smallest variance and compares it to the `onephase` estimation. The argument `pref.vartype` can be used to define what type of variance (g-weight or external) should be used for the comparison. Synthetic estimates (SYNTH and PSYNTH estimator) are not considered for the comparison under the default setting (`exclude.synth = TRUE`) since they usually have a much smaller variance at the cost of a potential bias.

```
R> mphase.gain(grisons.sae.table, pref.vartype = "g_variance")
    area var_onephase var_multiphase      method      estimator gain  rel.eff
1     A     1987.117     1017.6327 twophase psynth extended 48.8 1.952686
2     B     3175.068     1019.1913 twophase psynth extended 67.9 3.115281
3     C     1180.853      763.0731 threephase psynth extended 35.4 1.547496
4     D     2290.652     1110.2454 twophase psynth extended 51.5 2.063194
```

The function call returns a data frame containing the one-phase variance (`var_onephase`) and the variance of the best performing multi-phase estimator (`var_multiphase`). The multi-phase estimation procedure is again specified in the `method` and `estimator` column. The last two columns quantify the potential benefit of the multi-phase estimation. The `gain` is the reduction (if its value is positive) in variance when applying the multi-phase as alternative to the one-phase estimation.

For example, it is indicated that the two-phase extended PSYNTH estimation procedure for small area "B" leads to a 67.9 % reduction in variance compared to the one-phase procedure. The column `rel.eff` specifies the relative efficiency which is defined as the ratio between the one-phase variance and the multi-phase variance:

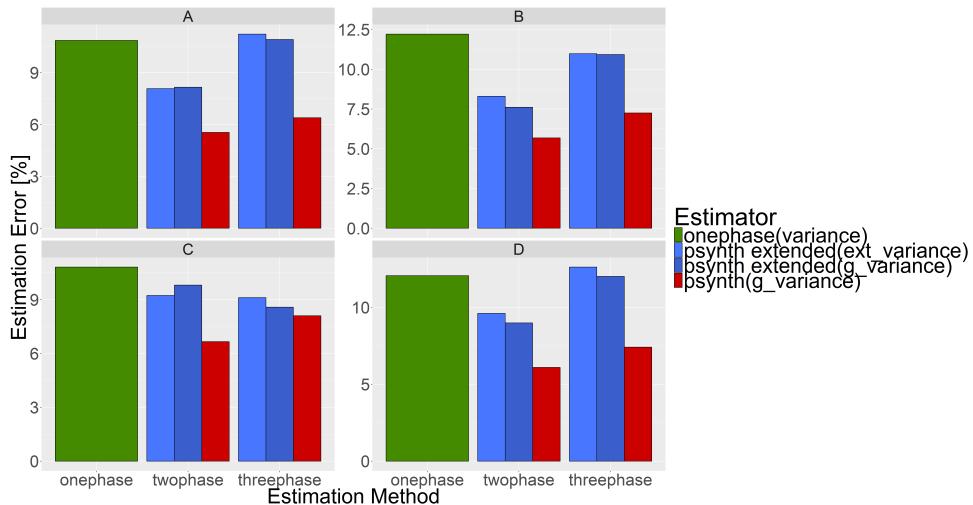
$$rel.eff[\%] = \frac{\hat{V}_{onephase}(\hat{Y})}{\hat{V}_{multiphase}(\hat{Y})} \cdot 100 \quad (1.31)$$

The relative efficiency can be interpreted as the relative sample size of the one-phase estimator needed to achieve the variance of the multi-phase estimator. For small area "B" we can thus see that we would have to increase the terrestrial sample size by factor 3 in the one-phase approach in order to get the same estimation precision as the two-phase EXTPSYNTH estimator. If the average costs for a terrestrial sample plot survey are known, the relative efficiency can thus be a simple means of quantifying the financial benefit of using multi-phase estimation for forest inventories.

### 1.7.2 Visualization

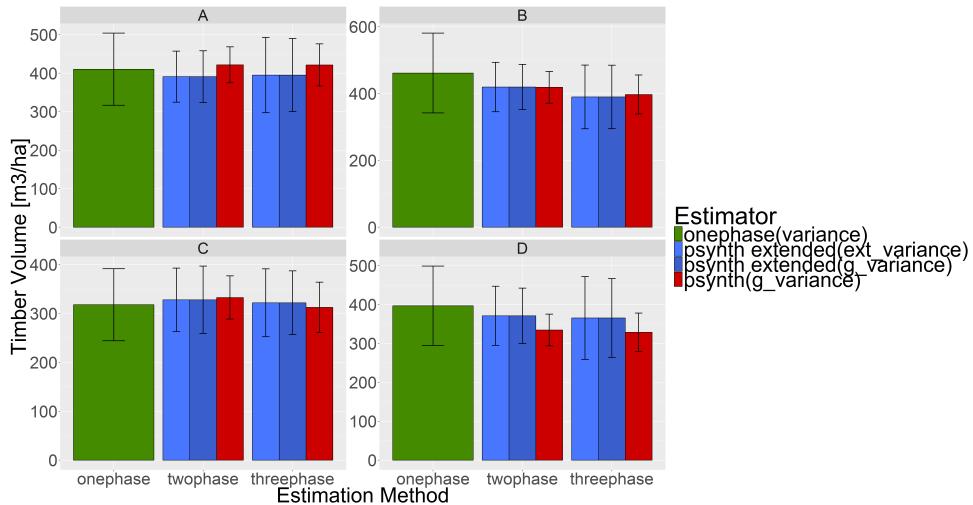
The **forestinventory** package also provides a S3 generic plot method based on the **ggplot2** package (Wickham, 2009) to visualize the estimation results in two ways: 1) the point estimates with overlayed confidence intervals, and 2) the estimation errors. Both plots can be obtained by passing the `esstable` object to the `plot()` function.

```
R> plot(grisons.sae.table, ncol = 2)
```



Whereas the estimation errors are plotted by default, the point estimates and confidence intervals are returned when setting the argument `yvar = "estimate"`. Note that the graphics can arbitrarily be extended by additional `ggplot2` parameterizations.

```
R> plot(grisons.sae.table, ncol = 2, yvar = "estimate") +
+   ylab("Timber Volume [m3/ha]")
```



## 1.8 Future plans

The **forestinventory** package currently provides a fairly well-rounded toolkit for forestry inventors to integrate auxiliary information into their estimates using the model-assisted methods under the design-based approach. Although 32 combinations of inventory scenarios, estimators and sample designs are covered, there are still potential improvements planned for the future. As this is an open-source project, everyone is encouraged to give feedback and/or make contributions on the package's development page on GitHub (Hill, 2017). Currently planned extensions include:

- Implement parallel procedures for efficiently calculating many small areas.
- Allow functions to accept objects of class `data.table` from the `data.table` package (Dowle & Srinivasan, 2017) to improve memory efficiency.
- Enable the user to choose other types of models than linear regressions fitted with OLS.

## Acknowledgements

We want to express our gratitude to Prof. H. Heinemann (Chair of Land Use Engineering, ETH Zurich) for supporting this study and providing the possibility of working on the package. We also want to thank Daniel Mandallaz for his support in completing the range of the already published estimators in the frame of the three-phase small area estimators, as well as many helpful discussions and advice throughout the implementation of our package. Our thanks also go to Meinrad Abegg for proofreading the manuscript, and to the Amt für Wald und Naturgefahren of the Swiss canton of grisons for providing the example data.



## Chapter 2

# Combining canopy height and tree species map information for large scale timber volume estimations under strong heterogeneity of auxiliary data and variable sample plot sizes

Andreas Hill<sup>1</sup>, Henning Buddenbaum<sup>2</sup>, Daniel Mandallaz<sup>1</sup>

<sup>1</sup>ETH Zürich

Department of Environmental Systems Science, Universitätstrasse 16, 8092 Zürich, Switzerland

<sup>2</sup>Trier University

Environmental Remote Sensing and Geoinformatics Department, Behringstrasse 21, 54286 Trier, Germany

Submitted to:

*European Journal of Forest Research* (in review).

- Henning Buddenbaum processed the airborne Laserscanning data and supported writing the manuscript.
- Daniel Mandallaz supported the statistical data analysis.

## Abstract

A timber-volume regression model applicable to the entire forest area of the federal German state of Rhineland-Palatinate is identified using a combination of airborne laser scanning (ALS)-derived metrics and information from a satellite-based tree species classification map available on the federal state level. As is common in many forest inventory datasets, strong heterogeneity in the LiDAR data due to different acquisition dates and misclassifications in the tree species classification map had noticeable effects on the regression model's performance. This article specifically addresses techniques that improve the performance of ordinary least square regression models under such restricting conditions. We introduce a calibration technique to neutralize the effect of misclassifications in the tree species variable that originally caused a residual inflation of 0.05 in adjusted  $R^2$ . Incorporating the calibrated tree species information improved the model accuracy by up to 0.07 in adjusted  $R^2$  and suggests the use of such information in forthcoming inventories. We also found that including ALS quality information as categorical variables within the regression model considerably mitigates issues with time lags between the ALS and terrestrial data acquisition and ALS quality variations (increase of 0.09 in adjusted  $R^2$ ). The model achieved an adjusted  $R^2$  of 0.48 and a cross-validated root mean square error ( $\text{RMSE}_{cv}$ ) of 46.7% under incorporation of the tree species and ALS quality information, and was thus improved by 0.12 in adjusted  $R^2$  (5% in  $\text{RMSE}_{cv}$ ) compared to the simple model only containing ALS height metrics (adjusted  $R^2=0.36$ ,  $\text{RMSE}_{cv}=51.7\%$ ).

## 2.1 Introduction

Forest inventory methods are the primary tools used to assess the current state and development of forests over time. They provide reliable evidence-based information that is used to define and identify management actions as well as to adapt forest management strategies to both national and international guidelines. Two methods that have become particularly attractive are so-called *double-sampling* (Mandallaz, 2008a, Ch. 5) and *mapping* (Brosofske et al., 2014) procedures. The core concept of these methods is to use predictions of the terrestrial target variable at additional sample locations where the terrestrial information has not been gathered. These predictions are produced by models that use explanatory variables derived from *auxiliary data*, commonly in the form of spatially exhaustive remote sensing data in the inventory area. Especially models to predict timber volume based on airborne laser scanning (ALS) have been extensively investigated for a long time (Næsset, 1997). The specific scope of double-sampling is to enlarge the terrestrial sample size by a much larger sample of predictions of the target variable in order to gain higher estimation precision without performing additional expensive terrestrial measurements. Model-dependent and design-based regression estimators are used in a broad range of double sampling concepts and methods (Gregoire & Valentine, 2007; Köhl et al., 2006; Schreuder et al., 1993; Saborowski et al., 2010; Mandallaz, 2013a,d) and have been applied to existing inventory systems (Breidenbach & Astrup, 2012; von Lüpke & Saborowski, 2014; Mandallaz et al., 2013; Magnussen et al., 2014; Massey et al., 2014). While double-sampling methods provide reliable estimates for a given spatial unit, e.g. a forest district, they do not provide information about the spatial distribution of the estimated quantity within this area. For this reason, the same modeling technique used in double-sampling procedures has also been intensively used to produce exhaustive prediction maps that provide pixelwise estimations of a target variable in high spatial resolution (Bohlin et al., 2017; Latifi et al., 2010; Tonolli et al., 2011; Hill et al., 2014; Nink et al., 2015).

To allow for an area-wide application of the prediction model, both double sampling and mapping methods require that the remote sensing data are available over the entire inventory area. This is usually not a limiting factor in *small-scale* applications. In the optimal case, the remote

sensing data are in principle collected in accordance to the specific study objective. Quality standards that have often been addressed are that *a*) the remote sensing data should be acquired close to or even at the time of the terrestrial inventory in order to ensure best possible comparability between the target variable on the ground and the remote sensing derived variables (McRoberts et al., 2015); *b*) the remote sensing technology and its spectral and spatial resolution should be chosen according to the modelling purpose (Köhl et al., 2006); and *c*) the variation in quality of the remote sensing data over the inventory area should be minimized in order to avoid artificial noise in the data (Næsset, 2014). Despite the increasing availability and decreasing costs of remote sensing data (White et al., 2016), these quality standards of the remote sensing data can often not be guaranteed for *large-scale* applications (Maack et al., 2016), and trade-offs must be accepted (Jakubowski et al., 2013). The prime objective is then to produce the best possible prediction model given the restrictions imposed by the available remote sensing information. The exploration of scarcely used remote sensing products and the optimization of prediction models under severe quality restrictions in the remote sensing data are thus one of the challenges in large-scale model-supported inventory applications.

Among the still rarely used remote sensing data in large scale applications, the integration of tree species information in prediction models - especially for timber volume estimation - has been stated as some of the most promising but often missing information (Koch, 2010; White et al., 2016). As timber volume estimations on the single tree level in forest inventories are often based on species-specific biomass and volume equations (Husmann et al., 2017; Zianis et al., 2005), the application of species-specific models is expected to be a key factor for improving estimation precision (White et al., 2016). This has been supported by studies from Breidenbach et al. (2008) who achieved a substantial improvement in accuracy of their timber volume prediction model when including a variable estimating the deciduous proportion derived from leaf-off ALS data. Similar gains in model performance were also reported by Straub et al. (2009) and Latifi et al. (2012) who used broadleaf and coniferous information based on color infrared orthophotos as a categorical explanatory variable. However, studies that explore the use of more species-specific information (i.e. a further discrimination of tree species) as explanatory variables have been rare. Further investigations are thus necessary especially in countries whose forests are characterized by a larger variety of tree species that may also occur in mixed and uneven-aged stands (McRoberts et al., 2010). The area-wide tree species information in most studies was obtained from satellite and airborne remote sensing sensors based on automatic classification methods. Whereas the presence of misclassifications has already been addressed (Latifi et al., 2012), an issue that has so far been neglected is how misclassifications actually affect the prediction model (Gustafson, 2003).

A frequently encountered problem in large scale forest inventories is the lack of temporal synchronicity between the remote sensing acquisition and the terrestrial survey. As a result, the available remote sensing data often exhibit notable time-lags with respect to the date of the terrestrial inventory. This has often been addressed as a major drawback, especially for the application of design-based change estimation (Massey & Mandallaz, 2015b).

Our study is embedded in the current implementation of design-based regression estimators (Mandallaz, 2013a; Mandallaz et al., 2013; Mandallaz, 2013d) for estimating the standing timber volume within the state and communal forest management units over the entire state of Rhineland-Palatinate (RLP, Germany). With respect to this overall objective, the aim of this study was to derive an ordinary least square (OLS) regression model to generate predictions of the standing timber volume associated with a sample location of the Third German National Forest Inventory (BWI3) over the entire state and communal forest area ( $6155 \text{ km}^2$ ). A merged ALS dataset from different acquisition years and a satellite-based tree species classification map for the five main tree species in RLP was available for the entire inventory area and consequently used to derive predictor variables. The major limiting factors for using these data in a regression analysis are (**i**) variation in the ALS data quality as well as time-lags of up to 10 years between the ALS acquisitions and

the terrestrial survey, (ii) misclassifications in the tree species classification map and (iii) the ambiguous choice of a suitable extraction area (*support*) for all remote sensing information under angle count sampling in the terrestrial survey (variable sample plot sizes). For this reason, we address the following specific research questions:

1. How can tree species map information be optimally used within a regression model that predicts timber volume? What effects do misclassifications have on the predictions and how can these effects be minimized?
2. What are the effects of quality restrictions and substantial time lags between the ALS- and terrestrial data acquisition on the regression model and how can these effects be mitigated?
3. Does support size influence model accuracy? What is the optimal support size and what are the determining factors?

## 2.2 Materials and Methods

### 2.2.1 Study Area

The German federal state Rhineland-Palatinate (RLP) is located in the western part of Germany and borders Luxembourg, France and Belgium (figure 2.1). With 42.3% (appr. 8400 km<sup>2</sup>) of the entire state area (19850 km<sup>2</sup>) covered by forest, RLP is one of the two states with the highest forest coverage among all federal states of Germany (von Thünen-Institut, 2014). The forest area of RLP is divided into three ownership classes, i.e. state forest (27%), communal forest (46%) and privately owned forest (27%). The most frequent tree species in RLP are European beech (*Fagus sylvatica*, 21.8%), oak (*Quercus petrea* and *Quercus robur*, 20.2%), Norway spruce (*Picea abies*, 19.5%), Scots pine (*Pinus sylvestris*, 9.9%), Douglas fir (*Pseudotsuga menziesii*, 6.4%), European larch (*Larix decidua*, 2.4%) and Silver fir (*Abies alba*, 0.7%). The share of broadleaf tree species is 58.7%. The forests of RLP further exhibit heterogeneous structures (von Thünen-Institut, 2014): around 82% of the forest area in RLP are mixed forest stands (i.e. at least two different tree species occur in the same stand) and 69% of the forest area exhibit a multi-layered vertical structure. While the average tree age is around 80 years, most of the forest area (20%) is occupied by trees between 40 and 60 years of age, whereas 27% of the trees are older than 100 years. Spatially variable climate conditions have a strong influence on the local growth dynamics as well as tree species composition and create a large variety of forest structures, ranging from characteristic oak coppices (Moselle valley), pure spruce, beech and Scots pine forests (e.g. Hunsrück and Palatinate forest) to mixed forests comprising variable proportions of oak, larch, spruce, Scots pine and beech. Accordingly, RLP has been divided into 16 bioclimatic growing regions that form homogeneous areas with respect to the afore mentioned characteristics (Gauer & Aldinger, 2005).

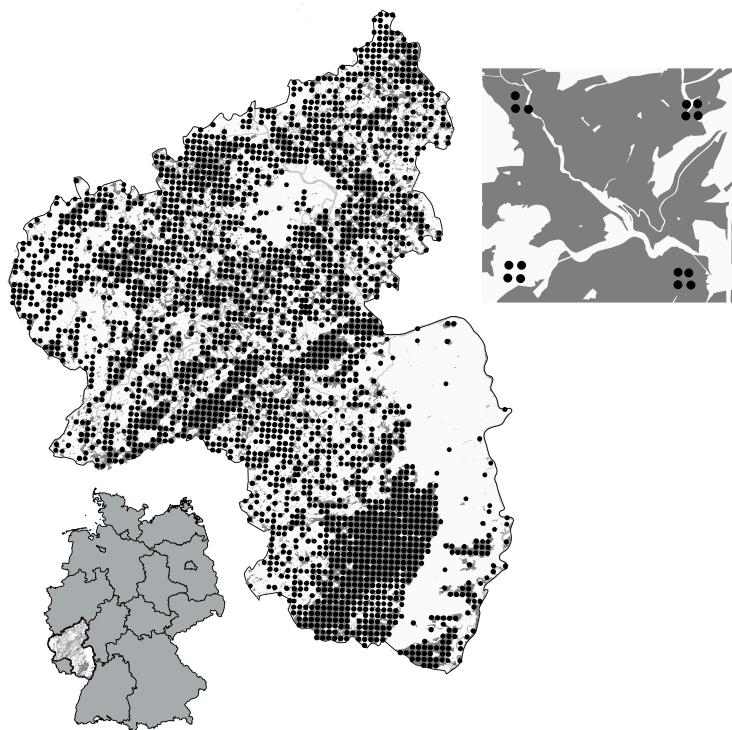


Figure 2.1: Spatial distribution of the BWI3 cluster samples over Rhineland-Palatinate

### 2.2.2 Terrestrial Inventory Data

The German National Forest Inventory (NFI) is carried out over the entire forest area of Germany in reoccurring time periods of 10 years. The most recent inventory (BWI3) has been conducted in the years 2011 and 2012. In this framework, Rhineland-Palatinate is covered by a 2x2 km grid that defines the sample locations for the terrestrial survey. A sample unit consists of four sample locations (also referred to as *sample plots*) that are arranged in squares (so called *clusters*) with a side length of 150 metres (figure 2.1). The number of plots per cluster can however vary between 1 and 4 depending on forest/non-forest decisions on the plot level (Bundesministerium fÃ¼r Ernährung, 2011). In the field survey of the BWI3, sample trees for timber volume estimations are selected according to the angle count sampling technique (Bitterlich, 1984), using a basal area factor (*BAF*) of 4 that is respectively adjusted for boundary effects at the forest border (Bundesministerium fÃ¼r Ernährung, 2011). A further selection criterion for a tree to be recorded is a diameter at breast height (*dbh*) of at least 7 cm. This sampling technique was applied to 8092 sample plots (2810 clusters) in RLP, resulting in the collection of 56561 sample trees for which the *dbh*, the tree diameter at 7 m (*D7*) and the tree species were recorded for all trees. Tree height measurements were conducted only for a subset of all sample trees and used to predict the height for the remaining sample. During the last inventory, all plot center positions were remeasured with differential global positioning system (DGPS) technique. Knowledge about the exact plot positions were considered crucial to provide optimal comparability between the terrestrial observations and the information derived from the auxiliary data. A detailed analysis by Lamprecht et al. (2017) indicated that horizontal DGPS errors do not exceed 8 meters for 80% of all plots in RLP. For 162 plots, the DGPS coordinates were replaced by their former target coordinates due to missing or implausible values. In order to derive a volume estimation for each sample tree, the BWI3 estimates a taper curve for each sample tree by calibrating the random effects term of linear mixed-effects taper models with

the set of diameters and corresponding height measurements taken from the respective sample tree (Kublin et al., 2013). The integration of the derived taper curves consequently lead to a volume prediction for each sample tree. Since the overall objective of the study was to subsequently use the identified regression model for design-based timber volume estimations of state and communal forest management units, we already restricted the sample plots used for modeling to the state and communal forest area (73% of the entire forest area of RLP). This provides the advantage that when the regression model is used as an *internal model* in design-based estimators, the model predictions hold the assumption on the residuals to be zero on average over the state and communal forest area by construction of OLS technique (Mandallaz, 2013a; Mandallaz et al., 2013; Mandallaz, 2013d). The dataset of this study hence comprised 5791 plots (2055 clusters). For this sample, the timber volume density per hectare on plot level,  $Y(x)$ , was calculated according to the formula of one-phase one-stage sampling (Mandallaz, 2008a, Ch. 4.2 ). The timber volume density per hectare on plot level was used as the response variable in the regression analysis.

Table 2.1: Descriptive statistics of the forest observed on NFI sample plots located iwithin communal and state forest area (n=5791).

Variable	Mean	SD	Maximum
Timber Volume (m <sup>3</sup> /ha)	300.86	195.55	1375.31
Mean DBH (mm)	354.90	137.22	1123.20
Mean height (dm)	239.60	72.43	497.43
Mean stem density per hectare	101.00	114.01	1010.31

### 2.2.3 Auxiliary Data

#### ALS Canopy Height Model

Between 2003 and 2013, the topographic survey institution of RLP acquired airborne laser scanning (ALS) data over the entire state of RLP at leaf-off condition (Figure 2.2). The objective of this campaign was to derive a countrywide digital terrain and surface model based on the acquired ALS point clouds. During the extended acquisition period, airborne laser scanning technology and data quality evolved significantly. The tiles recorded in 2002 and 2003 have a rather poor quality with about only 0.04 points per m<sup>2</sup> , while more recently acquired datasets contained about 5 points per m<sup>2</sup> . The data was delivered as two separate datasets comprising the Vegetation First Pulse (VEF) and Ground (GRD) points. All point clouds were stored as three-column (easting, northing, and height above sea level) ASCII files in tiles of 1 km<sup>2</sup>. In order to create a surface model (DSM) in a given raster resolution, the highest point of the combined VEF and GRD dataset was identified in each raster cell and saved as a thinned surface point cloud. For the elevation model (DEM), the mean of all GRD points in the cell was calculated, and the result was saved as a thinned ground point cloud. The thinned point clouds were then aggregated to larger tiles and interpolated to raster images using a Delauney interpolation in the Matlab software (Mathworks, 2017). The resulting DSM and DEM raster sets were then subtracted from each other to calculate a canopy height model (CHM) in raster format, providing discrete information about the canopy surface height of the entire forest area of RLP in a spatial resolution of 5 meters. The thinning process led to much smaller datasets that could be processed in larger tiles and considerably lowered processing times compared to the original dense point clouds. Since the data was recorded in leaf-off condition, the original point clouds contained many returns from within the crowns of deciduous trees. The thinned dataset provided the advantage that those measurements did not skew the vegetation height estimate in the final CHM.

As explanatory variables, the mean canopy height (*meanheight*) and the standard deviation

(*stddev*) were calculated as the mean and standard deviation of all raster values within a pre-defined circle (i.e. *support* of the explanatory variable, see section 2.2.4) around each sample plot center. In order to correct for edge effects at the forest border, each support area was previously intersected with the state and communal forest area, which was defined by a polygon mask provided by the forest service (figure 2.3b). Restricting the support area and thus the evaluation of the auxiliary data to the forest area is a means to optimize the coherence between explanatory variables computed at the forest boundary and the corresponding terrestrial response variable (Mandallaz et al., 2013). The tree height is one prominent predictor variable in the taper functions of the BWI3 that are used to calculate a timber volume value for each sample tree (Kublin, 2003; Kublin et al., 2013). A visual inspection of the tree volumes of all sample trees collected in the BWI3 within RLP against their tree heights also revealed the characteristic shape of an allometric relationship between these variables (Online Resource 1). It was hypothesized that this relationship on single-tree level is also apparent on the aggregated level of a sample plot and cluster, and can be used within the frame of regression modeling.

The strength of correlation between *meanheight* and timber volume on plot level was expected to show high variation according to the mentioned time-lag up to 10 years between ALS acquisition and terrestrial survey. The quality of the height information was also expected to vary according to changing sensor technologies and different point densities used over the years. For these reasons, the ALS acquisition year (*ALSyear*) for each sample plot was considered as a potential categorical explanatory variable to explain the variation in the data introduced by these factors. For this purpose, the acquisition year *2008* was further divided into *2008* and *2008\_1*. In the latter, the data quality turned out to be very poor due to sensor failures during the acquisition. Additionally, the years *2006* and *2007* as well as *2012* and *2013* were pooled in order to increase the number of observations per factor level for modelling reasons. As a result, the *ALSyear* variable comprised nine categories (*2002*, *2003*, *2007*, *2008*, *2008\_1*, *2009*, *2010*, *2011* and *2012*).

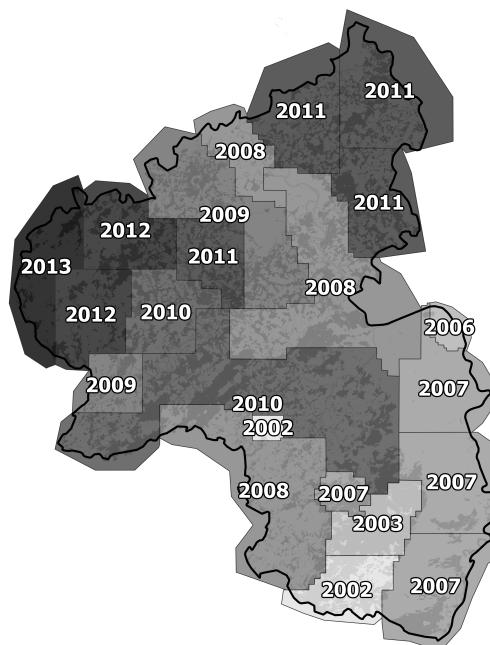


Figure 2.2: Separate ALS acquisitions in Rhineland-Palatinate over the years. The colors also indicate the quality of the data: *light*: low point densities ( $0.04/m^2$ ), *dark*: high point densities ( $>4/m^2$ ). Blue semitransparent layer: state and communal forest area.

## Tree Species Classification Map

A countrywide satellite-based classification map of the five main tree species (European beech, Sessile and Pedunculate oak, Norway spruce, Douglas fir, Scots pine) described in Stoffels et al. (2015) was used to derive tree species information on sample plot level. The classified tree species map has a grid size of 5 meters and predicts five of the seven tree species that are used in the BWI3 taper functions (Kublin et al., 2013) to calculate the timber volume of a sample tree. Due to unavailable satellite data for the classification, the tree species map excluded one patch with an area of 415 km<sup>2</sup> in the south-west part of RLP, and two further patches with an area of 76 km<sup>2</sup> and 100 km<sup>2</sup> in the northern part (Stoffels et al., 2015). The tree species information was consequently missing for 411 (7%) of the 5791 sample locations.

## Prediction of main plot tree species

A visual inspection of all BWI3 sample trees of RLP suggested that a stratification of the relation between tree height and timber volume according to these seven tree species may provide a considerable reduction in variation within the tree species groups (Online Resource 1). This led to the hypothesis that this tree species specific signal might also be apparent on sample plot and cluster level and can consequently be used to increase the accuracy of the prediction model. Based on the tree species classification map, the main tree species of each sample plot was calculated as an additional categorical explanatory variable (*treespecies*) with six categories following a similar approach as Latifi et al. (2012): one of the five tree species was assigned as the main plot tree species if its proportion within the edge-corrected support around the sample location exceeded a predefined threshold. If this threshold was not reached by any of the five tree species, the respective sample plot was assigned the category 'Mixed'. We hypothesized that the choice of the threshold-value might have an influence on the resulting classification accuracy and the regression model accuracy (section 2.2.5). We thus investigated the application of 5 threshold settings, i.e. 0%, 50%, 60%, 80% and 100%.

## Calibration

Our analyses revealed that the prediction of the main tree species for a sample plot can be subject to misclassifications (section 2.3.1). Errors in the explanatory variables of linear regression models can however lead to a bias of the regression coefficients in the direction of zero due to an artificial introduction of noise (Carroll et al., 2006, Ch. 3). This can cause an inflation of the residual variance and a consequent decrease of the model accuracy (Magnussen et al., 2010). In case of classification, the impacts of misclassifications on the model properties are even harder to predict (Gustafson, 2003, Ch. 3). While errors in the explanatory variables do not affect the unbiasedness of the estimators in the design-based framework, a reduction or elimination of the classification errors could provide an improvement of the regression model accuracy and thereby potentially lead to smaller prediction and estimation errors. We therefore addressed the effect of misclassifications in the *treespecies* variable categories as well as means to correct these errors.

We transferred the concept of *regression calibration* as known from classical measurement error statistics (Carroll et al., 2006) to the problem of misclassifications in the *treespecies* variable. In regression calibration, one considers an error-prone explanatory variable  $W$  that can be measured in high quantity, whereas  $X$  constitutes the same but error-free variable whose determination is however very expensive. In order to yield a corrected or less error-prone version of  $W$ , one can define a calibration model  $f_{calmod}(X, W)$  that predicts  $X$  as a function of  $W$ . After calibration on a training set,  $f_{calmod}()$  can then be applied to any observed  $W$  and yields the corrected, less error-prone variable  $W_{calib}$ . Using  $W_{calib}$  instead of  $W$  in the regression model then asymptotically

provides an unbiased estimate of the regression coefficients and thus corrects for the attenuation to zero.

We transferred this concept by using a random forest algorithm (Breiman, 2001) as calibration model. We considered the main tree species of the sample trees at each plot location  $x$  as the error-free variable  $treespecies_{terr}$ , that would also yield the highest model accuracies when used as predictor variable. The objective of the calibration model was thus to provide an improved classification accuracy of each predicted main plot tree species category with respect to  $treespecies_{terr}$ . The calibration model was considered to correct for potential systematic misclassifications and thus minimize the effect of misclassifications on the regression model when substituting the uncalibrated with the calibrated  $treespecies$  variable. The random forest algorithm is a machine learning algorithm that grows a large number of decorrelated classification trees by considering only a subset of all provided predictor variables for each split. In the case of classification, new data are thus predicted by aggregating the predictions of all trees using a majority vote. We calibrated the random forest algorithm ( $f_{RF}$ ) with a set of  $p$  predictor variables that comprised the initial prediction of the main plot tree species ( $treespecies$ ), the mean canopy height ( $meanheight$ ) and standard deviation ( $stddev$ ) derived from the CHM, the proportion of coniferous trees estimated from the tree species classification map ( $prop.conif$ ) and the bioclimatic growing region ( $wgb$ ) at the sample location (equation 2.1). An advantage for using those explanatory variables in the calibration model was that they also provided explanatory power in the regression model. This approach thus saved computation time and minimized data storage. The calibration model was implemented using the random forest algorithm (Liaw & Wiener, 2002) in the statistical software  $R$  (R Core Team, 2017). The algorithm was grown with 2000 trees, considering  $\sqrt{p} \approx 3$  of the predictors for each split.

$$treespecies_{terr}(x) = f_{RF}(treespecies, meanheight, stddev, prop.conif, wgb) \quad (2.1)$$

The calibration model was subsequently applied to the entire dataset. We then investigated the effect on the regression model performance (regression coefficients, model accuracy) when substituting the calibrated (less error-prone) for the uncalibrated (most error-prone) variable, and likewise for the actual (error-free) main plot tree species derived from the sampled trees of the respective sample plot under identical threshold settings.

#### 2.2.4 Choice of Support under Angle Count Sampling

One characteristic of angle count sampling applied in the BWI3 is that a sample plot does not have a fixed radius in which trees are selected (*fixed-radius plot*), but each tree generates an individual radius from the plot center depending on its diameter at breast height (*variable-radius plot*). This tree-individual radius is known as the *limiting distance* from the plot center where the tree would still be included in the sample. A consequence of the absence of a fixed plot radius is the question about the optimal support (Hollaus et al., 2007), i.e. the spatial extent around the plot center in which the auxiliary information is evaluated and transformed into an explanatory variable. It has widely been hypothesized that the best relationship between the target variable on the ground and any explanatory variable derived from the auxiliary information is obtained if the support is spatially identical to the sample plot extent. In case of angle count sampling, an individual extent for each sample plot can be approximated by regarding the maximum limiting distances of its sample trees as the outer plot radius. However, many design-based applications under double-sampling do not allow for a between-plot change of the support for a specific explanatory variable (Mandallaz, 2013d,a).

For this reason, the task is to find a unique support for each auxiliary information that leads to the best overall model accuracy. Deo et al. (2016) conducted extensive analysis to identify optimal supports for modelling standing timber volume for *variable-radius plot* designs in conifer forests. They analysed 24 different radii (i.e. circular supports) in which they extracted 57 metrics from a ALS derived point cloud with an average point density of 18 pulses per square meter. They successively evaluated the prediction performance of each support size by using the ALS metrics in a random forest algorithm and comparing the resulting model accuracies. In order to identify the best-performing supports for our explanatory variables, we followed a similar approach. The explanatory variables were calculated using *individual* (i.e. plot-varying) supports (*ind*), i.e. an individual support radius was used for each plot according to the maximum limiting distance of all sample trees associated to the respective sample plot. We then compared the model accuracies achieved by the individual supports against the model accuracies from a set of *fixed* (i.e. non plot-varying) supports. The extents of the fixed supports were chosen from the cumulative distribution function (ECDF) of the maximum limiting distances of all 5791 sample plots of the analysed forest area (Fig. 2.3a). We considered the 25<sup>th</sup> ( $q_{25}$ , 9 meters), 50<sup>th</sup> ( $q_{50}$ , 12 meters), 80<sup>th</sup> ( $q_{80}$ , 15 meters) and the 100<sup>th</sup> ( $q_{100}$ , 38 meters) percentiles, resulting in support radii of 18, 24, 30 and 76 meters (Fig. 2.3). While in this study we also used circular supports to extract the auxiliary information, also other support-shapes are possible (e.g. rectangles, hexagons). We also want to emphasize that the use of different support sizes for each explanatory variable is perfectly valid in the infinite population framework of design-based estimators (Mandallaz, 2013d,a).

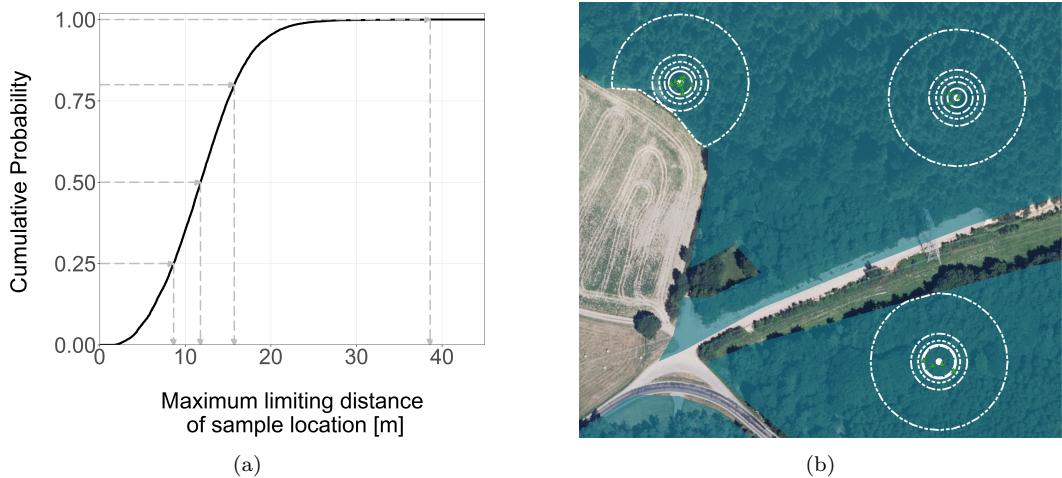


Figure 2.3: Identification (a) and visualization (b) of potential support radii used for calculating the predictor variables on plot level based on ECDF of maximum limiting distances of all BWI3 sample locations in RLP.

## 2.2.5 Model Building and Evaluation

In order to judge the quality of the *treespecies* variable, the user's accuracy for each classified species category and the overall accuracy of the classification scheme was calculated based on the confusion matrix (Congalton & Green, 2008). As reference data, we calculated the actual main plot tree species by applying the respective threshold to the sample trees of each sample plot. The classification accuracy was evaluated for all support sizes for both the calibrated and the uncalibrated *treespecies* variables. The measures of the regression model accuracy using both CHM- and *treespecies* variables were defined as the 10-fold cross-validated root mean square error (RMSE<sub>cv</sub>, equation 2.2) and the adjusted coefficient of determination (adjusted  $R^2$ ) of the multiple

linear regression model defined in equation 2.3. Additionally, we considered the interaction terms *meanheight:treespecies*, *meanheight<sup>2</sup>:treespecies*, *meanheight:ALSyear*, *stddev:ALSyear* and *meanheight:stddev* and performed a variable selection based on the Akaike Information Criterion (AIC) (Akaike, 2011) in order to minimize the number of variables in the model. Due to a pronounced unbalanced design in the *treespecies-ALSyear* strata (Online Resource 2), no interaction between *treespecies* and *ALSyear* was possible. We evaluated the model for all support combinations, considering the use of individual support sizes for each auxiliary information, using both the calibrated and the uncalibrated *treespecies* variable. The calibration model (section 2.2.3) for the *treespecies* variable was recalculated for each respective support-threshold setting.

206 sample plots included no sample trees and the timber volume density  $Y(x)$  was thus set to zero. These *zero-plots* were removed from the modeling dataset since they acted as leverage points in cases where the ALS height metrics were recorded long before the terrestrial survey. Together with the missing tree species information (section 2.2.3), the modeling dataset  $s$  was limited to  $n=5171$  observations.

$$RMSE = \sqrt{\frac{\sum_{x \in s} (\hat{Y}(x) - Y(x))^2}{n}} \quad (2.2a)$$

$$RMSE\% = \frac{RMSE}{\frac{1}{n} \sum_{x \in s} Y(x)} \quad (2.2b)$$

$$\begin{aligned} Y(x) = & \beta_0 + \beta_1 * meanheight + \beta_2 * meanheight^2 + \beta_3 * stddev + \\ & \beta_4 * ALSyear_1 + \dots + \beta_{12} * ALSyear_9 + \\ & \beta_{13} * treespecies_1 + \dots + \beta_{18} * treespecies_6 + \varepsilon(x) \end{aligned} \quad (2.3)$$

## 2.3 Results

### 2.3.1 Classification Accuracies

#### Effect of Support Size and Threshold

The lowest user's accuracies (*UA*) for the uncalibrated tree species variable were mostly realized using high thresholds of 80% and 100% (figure 2.4). A plausible reason for this is that raising the threshold to higher values (e.g. 80%, 100%) distinctively increases the probability of the reference class (based on the sample trees of the sample location) to be assigned as class 'Mixed', while the much coarser spatial resolution of the tree species map causes the *predicted* class to remain classified as one of the five tree species. However, as the support size is increased, so does the number of tree species raster cells to be evaluated at the sample location, thereby increasing the probability that the predicted class will be 'Mixed'. For this reason, most tree species exhibit an increase in user's accuracy under higher thresholds with higher support sizes. This scale-threshold dependency of the user's accuracy particularly affects tree species that most commonly occur in mixed forest stands in Rhineland-Palatinate (*Scots pine*, *oak* and *beech*), whereas the user's accuracies for tree species that are mostly prominent in pure forest stands (*spruce*, *Douglas fir*) logically turned out to be much more robust to changes in the thresholds and support sizes.

Among the uncalibrated tree species predictions, *beech* and *spruce* produced the best predictions achieving UAs of up to 70% and 80%. Although the predictions for *Douglas fir* and *Scots pine* generally performed less well than *beech* and *spruce*, similar UAs can be produced by adjusting the threshold and support choices. UAs for *oak* never performed better than 50%. A detailed table of the user's and overall accuracies is provided in Online Resource 3.

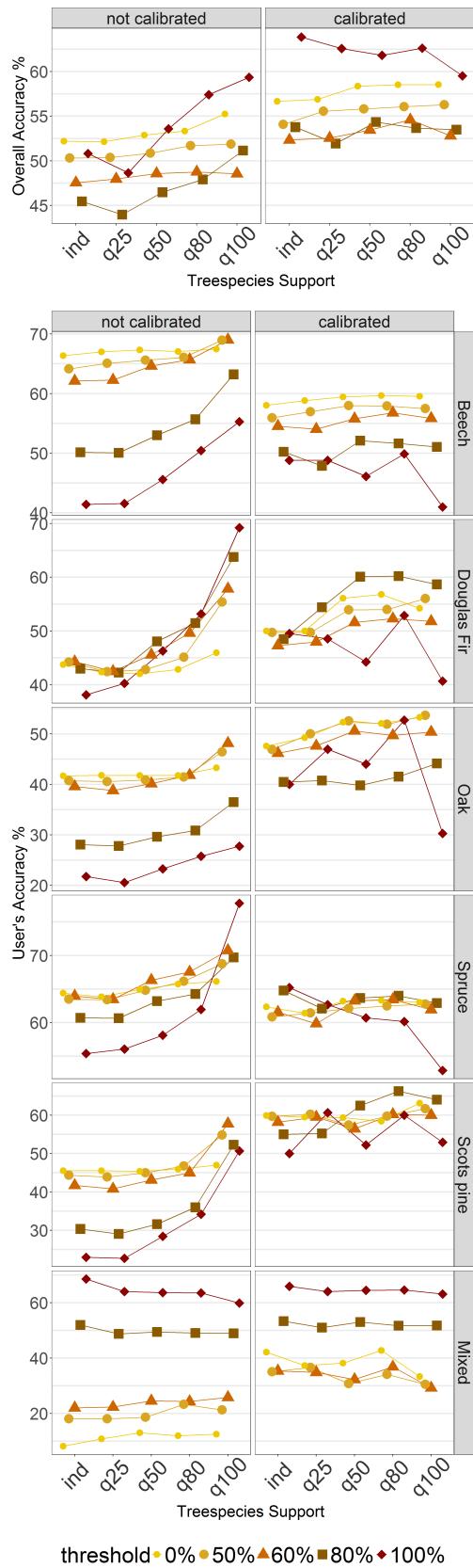


Figure 2.4: Classification accuracy for the main tree species of a sample location *before* and *after* calibration: *top*) overall accuracies. *bottom*) user's accuracies. *ind*: plot individual support sizes.

## Calibration

Calibration substantially diminished the effect of the scale-threshold dependency for the five tree species and also increased the UAs for *Scots pine* and *oak*. Whereas the UAs for *beech* and *spruce* were found to be slightly lower after calibration, the overall accuracy under each support choice was always considerably increased by calibrating the tree species prediction (figure 2.4). With respect to the calculated random forest models, the initial tree species prediction (*treespecies*) and the information about the growing region (*wgb*) turned out to be the most valuable information, followed by the estimated proportion of coniferous trees (*prop.conif*) and the mean canopy height (*meanheight*).

### 2.3.2 Regression Model Accuracies

#### Effect of Support Size and Threshold

Figure 2.5 shows the accuracies of the regression model (equation 2.3) achieved under all possible combinations of support sizes for the auxiliary data. The stepwise selection procedure always included all considered single and interaction terms. In terms of adjusted  $R^2$  and  $\text{RMSE}_{cv}$ , the analysis revealed that the choice of the CHM support size controls the overall level of the model's accuracy. The information about the main plot tree species can then be used to further improve the model fit under suitable *treespecies* support and threshold settings. When using the uncalibrated *treespecies* variable, an increase of the *treespecies* support size causes an increase in the model performance if low thresholds are used, whereas high thresholds (80%, 100%) cause a decrease in the model performance. This threshold-dependency could be removed by calibrating the *treespecies* variable. The highest adjusted  $R^2$  and the lowest  $\text{RMSE}_{cv}$  were realized using the *q50* support for both the CHM and calibrated *treespecies* variables in combination with a *treespecies* threshold of 100%, resulting in (adjusted  $R^2$  of 0.48 and  $\text{RMSE}_{cv}$  of 136.62 m<sup>2</sup>/ha (43.8%). However, various support and threshold combinations for the CHM and *treespecies* variables can be used to yield almost identical  $\text{RMSE}_{cv}$  and adjusted  $R^2$  values. A detailed table of the model accuracies is given in Online Resource 4.

#### Effect of Misclassifications

We accessed the magnitude of the misclassification effect for all models that were analysed in section 2.3.2, i.e. for all possible support and threshold combinations for the CHM and *treespecies* predictor variables. We first compared the adjusted  $R^2$  of each model when using the uncalibrated *treespecies* variable against the adjusted  $R^2$  using the actual, i.e. error-free variable. We then did the same comparison for the model using the calibrated *treespecies* predictor variable. Figure 2.6 provides a visualization of this comparison. Note that only the model with the predicted tree species variables can be applied to additional sample locations where no terrestrial survey has been carried out.

As expected, the highest adjusted  $R^2$  for every evaluated model was always achieved using the error-free tree species variable, whereas the missclassifications in the tree species variable led to a systematic decrease of the model accuracy. The calibration of the initially predicted main plot tree species using the random forest classification algorithm (section 2.2.3) turned out to not only improve the classification accuracies (section 2.3.1), but also to considerably decrease the effect of the missclassifications on the regression model predictions and accuracy. Figure 2.6 (right) shows that the adjusted  $R^2$  under the actual and the calibrated predicted tree species variable are in general much closer to, and in many cases even on the identity line. The differentiation into two distinct point clouds results from the poor model performance under support size *q100* for the CHM variables (i.e. the lower point cloud). Whereas the missclassifications in the uncalibrated

*treespecies* variable led to a residual inflation of 0.01 - 0.05 in adjusted  $R^2$ , it was only between 0 and 0.01 after calibration. Further analysis revealed that when using the calibrated *treespecies* variable, the regression coefficients were almost identical to the ones received using the actual main plot tree species.

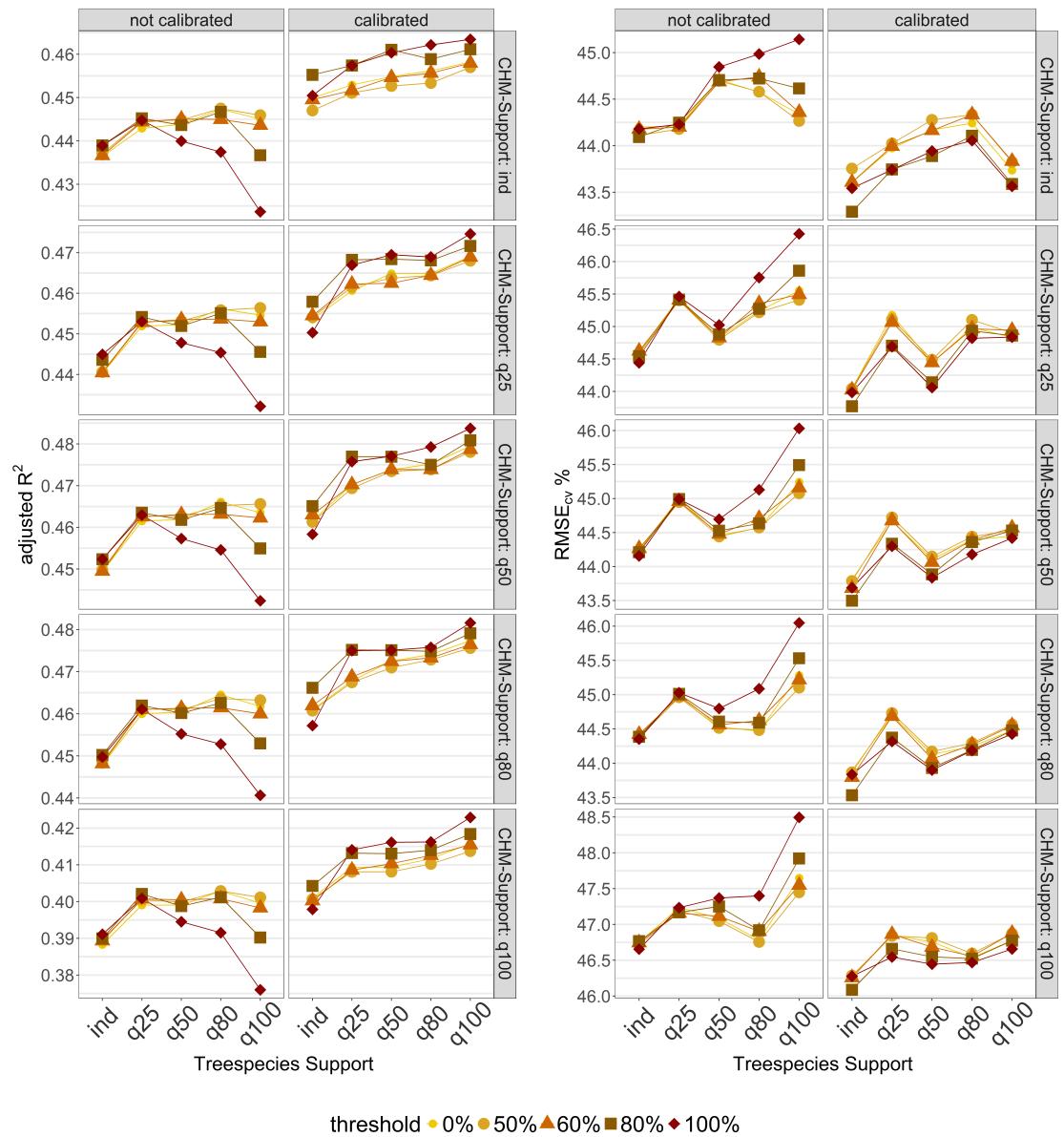


Figure 2.5: 10-fold RMSE<sub>cv</sub>[%] and adjusted  $R^2$  realized under various support choices for the CHM and *treespecies* explanatory variables

### 2.3.3 Final Regression Model

In order to address research questions 1 and 2 (i.e. the gain in model accuracy by tree species information and effect of heterogeneity in the ALS data), we investigated the model properties in more detail. For this purpose, we decided to use the best found model that was achieved under

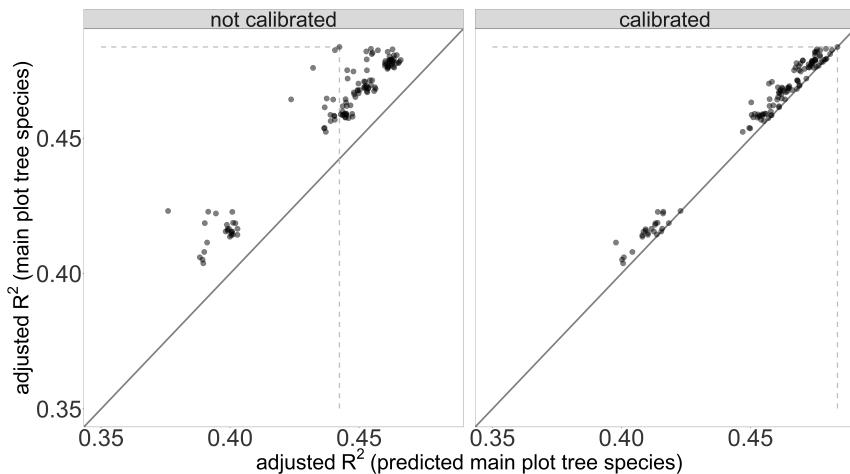


Figure 2.6: Effect on the adjusted  $R^2$  when substituting the actual main tree species with the predicted main tree species of a sample plot. The dotted line tracks the the model with the highest adjusted  $R^2$  under the use of the error-free *treespecies* variable. Semitransparent colours for the data points are used to visualize overlap.

the support settings of  $q50$  for both auxiliary data with a threshold of 100% for the tree species variable as the regression model of choice. The reason for inspecting this model was that *a)* the model provided the highest adjusted  $R^2$  among all validated models while reducing the data handling complexity for upcoming applications (i.e. identical support sizes for all remote sensing data) and *b)* the calibration neutralized the effects of misclassifications on the model predictions. The interaction term between *meanheight*<sup>2</sup> and *treespecies* (i.e. considering separate curvatures for each tree species) turned out not to have a significant influence on the model accuracy and was dropped, resulting in an adjusted  $R^2$  of 0.48 and a slightly increased  $RMSE_{cv}$  of 140.62 m<sup>2</sup>/ha (46.7%). The final model thus comprised 39 parameters (regression coefficients), i.e. the intercept, 3 main effects for continuous variables, 13 main effects for categorical variables and 22 interaction parameters (table 2.2).

We also conducted an analysis for detecting influential data points or outliers for the final regression model. We here considered the commonly applied criteria of leverages and Cook's Distance as amongst others described in Fahrmeir et al. (2013, p. 160-167). The critical threshold of  $2p/n$  (i.e. twice the average of the hat matrix' diagonal entries) was exceeded by 10% of the observations. However, only 3% of these leverage points were assigned to studentized residuals with values  $> 1$  or  $< -1$ . Removing these observations from the dataset and refitting the model led to an adjusted  $R^2$  of 0.49 compared to 0.48 when including them. Additionally, Cook's Distance values  $D_i$  did not exceed a value of 0.019, and were thus far apart from the commonly used critical threshold of  $D_i > 0.5$  that indicate a considerably change of the regression model results when omitting them. We thus decided not to remove any observations from the modelling dataset. We thus decided not to remove any observations from the modelling dataset.

### Interpretation of Final Regression Model

Figure 2.7 provides a visualisation of the tree species prediction functions separated by the ALS acquisition years. Sample plots classified as *oak* and *Scots pine* revealed to have an almost identical relationship (nearly identical slopes) for the mean canopy height - timber volume relationship. They only differ by a marginally higher intercept for *Scots pine* plots, meaning that given the same mean canopy height a sample plot dominated by *Scots pine* yields a marginally higher timber volume on the plot level than a plot dominated by *oak*. *Beech*-dominated sample plots tend to achieve

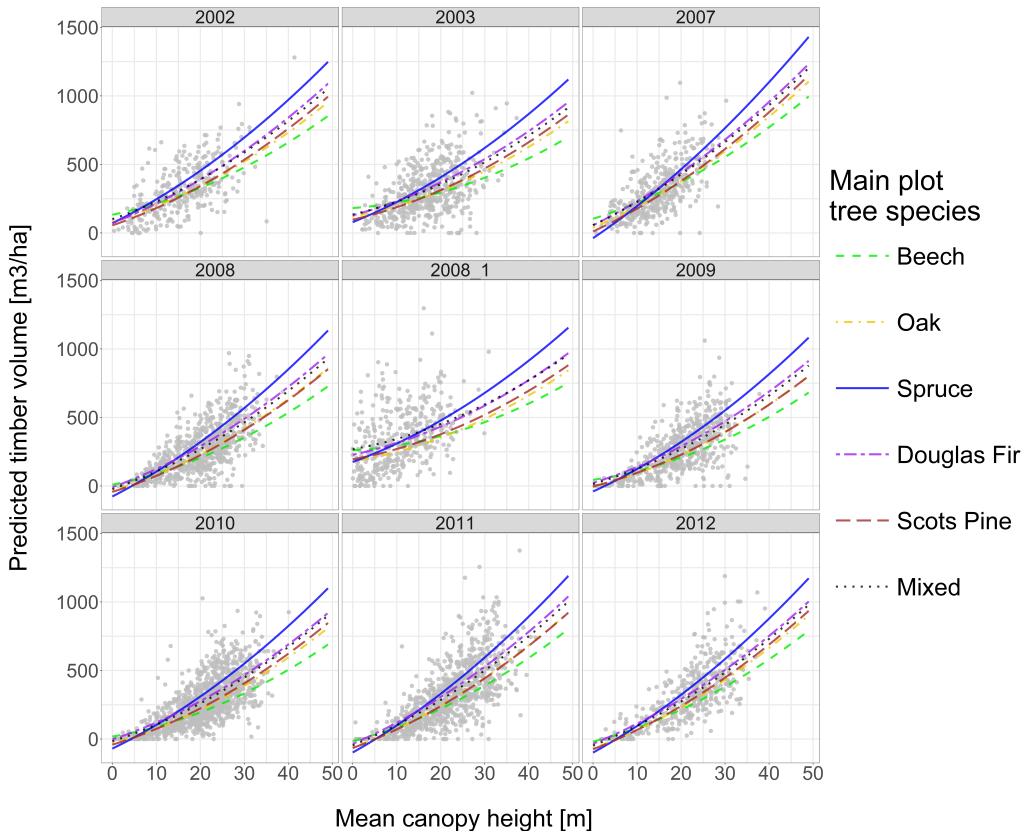


Figure 2.7: Visualization of the timber volume prediction function (*final regression model*) on sample plot level for each main plot tree species and ALS acquisition year. For visualization purposes, the predictor variable *stddev* was set to its average value within the respective *treespecies* and *ALYear* categories. The terrestrially observed timber volume values are plotted in the background.

Table 2.2: Accuracy metrics for submodels of final OLS regression model.  $p$  gives the number of parameters for each model. Interaction terms are indicated by ':':

model terms	model	$p$	$R^2_{adj}$	$RMSE_{cv}$	$RMSE_{cv}\%$
meanheight + stddev + meanheight <sup>2</sup> + treespecies + ALYear + meanheight:treespecies + meanheight:ALYear + meanheight:stddev + stddev:ALYear	final model	39	0.48	140.62	46.69
meanheight + stddev + meanheight <sup>2</sup> + meanheight:stddev	submodel 1	5	0.36	155.54	51.65
meanheight + stddev + meanheight <sup>2</sup> + ALYear + meanheight:ALYear + meanheight:stddev + stddev:ALYear	submodel 2	29	0.45	145.62	48.35
meanheight + stddev + meanheight <sup>2</sup> + treespecies + meanheight:treespecies + meanheight:stddev	submodel 3	15	0.40	150.32	49.92

a higher timber volume than *oak* and *Scots pine* for canopy heights below 20 meters, but realize the lowest timber volumes for canopy heights above 20 metres. Sample plots dominated by any of the remaining coniferous tree species (*Douglas fir*, *spruce*) revealed to have higher slopes than broadleaf classified plots. This indicates that given the same mean canopy height, sample plots dominated by *Douglas fir* and *spruce* yield higher timber volume values than broadleaf- or *Scots pine* dominated sample plots, and this difference becomes more pronounced with increasing mean canopy heights. Within the group of coniferous-dominated sample plots, *spruce* turned out to have the highest slope, thereby yielding the highest timber volume values for mean canopy heights above 15 meters. An undesired characteristic of the model is that the predicted timber volume can in some cases (< 1%) take negative values for low canopy heights (e.g. for *spruce*-dominated plots with *meanheight* below 5 meters and *stddev* of 4 meters). However, we chose not to use a log-transformation of the response variable. Doing so would have prevented the subsequent calculation of the g-weight variance of the design-based estimators (Mandallaz, 2013a; Mandallaz et al., 2013), which is only possible for response variables on the original scale. The g-weight variance provides the benefit of a better variance estimate for internal models by considering the dependency of the regression coefficients on the realized sample. The rare occurrence of negative predictions were however not considered to have an influence on subsequent design-based estimates when averaging multiple predictions within given spatial domains.

### Effect of Time-Lags and Heterogeneity in ALS Data

Incorporating the ALS acquisition year as a categorical variable (*ALSpyear*) in the regression model substantially accounted for the variability in the data introduced by *a*) the time-lags between ALS acquisition and terrestrial survey, and *b*) variation in ALS data quality which are due to sensor- and post processing techniques (table 2.2). Whereas the adjusted  $R^2$  for the regression model without considering the ALS acquisition year as additional predictor variable (*submodel 1*) was 0.36, it could already been increased to 0.40 by including the tree species variable (*submodel 2*). A further stratification by the ALS acquisition year increased the adjusted  $R^2$  of *submodel 1* from 0.36 to 0.45, and the adjusted  $R^2$  of *submodel 3* from 0.40 to 0.48.

We further analysed the model residuals within each ALS acquisition year (within-group variation) for the final model and nested submodels. It turned out that the  $R^2$  values vary distinctly between the ALS acquisition year strata (table 2.3). More precisely, the within-group  $R^2$  can be higher and lower than the overall  $R^2$  of the respective model. Figure 2.8 shows that a stratification according to the ALS acquisition years (*submodel 2*) can already increase the  $R^2$  in most acquisition year strata, compared to the basic model using only the ALS height metrics as predictor variables (*submodel 1*). In the ALS acquisition year stratum 2007, the increase in  $R^2$  even reached 0.08.

Table 2.3:  $R^2$ , RMSE and RMSE% of final regression model within ALS acquisition year strata (*ALSpyear*).  
*Area<sub>ALSpyear</sub>*: Area covered by ALS acquisition given in km<sup>2</sup>. *n*: number of validation data.

<i>ALSpyear</i>	<i>Area<sub>ALSpyear</sub></i>	$R^2$	RMSE	RMSE%	<i>n</i>
2012	2807	0.61	135.84	44.87	408
2011	4361	0.57	146.21	48.29	883
2010	4182	0.51	120.90	39.93	1171
2009	2100	0.42	133.42	44.07	559
2008	2968	0.48	130.38	43.06	701
2008_1	2116	0.33	175.43	57.94	394
2007	3498	0.46	136.47	45.08	418
2003	602	0.27	154.48	51.02	529
2002	775	0.44	141.55	46.75	314

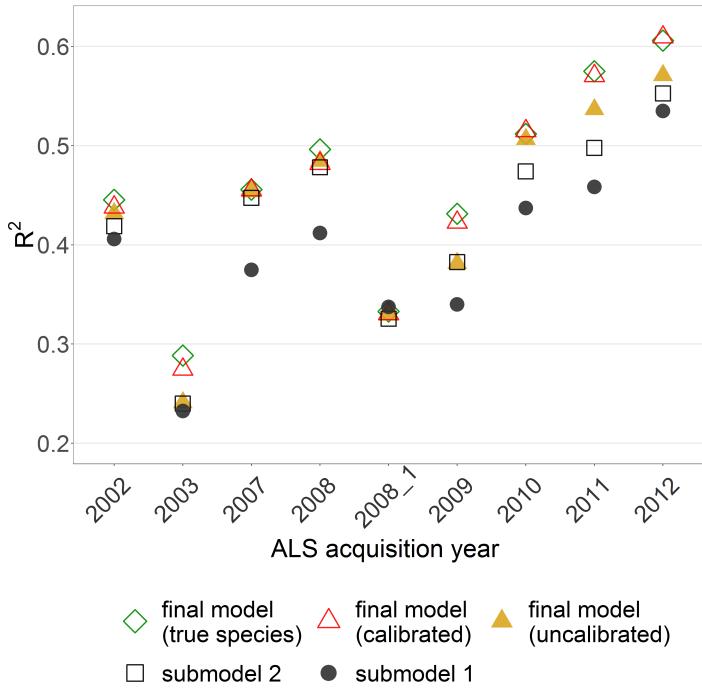


Figure 2.8:  $R^2$ -values of the final regression model, submodel 1 and submodel 2 achieved *within* the ALS acquisition year strata.

### Added Value of Tree Species Map Information

Introducing the predicted main tree species of a sample plot as an additional categorical variable to submodel 2 yielded a further increase in the adjusted  $R^2$  of 0.03 (table 2.2). However, the improvement was even more pronounced in ALS acquisition years close or identical to the year of the terrestrial inventory (figure 2.8). We observed an increase of 0.06 in  $R^2$  for ALS acquisition year 2012, and of 0.07 for ALS acquisition year 2011. The analysis illustrated once more that misclassifications in the tree species variable generally reduce model accuracy compared to using error-free tree species information. The residual inflations caused by the misclassifications in the uncalibrated *treespecies* variable within the *ALSpyear* strata were up to 0.05 in  $R^2$ . However, the calibration was able to substantially decrease or even remove the effects of misclassifications on the model accuracy in all ALS acquisition year strata.

## 2.4 Discussion

### 2.4.1 Stratification according to ALS Acquisition Years and Tree Species

Incorporating the main tree species of a sample location in the timber volume regression model increased the model accuracy and revealed strong evidence for the existence of a tree species specific behavior concerning timber volume on the plot level. This result seems reasonable regarding the species specific taper functions on single-tree level applied in the BWI3 (Kublin, 2003; Kublin et al., 2013). These findings also agree with those of Latifi et al. (2012) who found an almost identical improvement in RMSE of 2% when stratifying to broadleaf and coniferous tree species. The overall RMSE of their model was however 10% smaller than in our study. This might be due to

a more heterogeneous dataset of much smaller sample size, but also because the temporal alignment between the auxiliary data acquisition and the terrestrial survey was much better. Additionally, the number of different tree species present in their dataset was lower than in our case and only comprised Scots pine, European beech and oak. The individual effects of spruce and Douglas fir indicated by our model also support the findings of Breidenbach et al. (2008), who found a higher percentage of coniferous trees in a sample plot to increase the timber volume predictions. This was not true for Scots pine and oak whose effects turned out to be very similar for our dataset. However, in our study the stratification according to the ALS acquisition years severely limited the flexibility of species-specific prediction functions and model interpretability. In particular, using the ALS acquisition years as categorical variables led to highly unbalanced datasets when stratifying according to the main plot tree species. This prevented the use of further stratification variables such as bioclimatic growing regions due to confounding effects and consequent singularities in the design matrices. It also implied an artificial increase in the number of parameters in the OLS regression model, which was however not regarded as critical with respect to overfitting issues due to the high amount of observations used for fitting the regression coefficients (Draper & Smith, 2014, Ch. 15.1). A stratification to the ALS acquisition years however proved to be an effective means in accounting for the artificially introduced noise in the data caused by quality variations and the large time-lags between the remote sensing and terrestrial data. It allowed for a model accuracy that was very close to those reported by Maack et al. (2016) who conducted a very similar study in the German federal state of Baden-Württemberg. Model accuracies were also particularly higher in ALS acquisition year strata in which the data showed considerably less noise or were closer to the date of the terrestrial survey. This effect was significantly reduced or even removed when merging several ALS acquisition year strata. Promising steps with respect to more up-to-date canopy height information have already been made, as the topographic survey institution of RLP is currently processing a canopy height model from aerial imagery acquisitions for 2011 and 2012 covering the entire federal state. These aerial photography acquisitions will in the future be conducted in a two-year period, allowing to derive up-to-date canopy height information in the framework of future forest inventories. For a smaller study area, Kirchhoefer et al. (2017) have already demonstrated that similar model accuracies for German NFI data can be achieved using imagery-based canopy height models.

Incorporating the calibrated tree species information further improved the model accuracy by 0.03 in adjusted  $R^2$ . Compared to the simple model only containing ALS height metrics, including the ALS quality and calibrated tree species information increased the adjusted  $R^2$  by 0.12 in total. A differentiated evaluation of the final regression model revealed that the highest  $R^2$ -values were achieved within ALS acquisitions year strata close or identical with the year of the terrestrial survey, showing differences of up to 0.3 between the  $R^2$ 's. Also the gain in  $R^2$  by including the tree species information was largest (i.e. 0.07) in combination with ALS information acquired in the year of the terrestrial inventory. These insights were particularly interesting with respect to the further use of the regression model for small area estimations. Small area estimators generally gain modeling strength by defining the prediction model *globally* (i.e. using all data in the inventory area), and then applying the so-derived prediction model to a subset of observations located within the area of interest (Mandallaz, 2013a). Consequently, the proposed stratification technique in the prediction model is expected to yield a gain in model accuracy and a reduction of the small area estimation errors if the small area domain mostly includes data from strata that have high within-strata model accuracies. Findings of Breidenbach et al. (2008) indicated that a further increase in model accuracies could possibly be achieved when incorporating these categorical variables as random rather than fixed-effects in linear mixed-effects models (Pinheiro & Bates, 2000). The reason we did not apply this family of models was that small area regression estimators subsequently applied in RLP (Mandallaz, 2013a; Mandallaz et al., 2013) require the internal models to be fitted by OLS technique.

## 2.4.2 Calibration of Tree Species Map Information

The accuracy assessment of the initially derived main plot species from the classification map revealed the presence of misclassifications that led to a decrease in model accuracy. This is in agreement with the potential effects of erroneous explanatory variables discussed in Carroll et al. (2006) and Gustafson (2003), i.e. an increase of variability (noise) in the data that can increase the amount of unexplainable variance and thereby reduce the model accuracy. One reason for the misclassifications were that the classification algorithm of Stoffels et al. (2015) was exclusively trained in pure stands with the objective to predict the *dominant tree species* of a forest stand. Thus, our requirements on the classification map differed considerably from the ones imposed by Stoffels et al. (2015) and have to be considered as far more difficult to meet. Firstly, the reference data used in the accuracy assessment also included understory trees that were recorded in the BWI3 sample. Secondly, determining an exact spatial validation unit for a sample location (support) is not possible due to the properties of angle count sampling (section 2.2.4). Thirdly, distinct discrepancies in the spatial scale between the reference data and the classification map severely hamper exact predictions of the main plot tree species especially in mixed forest stands. The latter issue caused a pronounced dependency of the user's accuracy on the support and threshold choice, particularly for tree species that most commonly occur in mixed forest structures, i.e. *Scots pine* (91%), *oak* (90%) and *beech* (85%) (von Thünen-Institut, 2014). With respect to this set-up, the application of our calibration method proved to be of high value. It led to an increase in the classification accuracies, particularly for those tree species that performed worse in the uncalibrated setup, and thereby successfully minimized and even removed the deleterious effect of misclassifications on model accuracy and regression coefficients. Whereas the extensive analysis in our study deepened the understanding of the afore mentioned scale-effects, an alternative method for future applications could be to use map-derived percentages of each tree species as predictor variables in the random forest algorithm in order to directly predict the terrestrially observed main plot tree species.

## 2.4.3 Choice of Support under Angle Count Sampling

The validation of different support sizes underlined that the support choice can impact the accuracy of a prediction model, and thus confirmed the findings of Deo et al. (2016). In the present study, differences in the model accuracies however turned out to be small for most support choices. An exception was the choice of the  $q100$  support for the CHM derived variables (76 meter radius), where the model accuracy was considerably worse than under the optimal settings. Contrary to our hypothesis, the use of plot-individual supports did not yield the best prediction performance overall. Kirchhoefer et al. (2017) recently came to the same result when they transferred the angle-count sampling technique to a pixel-wise selection method of the auxiliary data that resembles the sample tree selection even more precisely. In their study, the application of fixed support sizes did also not perform worse than under variable supports. We consider two plausible reasons for the joint findings: first, the determination of an exact spatial extent that can be transferred to auxiliary data extraction remains technically infeasible under angle count sampling. Thus, angle count sampling does not seem to be adequate when linking inventory information with remote sensing data. Secondly, inaccuracies in the DGPS-measurements of the plot center locations as reported by Lamprecht et al. (2017) may have an increased impact on the model accuracy the more exact the auxiliary data derivation spatially corresponds to those of the field survey. However, the extensive analysis carried out in our study also indicated that the optimal support size does not only depend on the spatial extent of the field plots, but also on the spatial resolution of the remote sensing data as well as the context in which the derived information is used in the prediction model. In the case of transforming the tree species information map into a suitable categorical

predictor variable, the use of a large support size of 76 meter radius turned out to yield the best model accuracy. However, only few sample locations in the study area were actually characterized by limiting circles of that particular size. An analysis to find the best support settings therefore seems to be advisable prior to further applications of design-based or model-dependent inventory methods so as not to lose model accuracy by unsuitable support choices. The concept of the demonstrated analysis method for identifying suitable supports can be transferred to any kind of auxiliary information, predictor variable and prediction model.

## 2.5 Conclusion

We draw three major conclusions from our study: (1) our analyses strongly indicated that the acquisition of auxiliary data close to the date of the terrestrial survey is a key factor to achieve good model accuracies. Particularly for large-scale inventory applications, this requirement is often difficult to meet. In such cases, we consider that the proposed method of including quality information about the auxiliary data in a prediction model can be an effective technique for improving the prediction accuracy. Ongoing studies investigate whether this modelling technique can also lead to smaller estimation errors of design-based estimators. (2) Our study also indicated that the relationship between the field measured timber volume and remote-sensing derived height information is tree species specific. We expect that using the tree species information in a timber volume model would even lead to higher prediction accuracies when combined with explanatory variables that can further explain the variation within each tree species group, such as bioclimatic growing conditions, soil properties and stand density on the plot level. (3) We consider the demonstrated calibration technique to be a valuable method for future studies where an external tree species map (i.e. the map was not created for the specific study objective) is used in prediction models. The application of a calibration model can also be transferred to any error-prone explanatory variable and be a simple means to clean the data set from noise and thus increase the model accuracy.

## Acknowledgements

We want to express our gratitude to Prof. H. Heinimann (Chair of Land Use Engineering, ETH Zurich) for supporting this study. We want to explicitly thank Dr. Johannes Stoffels from the Environmental Sensing and Geoinformatics Group of Trier University for providing the tree species classification map as well as for constructive discussions when it came to interpreting the results. Special gratitude is owed to the State Forest Service of Rhineland-Palatinate, in particular Dr. Joachim Langshausen, Jürgen Dietz and Claus-Andreas Lessander, for collaboration and providing the forest inventory and geodata. We also want to thank Kai Husmann and Christoph Fischer from the Northwest German Forest Research Institution Göttingen for their advice in processing the terrestrial inventory data, and Alexander Massey and Michael Hill for proofreading.



## **Chapter 3**

# **A double-sampling extension of the German National Forest Inventory for design-based small area estimation of timber volume resources on forest district levels**

Andreas Hill<sup>1</sup>, Daniel Mandallaz<sup>1</sup>, Joachim Langshausen<sup>2</sup>

<sup>1</sup>ETH Zürich

Department of Environmental Systems Science, Universitätstrasse 16, 8092 Zürich, Switzerland

In progress

- Daniel Mandallaz developed the design-based estimators. He also supported the statistical analysis and writing of the manuscript.
- Joachim Langshausen supported writing of the manuscript.



## Chapter 4

# Accuracy Assessment of Timber Volume Maps using Forest Inventory Data and LiDAR Canopy Height Models

Andreas Hill<sup>1</sup>, Jochen Breschan<sup>1</sup>, Daniel Mandallaz<sup>1</sup>

<sup>1</sup>ETH Zürich

Department of Environmental Systems Science, Universitätstrasse 16, 8092 Zürich, Switzerland

Published in:  
*Forests* 5 (2014): 2253-2275  
(DOI: 10.3390/f5092253)

- Jochen Breschan supported the formulation of the optimization model.
- Daniel Mandallaz supported the statistical data analysis.

## Abstract

Maps of standing timber volume provide valuable decision support for forest managers and have therefore been the subject of recent studies. For map production, field observations are commonly combined with area-wide remote sensing data in order to formulate prediction models, which are then applied over the entire inventory area. The accuracy of such maps has frequently been described by parameters such as the root mean square error of the prediction model. The aim of this study was to additionally address the accuracy of timber volume classes, which are used to better represent the map predictions. However, the use of constant class intervals neglects the possibility that the precision of the underlying prediction model may not be constant across the entire volume range, resulting in pronounced gradients between class accuracies. This study proposes an optimization technique that automatically identifies a classification scheme which accounts for the properties of the underlying model and the implied properties of the remote sensing support information. We demonstrate the approach in a mountainous study site in Eastern Switzerland covering a forest area of 2000 hectares using a multiple linear regression model approach. A LiDAR-based canopy height model (CHM) provided the auxiliary information; timber volume observations from the latest forest inventory were used for model calibration and map validation. The coefficient of determination ( $R^2 = 0.64$ ) and the cross-validated root mean square error ( $RMSE_{cv} = 123.79 \text{ m}^3/\text{ha}$ ) were only slightly smaller than those of studies in less steep and heterogeneous landscapes. For a large set of pre-defined number of classes, the optimization model successfully identified those classification schemes that achieved the highest possible accuracies for each class.

## 4.1 Introduction

### 4.1.1 Context and Problem

Among the multitude of information that forest inventories are expected to provide (McRoberts et al., 2010), knowledge about available standing timber volume on the national, regional, as well as enterprise level is still of high interest. Since on these spatial levels, a full census is too cost-intensive and, in most cases, even practically unfeasible, a broad range of methods in the framework of sampling theory has been developed and applied to estimate this quantity (Gregoire & Valentine Harry, 2008; Mandallaz, 2008b; Schreuder et al., 1993). The strength of forest inventory methods relying on design-based procedures is that (at least asymptotically) unbiased point and variance estimates can be obtained, and this without assuming the applied prediction models to be correct in the classical statistical (model-dependent) sense. An important advancement in increasing this accuracy without, at the same time, increasing the number of costly terrestrial samples has been achieved by combining terrestrial samples with auxiliary information provided by remote sensing data, so-called two-phase or double-sampling procedures (Gregoire & Valentine Harry, 2008; Mandallaz, 2008b; Cochran, 2007; Köhl et al., 2006). In this context, especially airborne laser scanning (ALS) data have proven to provide a high degree of information for timber volume estimation (Holmgren, 2004; ?; Næsset, 2002). It has recently been shown that the efficiency of two-phase sampling can be further increased by extending this procedure to stratification (Saborowski et al., 2010; von Lüpke, 2013) or by using part of the auxiliary information exhaustively when remote sensing data are covering the entire inventory area (Mandallaz et al., 2013). The two-phase procedure is thus not restricted to large forest areas, but has also been applied in the context of small area estimation (Breidenbach & Astrup, 2012). Given that the number of terrestrial samples in the small area is sufficiently large (i.e., one is not restricted to the application of synthetic estimations), even for small areas the accuracy specifications are ensured to be unbiased (Mandallaz, 2013a). While these forest inventory methods have the advantage of supplying reliable accuracy specifications for their estimates, they do not provide information about the spatial distribution of

the estimated quantity. However, the availability of spatially explicit stand information is of prime importance for efficiently locating forest management operations.

Accordingly, mapping the spatial distribution of standing timber volume has been the subject of various recent studies. The statistical models that have been used for mapping can be divided into parametric models, particularly linear regression models (Tonolli et al., 2011; Van Aardt et al., 2008), and non-parametric models (Franco-Lopez et al., 2001; Latifi et al., 2010; Nothdurft et al., 2009). Among the non-parametric models, k-NN imputation has become increasingly popular due to its simplicity and easy implementation (Magnussen et al., 2014). k-NN approaches have been investigated and applied in the model-dependent framework of forest inventory with promising results (McRoberts et al., 2007) and have also been used for the mapping of various forest attributes (Beaudoin et al., 2014; Chirici et al., 2012; Tomppo, 2006). Haara & Kangas (2012) compared the k-NN method to linear regression in a simulation study and found the two methods to perform similarly well. Especially in the case where the relationship between observations and the auxiliary variable followed a linear trend, the regression model performed better than the k-NN approach. Fehrmann, et al. (Fehrmann et al., 2008) came to a similar result when comparing linear and linear mixed effect models to an instance-based k-NN approach for single-tree biomass estimation. Also in their case, the performance of the k-NN approach and the linear mixed model only differed marginally, and both methods were slightly superior to simple linear regression. On the other hand, they also confirmed that the application of k-NN methods can be an effective and promising method if no a priori knowledge about the relationship between target and auxiliary variable(s) exists, particularly if the relationship is considered to be complex due to random and interaction effects. However, they also raised the question of whether a k-NN approach should be used in situations where the functional relationship among variables is approximately known. Additionally, the performance of k-NN estimation and its potential superiority to already existing methods has also been investigated in the design-based framework of forest inventory (Baffetta et al., 2010, 2009). While in several cases, the proposed k-NN estimator of the population mean achieved smaller errors than the Horwitz-Thompson estimator (Mandallaz, 2008b), the result also turned out to be dependent on the underlying model of the investigated population.

Irrespective of the model choice, a core issue of these mapping approaches is to characterize the accuracy of the resulting maps. If map predictions are made on a continuous prediction scale, the map precision is commonly characterized by quality parameters of the applied prediction model, such as the cross-validated root-mean-squared error (RMSE) and the coefficient of determination ( $R^2$ ). However, the derived map predictions are often visualized using constant class intervals in order to provide users, such as forest managers, with a better visual interpretation of the map. In this case, it could be misleading to still use the previously mentioned RMSE and  $R^2$  in order to provide information about the accuracy of resulting timber volume classes. This is because these parameters only describe the overall model performance on a continuous prediction scale, but do not quantify the accuracy of individual timber volume ranges (classes). A more appropriate validation strategy would then be to adopt the concept of confusion matrices, which provides a differentiated accuracy assessment (user's and producer's accuracy) for each particular volume class, as well as the complete mapping system. Franco-Lopez et al. (2001), for example, used these metrics. However, their classification scheme of constant class intervals exhibited a strong gradient of degrading class accuracies towards higher volume classes (most likely due to saturation effects in the remote sensing data). Such a severe gradient in class accuracies, however, reveals the following problems: (1) it implies that the chosen classification scheme with constant class intervals is not accounting for the fact that the performance of the underlying model may not be constant across the entire volume range; and (2) it severely hampers the usability of the maps in forest practice due to the high uncertainty within higher timber volume classes.

The motivation of this study was to improve the usability of volume maps for forest management operations by avoiding classification schemes of this kind. We hypothesized that this can be

achieved by optimizing the class intervals with respect to the accuracy potential of the underlying prediction model. This implies using smaller class intervals in those volume ranges where the model ensures precise prediction performance and enlarging these intervals in ranges where the model performs worse. If the class boundaries are allocated according to this concept, it becomes possible to design classification schemes that provide highest possible accuracies for each class, while avoiding a severe gradient between class accuracies. This concept was investigated by implementing an optimization algorithm which can be applied to any type of prediction model that provides estimates on a continuous scale. Implicitly, the method also provides an additional option for evaluating the precision of prediction models.

We demonstrate the method in a case study in the canton of Grisons using a LiDAR-derived canopy height model and regional forest inventory data. The workflow included: (1) the production of a map of estimated standing timber volume for the entire study area; (2) the calculation of reliable accuracy metrics for this map; and (3) the application of the proposed optimization algorithm in order to identify the classification schemes which provide the highest possible accuracies. In this particular case, we decided to use a multiple linear regression model, because most auxiliary variables exhibit a pronounced linear relationship to the terrestrial inventory (Hill, 2013). Additionally, the number of available terrestrial observations in our study was small ( $n = 67$ ) compared to similar studies, whereas it has been indicated that a good performance of k-NN requires larger datasets (Fehrman et al., 2008; Magnussen et al., 2010).

#### 4.1.2 Background on Heuristic Search Methods

Heuristic Search Methods (HSMs) are often used when coping with combinatorial optimization problems (Pirlot, 1996) or, in general, problems whose structure cannot be satisfactorily represented and performed by means of classical optimization techniques, such as Linear Programming (Rayward-Smith et al., 1996). Basically, HSMs aim at improving an objective function by subsequent inspection and adoption of neighboring solutions. Inspection rules are often inspired by nature (Pirlot, 1996) and have the property of occasionally accepting inferior solutions for further inspection to avoid getting trapped within a local minimum or maximum. Simulated Annealing (SA) (Kirkpatrick et al., 1983) is such an HSM, borrowing its accepting rule for inferior solutions from metallurgy. It is based on the assumption that a configuration of atoms in a metal can move to configurations of higher energy (i.e., inferior solution) with a certain probability at a given temperature. Given that probability is a function of temperature and energy difference to the inferior solution, SA adopts a cooling scheme that aims at annealing the metal to the point of minimum configuration energy (i.e., objective function). As opposed to classical optimization methods, optimality of HSM-derived solutions cannot be proven. However, one can assume to find a solution close to the true but unknown optimum if the heuristic is appropriately parameterized.

## 4.2 Materials and Methods

### 4.2.1 Materials

#### Study Area

The methods proposed in this article were applied to a study site located in the canton of Grisons, Eastern Switzerland (Figure 1). The site extends in the north-south direction between Klosters and Davos and covers a total area of 2887.39 hectares. According to a forest mask (in raster format) of the study site derived by the use of the Swiss TLM3D (Swiss Topographic Landscape Model) data with the approval of the Swiss Federal Institute of Topography (for details, see Hill (2013)), the forest area of the study site comprises 1974.49 hectares (68.4%). The study site is located

at an altitude between 900 and 2200 meters above sea level, and its relief is mainly characterized by rough terrain and steep slopes. Classifying the forest area of the study site according to the scheme given by (Ott et al., 1997) revealed 49.7% of the forest area belonging to the high montane vegetation zone, 49.5% to the sub-alpine zone and 0.8% to the upper sub-alpine vegetation zone. Consequently, the forests within the study area were assumed primarily to consist of coniferous tree species, especially Norway spruce (*Picea Abies*).

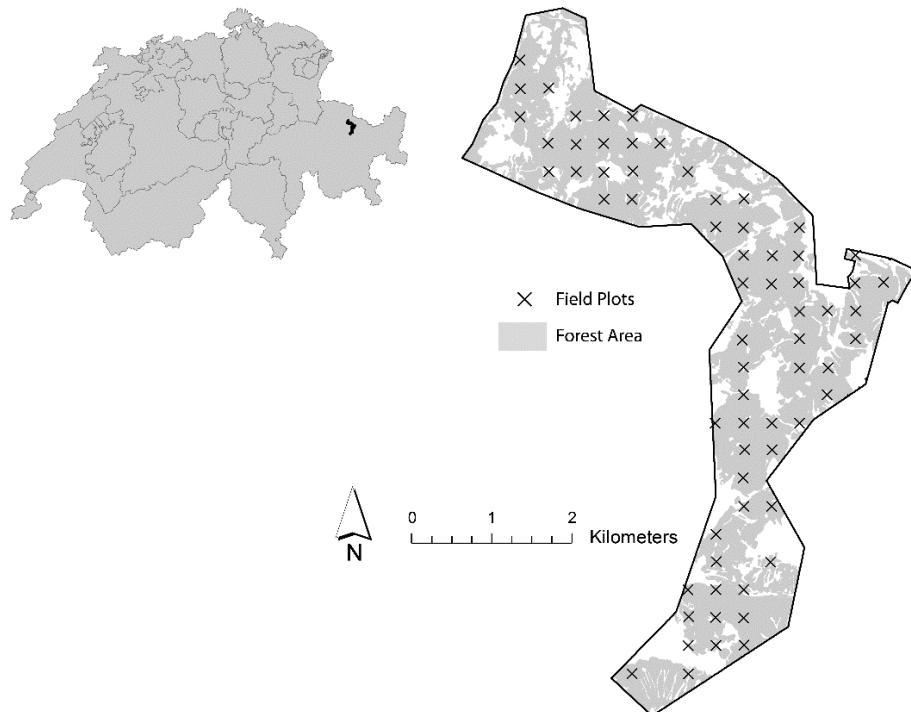


Figure 4.1: Study site, including the distribution of the 67 field plots (regional forest inventory) that are part of the forest area derived from the TLM3D (Topographic Landscape Model) data (with approval of swisstopo JA100120/JD100042).

### **Timber Volume Densities from Field Inventory**

The methods applied in this study are based on timber volume densities obtained from terrestrial field surveys. For the study site, the timber volume densities were provided from field surveys of the latest regional forest inventory at the canton of Grisons and provided by the canton's forest service. Any uncertainty of the timber volume densities associated with their acquisition was ignored (e.g., potential measurement errors). The forest area of the study site comprised 67 terrestrial plots, which had been surveyed in the year 2007 (Figure 4.1). Part of the survey was also the re-measurement of the plot centers with GPS technique. Unfortunately, no reliable information about the positional accuracy could be provided. The field surveys were conducted using circular sample plots with their center (i.e., the sample point) belonging to a regional permanent systematic inventory grid with a mesh size of 500 m. This regional sampling scheme thus constitutes a sub-grid of the nation-wide terrestrial inventory grid of the National Swiss Forest Inventory (NFI) with a mesh size of 1.4 km. Each sample plot consists of two concentric circles with a plot area of 200 and 500 m<sup>2</sup> around the sample center. Within the inner circle (radius of 7.98 m), all trees with a diameter at breast height (DBH) larger than 12 cm were selected, whereas in the second circle (radius of 12.62 m), all trees with a DBH larger than 36 cm were included in the sample. Boundary and slope adjustments were performed on plot level. The explicit survey methods and

the evaluation of the regional inventory surveys were identical to those of the NFI and can be found in detail in Brassel & Lischke (2001) and Keller (2011). To obtain the standing timber volume on plot level, the overbark timber volume of each sample tree was estimated by measuring its DBH and using it as the main predictor variable in the tariff models provided by the NFI. These tariff models are based on the general function proposed by Hoffmann (1982) and have been extended by further explanatory variables, such as the production region and additional tree and plot attributes (Brassel & Lischke, 2001). The standing timber volume for each plot was then estimated according to the Horwitz-Thompson estimator, which provides an unbiased estimation of the actual timber volume on plot level (Mandallaz, 2008b). The volume distribution of the 67 terrestrial observed field plots is illustrated in Figure 4.2, and a brief statistical summary is given in Table 4.1.

Table 4.1: Summary statistics of the timber volume observed at the 67 terrestrial sample plots (given in  $\text{m}^3/\text{ha}$ ).

Range	Mean	Median	SD	n
7.3 - 869.57	399.4	386.9	194.94	67

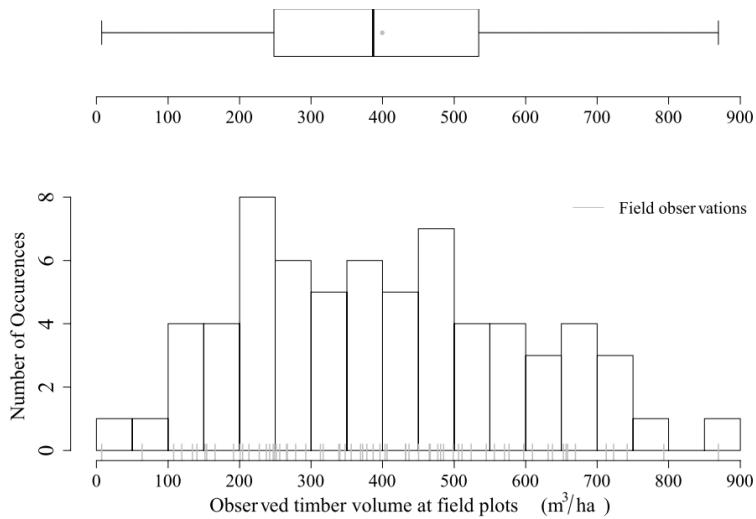


Figure 4.2: Terrestrial observed timber volume on plot level for all 67 sample plots of the study area (histogram bandwidth =  $50 \text{ m}^3/\text{ha}$ ; origin =  $0 \text{ m}^3/\text{ha}$ ).

## LiDAR Data

A LiDAR dataset covering the entire study area was acquired with a Riegl LMS Q 560 laser scanning system in the period September 11-15, 2010. The LiDAR acquisition was conducted as a part of a campaign to gather data for the Swiss National Park. A digital terrain model (DTM) and a digital surface model (DSM) with a spatial resolution of 0.5 m were computed by the provider Toposys by application of their company-internal processing software TopPit. Gaps in the DTM due to the absence of last echoes had not yet been interpolated. The average flight height was 700 m above ground, and the average echo density was 27 points  $\text{m}^2$ . The provider specified the positional accuracy as  $< \pm 0.50 \text{ m}$  and the height accuracy as  $< \pm 0.15 \text{ m}$ . Further specifications of the LiDAR acquisition are summarized in Table 4.2.

Table 4.2: Summary statistics of the timber volume observed at the 67 terrestrial sample plots (given in m<sup>3</sup>/ha).

Beam deflection	Rotating mirror
Pulse Repetition Frequency (kHz)	70
Average Flying Altitude (m above ground)	700
Max. scan angle (°)	± 15
Wavelength (nm)	1550
Beam divergence (mrad)	≤ 0.5
Average echo density (m <sup>-2</sup> )	27.4

#### 4.2.2 Methods

The conceptual model in Figure 4.3 captures the general workflow of creating the timber volume map for the study area and the subsequent accuracy assessment. It consists of the following steps.

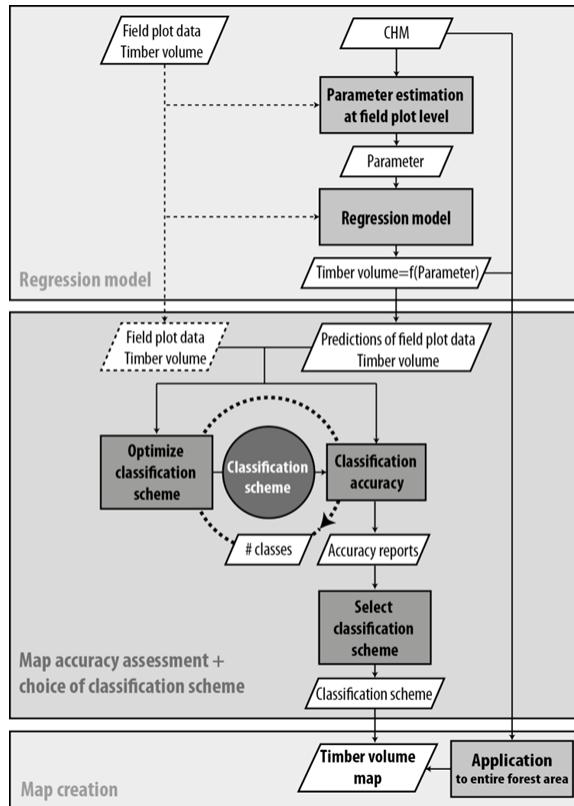


Figure 4.3: Conceptual model illustrating the workflow of producing the timber volume map and assessing its accuracy.

Step I: A regression model for predicting the standing timber volume in m<sup>3</sup>/ha is formulated using the terrestrial observed timber volume of the field plots as the response variable and parameters extracted from the canopy height model (CHM) at the respective plot locations as predictor variables. The model then allows for predicting the timber volume at each point of the CHM that lies within the inventory area.

Step II: The model predictions (Step I) at the plot locations are then compared to the corresponding field data. These comparisons are used to assess accuracy metrics of the timber volume map. Regarding the classification accuracies, which can be estimated under class representations of the timber volume map (i.e., timber volume classes), this part also comprises the application of an optimization model aiming to find an optimal classification scheme, i.e., the choice of class boundaries and class widths for a given number of classes.

Step III: The regression model (Step I) is applied to the entire study area, and the result is then represented using a classification scheme that best satisfies the required accuracies for map users (Step II).

### **Step I: Computation of Canopy Height Model**

A canopy height model (CHM) with a spatial resolution of 0.5 m, completely covering the study site, was calculated by subtracting the LiDAR-acquired digital terrain model (DTM) from the digital surface model (DSM). As the DSM represents the elevation characteristics of the surface including vegetation and man-made structures, whereas the DTM describes only the elevation of the terrain, this operation equals removing the underlying terrain information from all features in the DSM. The height information of all objects in the CHM hence describes their estimated object height. Before calculating the CHM, an interpolation step for the DTM was crucial, since it exhibited a considerable amount of missing height information-most likely due to the absence of last pulse LiDAR returns over densely-covered forest areas. Missing height values in the DTM were predicted by applying an inverse distance weighting (IDW) interpolation algorithm [38]. Due to a locally varying number of adjacent missing raster values, the IDW algorithm was iteratively applied five times using varying neighborhood distances within which available height values were considered for interpolation. The use of small neighborhood distances aimed at providing a high amount of local precision where sufficient height values are available in the direct neighborhood of a missing value, whereas the use of large neighborhood distances was necessary for the interpolation of large gaps of missing values. Starting with a threshold distance of 2 neighbors, i.e., 1 m, and iteratively increasing the distance up to 20 neighbors, i.e., 10 m, all missing height values of the DTM were replaced by height predictions.

### **Step I: Regression Model**

To predict the standing timber volume (TV) at location  $x$ , a multiple regression model was formulated using the observed timber volume of the field plots as the response variable and certain CHM metrics as predictor variables. The CHM metrics were extracted at all 67 field plots within squares of constant size, which were centered at the respective plot centers. To ensure high spatial consistency between the circular field plots and the CHM metrics, the square extent was chosen in order to tangentially circumscribe a field plot. The side length of a square was 25 m, resulting in an area of 625 m<sup>2</sup> compared to the field plot area of 500 m<sup>2</sup>. In the following, these squares are also referred to as the support of the estimates. By analyzing the distribution of the raster values of the CHM within each square, we calculated the following metrics of the LiDAR CHM at each plot location: the MEAN, the standard deviation (SD), the maximum value (MAX), as well as the 25%, 75% and 90% quantiles (Q25, Q75, Q90). The reason for the choice of these parameters was that they have often been used as predictors for estimating timber volume of forest stands (Holmgren, 2004; ?; Lefsky et al., 1999; Magnussen et al., 1999). The calculations of these variables were also adjusted for boundary effects by only considering those raster values within a square that are covered by the forest mask (Section 4.2.2). Table 4.3 provides the main statistics for the observed CHM metrics at the 67 terrestrial sample plots. The ordinary least square regression model

containing the maximum number of predictor variables conclusively reads as:

$$TV(x) = \beta_0 + \beta_1 MEAN(x) + \beta_2 SD(x) + \beta_3 MAX(x) + \beta_4 Q25(x) \\ + \beta_4 Q75(x) + \beta_4 Q90(x) + \varepsilon(x) \quad (4.1)$$

A variable selection procedure was applied in order to restrict the regression model to the most meaningful variables and to avoid an overfitting effect of the model (Draper & Smith, 2014) due to the small number of observations. As most of the predictor variables were considered to be correlated to each other (collinearity), criterion-based selection procedures, by means of AIC (Akaike information criterion) (Akaike, 1992) and adjusted R-square criteria (Srivastava et al., 1995), as well as Mallow's Cp statistic (Mallows, 2000), were preferred to the testing-based selection procedures relying on p-values.

Table 4.3: Summary statistics of the CHM metrics calculated at the 67 terrestrial sample plots (in meters).

Metrics	Range	Mean	Median	SD
MEAN	2.26-26.03	12.07	11.31	5.89
SD	3.71-15.97	8.93	8.64	2.73
MAX	17.03-45.35	32.63	32.74	7.05
Q25	0-22.88	4.25	0.67	6.35
Q75	1.43-34.21	18.92	18.81	7.96
Q90	8.11-37.78	23.77	23.48	7.25

## Step II: Assessment of Map Accuracy

### *Prediction Performance*

Assessing accuracy metrics of the timber volume map was realized by validating the timber volume predictions made by the regression model at the field plot locations using the corresponding observed timber volume of the field plots as reference data. Since the terrestrial sample was considered too small to be split into separate calibration and validation subsets, a leave-one-out cross-validation was performed to estimate the root-mean-squared error (RMSE) as a measure for the prediction performance of the timber volume map.

### *Classification Accuracy*

We applied the concept of representing the timber volume map by prior defined timber volume classes (i.e., assigning the prediction of each raster cell to a prior defined interval). This representation has commonly been used by various research studies producing maps of forest parameters, such as basal area or timber volume (Tonolli et al., 2011; Latifi et al., 2010; Clementel et al., 2012), in order to facilitate the interpretability, as well as the readability of a map. We propose to treat this map representation as a classification procedure and, consequently, adopted the concept of confusion matrices to provide a differentiated assessment for each particular class as well as the complete mapping system. Using the available field data as reference data, we estimated the following accuracy metrics for the resulting classified timber volume map (for details, see Congalton & Green (2008) and Richards & Richards (1999)):

- The overall accuracy (OAA) is the proportion of correctly classified pixels of the entire map. The true overall accuracy of the map is unknown, since we only have references for the classified raster cells at a small subspace of the map. The OAA is therefore estimated by the ratio

of the total number of correctly classified pixels and the total number of reference/classified pixels. The 95% confidence interval for the OAA was calculated according to the binomial distribution.

- The producer's accuracy (PA) is a measure of the classification performance. It indicates the probability that if a ground observation belongs to a certain class, this class will be reflected in the map. The producer's accuracy can be estimated for each class by dividing the number of correctly classified pixels of a class by the total number of reference pixels in this class.
- The user's accuracy (UA) of each class is the most interesting information for a user of the map. It indicates the probability that if the map shows a certain class, this class will actually be validated by a terrestrial survey. The producer's accuracy is estimated by the number of correctly classified pixels in a class divided by the total number of classified data in this class.
- Cohen's kappa coefficient is a measure to assess to what degree the classification accuracy was realized by a chance agreement. The kappa coefficient ranges between -1 (accuracy was realized under pure chance agreement) and 1 (accuracy was reached by no chance agreement)

### ***Class Selection Problem***

One of the main benefits of classifying the timber volume map and assessing its classification accuracy is to provide information on the accuracies for individual timber volume ranges. However, classifying the model predictions into classes produces the problem of having to choose an appropriate classification scheme, i.e., choosing the class boundaries of the timber volume classes. A classic approach would be to use equally-sized classes with origin at zero and constant class width, but we consider three reasons not to do so: a constant class width (i) is likely to create classes for which no reference data are available, especially if the class width is chosen small; for such classes, PA and UA cannot be estimated, and the overall accuracy would give the user of the map an overoptimistic impression of the actual map precision; (ii) may separate a reference from its prediction (or vice versa) even if the two values were almost identical (i.e., their difference is very small and even negligible from a user's point of view); (iii) does not account for saturation effects in the remote sensing data, occasionally leading to a strong gradient of degrading class accuracies towards higher volume classes. To overcome these problems, we propose a locally-adaptive selection of class boundaries which satisfies the following rules (Class Selection Problem):

**Rule I:** Choose the class boundaries to ensure that each timber volume class at least contains a minimum number of reference data. This also aims at using a smaller class width where a sufficient number of references is available (thus providing locally higher detail), whereas the class width is increased for regions where references are rare.

**Rule II:** Avoid cases where a reference and its (closely located) corresponding prediction are separated by a class boundary. This implies not only taking into account the distribution of the reference data, but also considering the distribution of their corresponding predictions. A slight adaptation of the class width may thereby increase the classification accuracy.

### ***Optimization Model***

The class selection problem (CSP) can be solved by repeatedly moving the boundaries and evaluating the resulting classification schemes until Rules I and II of the CSP are optimally satisfied. Since the number of combinations of alternative classification schemes can become big (e.g.,

2.3 million alternatives to distribute four boundaries along a range of 10-900 m<sup>3</sup>/ha specified by a lower and an upper boundary and discretized into 88 steps of 10 m<sup>3</sup>/ha), it risks becoming too computationally intensive to evaluate all alternatives. For this reason, we formulate the following multi-objective optimization model to automate the design of an (approximately) optimal classification scheme and solve it using Simulated Annealing (Section 4.1.2).

Its variables are specified as follows:

$$x_j : \text{class boundary value } j \text{ (m}^3/\text{ha)} \quad (4.2)$$

$$\text{class } j = (x_j, x_{j+1}), j=1, \dots, m+1 \quad (4.3)$$

$$n_{ref,j} : \text{number references / number plots in class } j \quad (4.4)$$

$$y_{ij} = \begin{cases} 1, & \text{if } V_{tp} \text{ and } V_{pred} \text{ at plot } i \text{ in same class } j \\ 0, & \text{otherwise} \end{cases} \quad (4.5)$$

Where  $V_{tp}$  is the timber volume of terrestrial plot and  $V_{pred}$  the predicted timber volume. The decision variables  $x_j$  define the class boundary positions. A class  $j$  is then defined via the interval which is bound to  $x_j$  and  $x_{j+1}$  (Eq. 4.3). The number of pre-defined classes is  $m$ . The number of reference data for class  $j$  is given by  $n_{ref,j}$  and the binary indicator variable  $y_{ij}$  (Eq. 4.5) captures whether  $V_{tp}$  and  $V_{pred}$  at plot  $i$  are assigned to the same class  $j$  (yes = 1, no = 0).

The optimization model can then be formulated as follows when  $n$  plots are subject to assignment:

$$MAX \sum_{j=1}^m \sum_{i=1}^n y_{ij} - w_1 \sum_{j=1}^m (x_{j+1} - x_j)^2 - w_2 \sum_{j=1}^m \left( \frac{n}{m} - n_{ref,j} \right)^2 \quad (4.6)$$

subject to:

$$x_{j+1} - x_j \geq t_w \quad (4.7a)$$

$$x_j \geq 0, y_{ij} \in \{0; 1\} \quad (4.7b)$$

The objective function (Equation 4.6) implements Rules I and II with three terms. The first term captures Rule II by maximizing for the number of cases where terrestrial observed volumes  $V_{tp}$  and corresponding estimated volumes  $V_{pred}$  are assigned to the same class  $j$ . In principle, the first term can be maximized by proposing a 'super class' covering almost the entire range. The second term is thus introduced to penalize the selection of a 'super class' by minimizing for local class widths, thereby also accounting for Rule I. This is realized by minimizing the square class width for each class  $j$ . Equation 4.7a limits the minimum class width to a threshold  $t_w$  (m<sup>3</sup>/ha) and concurrently satisfies that  $x_{j+1}$  is chosen bigger than  $x_j$ . The third term captures Rule I and aims at equally distributing the reference data over all classes by minimizing the squared difference between the average number of plots per class and the actual number of reference data for each class  $j$ .

Both the second and the third terms are implemented as penalty terms for the first term that aims at maximizing the number of correctly classified plots. The corresponding weights  $w_1$  and  $w_2$  can be used to control the emphasis of both penalty terms and provides the user with the

possibility to prefer one of the three objectives within the optimization process.

In order to find an appropriate choice for the weight factors, several weights have to be evaluated to achieve satisfactory classification schemes. In our study, we decided to give the weights equal emphasis ( $w_1 = w_2 = 2$ ). Class boundary selection was restricted to a range discretized into 10 m<sup>3</sup>/ha units to reduce the problem size and simultaneously create useable class boundaries in the final classification scheme. A satisfying alternative was then identified by picking out the best overall solution from 100 runs of Simulated Annealing, where each run included the computation of 1000 alternatives.

### Step III: Computation of the Timber Volume Map

After model selection, the regression model (Section 4.2.2) was applied at any location  $x$  of the CHM, i.e., over the entire study area. We chose the design of the timber volume map in accordance with the setup of the regression model: the spatial resolution of the map was defined as the size of the support used for ground calibration, i.e., 25x25 m. We further orientated the map in such a way that the supports at the field plots actually became an almost exact subset of all raster cells of the timber volume map. Within each of the raster cells, the CHM metrics used in the regression model were calculated by the same procedure as described in Section 4.2.2 (i.e., by the same support and technique) and then used to predict the timber volume. The estimated value was then assigned to the entire raster cell, resulting in a timber volume map with a spatial resolution of 25x25 m. The entire procedure is again illustrated in Figure 4.4.

The design of the map as proposed here has two main advantages that allow for relying on the provided classification accuracy metrics: (i) each estimate of a raster cell is based on exactly the same support that was used to calibrate the prediction model; and (ii) as the reference data are an (almost) exact subset of the map raster cells, this allows for a valid accuracy assessment of the classified timber volume map (Congalton & Green, 2008). Another advantage of the map design is that it is in perfect agreement with the inventory design of the generalized two-phase regression estimator proposed by Mandallaz et al. (2013), where part of the auxiliary information is derived exhaustively over the entire inventory area. It thereby provides a link between the derived timber volume map and sample designs of classical forest inventory.

Once the timber volume map was calculated for the entire study area, each raster cell of the map, still carrying estimates on a continuous scale, was classified into timber volume classes according to a classification scheme (Step II, Section 4.2.2).

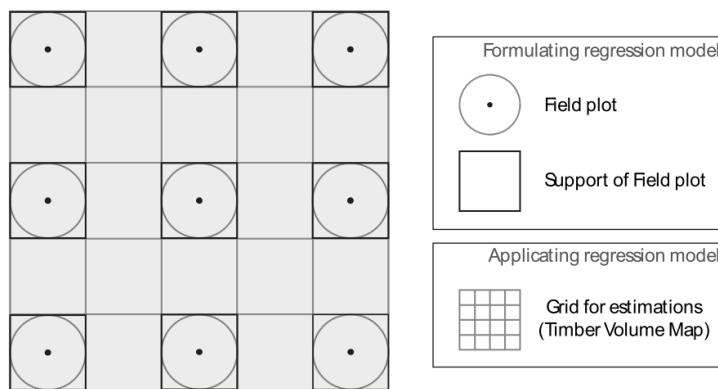


Figure 4.4: Schematic design of the timber volume map.

## 4.3 Results

### 4.3.1 Regression Model

Performing individual simple linear regressions revealed that all predictor variables derived from the CHM were correlated to the field-obtained timber volume on plot level. Here, the highest coefficient of determination ( $R^2$ ) was achieved by the variable mean ( $R^2 = 0.5$ ). A forward, backward, as well as bidirectional selection procedure according to minimize the AIC all revealed mean, standard deviation, maximum value and the 75% quantile as the model of choice (adjusted  $R^2 = 0.62$ , AIC = 646.4). These predictor variables also revealed a significant influence on the field-obtained timber volume (individual parameter t-test, 5% significance level). The model was also suggested by Mallow's Cp and the adjusted R2-criterion selection procedure. Summary statistics of the model are presented in Table 4.4. The estimated timber volume  $\hat{TV}(x)$  at location  $x$  was consequently calculated according to the following regression model formula:

$$TV(x) = 228.83 + 62.76 * MEAN(x) + 76.65 * SD(x) + 19.59 * MAX(x) + 33.44 * Q75(x) \quad (4.8)$$

Figure 4.5 shows the predicted timber volume on plot level plotted against the observed timber volume of the field plots. The leave-one-out cross-validated RMSE<sub>cv</sub> of the regression model was 123.79 m<sup>3</sup>/ha and, thus, only slightly larger than the RMSE without cross-validation (115.5 m<sup>3</sup>/ha).

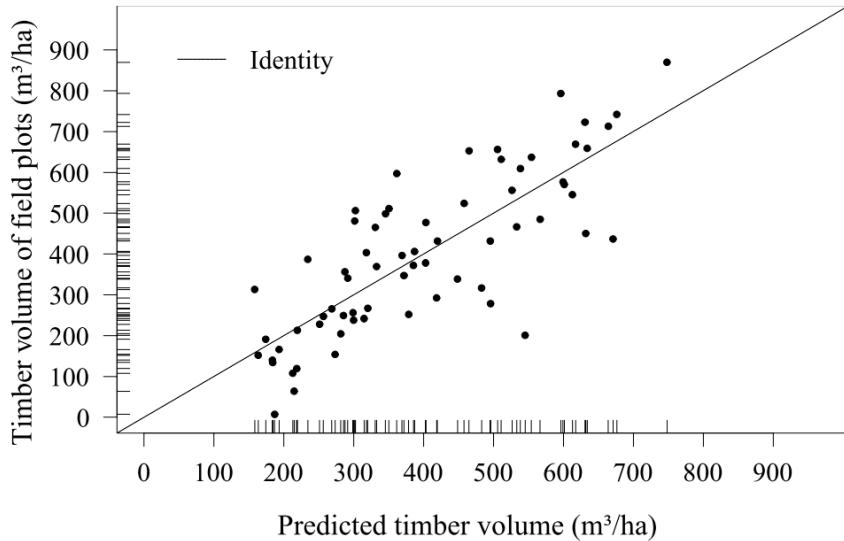


Figure 4.5: Model-predicted timber volume on plot level against the observed timber volume of the 67 field plots; the distribution of the predictions and observations are also indicated on the respective axis.

Table 4.4: Summary statistics of the regression model.

Predictors	AIC	R <sup>2</sup>	Adjusted R <sup>2</sup>	Mallow's C <sub>p</sub>	RMSE <sub>cv</sub>	Model Range
MEAN*, SD*, MAX*, Q75*	646.6	0.64	0.62	2.56	123.79	0 - 900 m <sup>3</sup> /ha

### 4.3.2 Assessment of Map Classification Accuracy

The classification accuracies described in Section 4.2.2 were estimated for nine possible constant class widths (100, 125, 150, 175, 200, 225, 250, 275 and 300 m<sup>3</sup>/ha), which were assumed to be of interest for representing the timber volume map. We applied the optimization model (Section 2.2.3) for each of those constant class widths in order to find a better classification scheme using the same number of corresponding classes, but locally adaptive class widths. The accuracies were then also estimated for the best-found classification schemes. Figure 4.6 provides a graphical summary of the results, giving the overall accuracies with their 95% confidence intervals, as well as the kappa coefficients. In all cases but one, the optimized locally-adapted class widths led to a higher overall classification accuracy compared to the corresponding constant class width. However, the 95% confidence intervals revealed that the overall accuracies using the constant class widths were, in all cases, not significantly different from their corresponding optimized alternatives (i.e., the confidence intervals are overlapping). In other words, even if the differences between the two accuracies seemed in several cases to be quite distinct, the true (but unknown) overall accuracies of the corresponding maps could actually be identical. However, with respect to the overlapping areas of the confidence intervals, the probability of acquiring a better overall accuracy by using the locally-adapted class widths was highest for smaller constant class widths (i.e., for larger numbers of classes). Interestingly, at the same time, a constant class width of 225 m<sup>3</sup>/ha (four classes) was also the best-found solution for locally-adapted class widths. In all cases but one, the kappa coefficients were higher under the application of the optimized class widths, even more distinct for larger numbers of classes (i.e., smaller constant class widths).

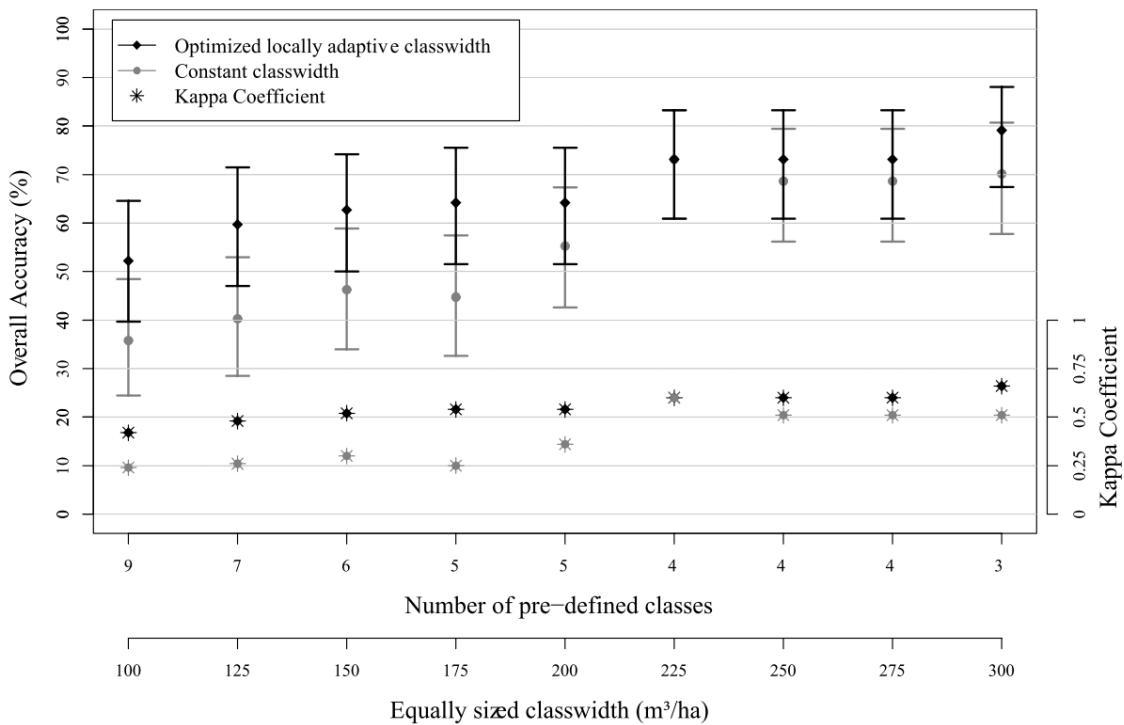


Figure 4.6: Overall accuracies with 95% confidence intervals and corresponding kappa coefficients for a pre-defined number of classes, calculated for constant and locally-adapted class widths.

We further compared the properties of the confusion matrices obtained by the use of constant and corresponding locally-adapted class widths, especially with respect to the estimated producer's

and user's accuracies. The producer's, as well as the user's accuracies under the optimized classification schemes were consistently higher than those of the corresponding constant class width approach. A phenomenon that, in several cases, occurred for the constant class width approach was that some of the classes did not include at least one prediction. This phenomenon especially concerned the classes comprising the upper scale region of the field data where the number of reference data was small. Consequently, the producer's accuracy for these classes was zero, whereas the user's accuracy is not even defined for this case (Section 4.2.2). The reason for this undesired property of the confusion matrix is likely to be that the fit of the regression model for the upper timber volume range was rather poor due to a limited number of field data or/and a saturation effect in the CHM (i.e., beyond a certain canopy height, different timber volumes cannot longer be discriminated by the regression model). The constant class width approach did, however, not account for this effect, whereas the optimization model proposed larger class widths for the upper-range classes, including a larger number of field data and predictions. An illustration of the findings described can be found in Table 4.5 and Table 4.6, using a constant class width of 200 m<sup>3</sup>/ha (five locally-adapted classes, respectively). The described locally-adaptive classification scheme provided a satisfactory trade-off between the number of classes and classification accuracies and was therefore exemplary used for the visualization of the timber volume map (Figure 4.7).

Table 4.5: Constant class width of 200 m<sup>3</sup>/ha (five timber volume classes); OAA ( $\pm$  95% confidence interval) (%): 55.22 (42.58, 67.4); kappa: 0.36; square sum of class width: 200'000.

Classes	Class Width	Producer's Accuracy	User's Accuracy	No. of References
(0,200)	200	60	85.71	10
(200, 400)	200	72	60	25
(400, 600)	200	40	40	20
(600, 800)	200	45.45	50	11
(800, 1000)	200	0	-	1

Table 4.6: Optimized class width for five timber volume classes; OAA ( $\pm$  95% confidence interval) (%): 64.18 (51.53, 75.53); kappa: 0.54; square sum of class width: 176'600. Class widths smaller than 200 m<sup>3</sup>/ha are indicated by \*.

Classes	Class Width	Producer's Accuracy	User's Accuracy	No. of References
(0,220)	220	76.92	90.91	13
(220, 330)	110*	61.54	50	13
(330, 450)	120*	57.14	53.33	14
(450, 660)	210	66.67	66.67	21
(660, 900)	240	50	75	6

Additionally, we investigated how often the locally-adaptive class widths were capable of satisfying Rules I and II of the optimization model (Section 4.2.2) more successfully than their constant class width counterparts. This was done by comparing the respective sum of squared class widths (second term of Equation 4.6), as well as the respective squared difference between the average number of plots per class and the actual number of reference data for each class (third term of Equation 4.6). It turned out that, particularly for the larger constant class widths (175, 200, 225, 250, 275 and 300 m<sup>3</sup>/ha), the locally-adapted approach worked very successfully for both

objectives: in six out of nine cases, each objective was better solved by the locally-adaptive approach, and in five out of nine cases, the optimization approach even succeeded in both objectives simultaneously.

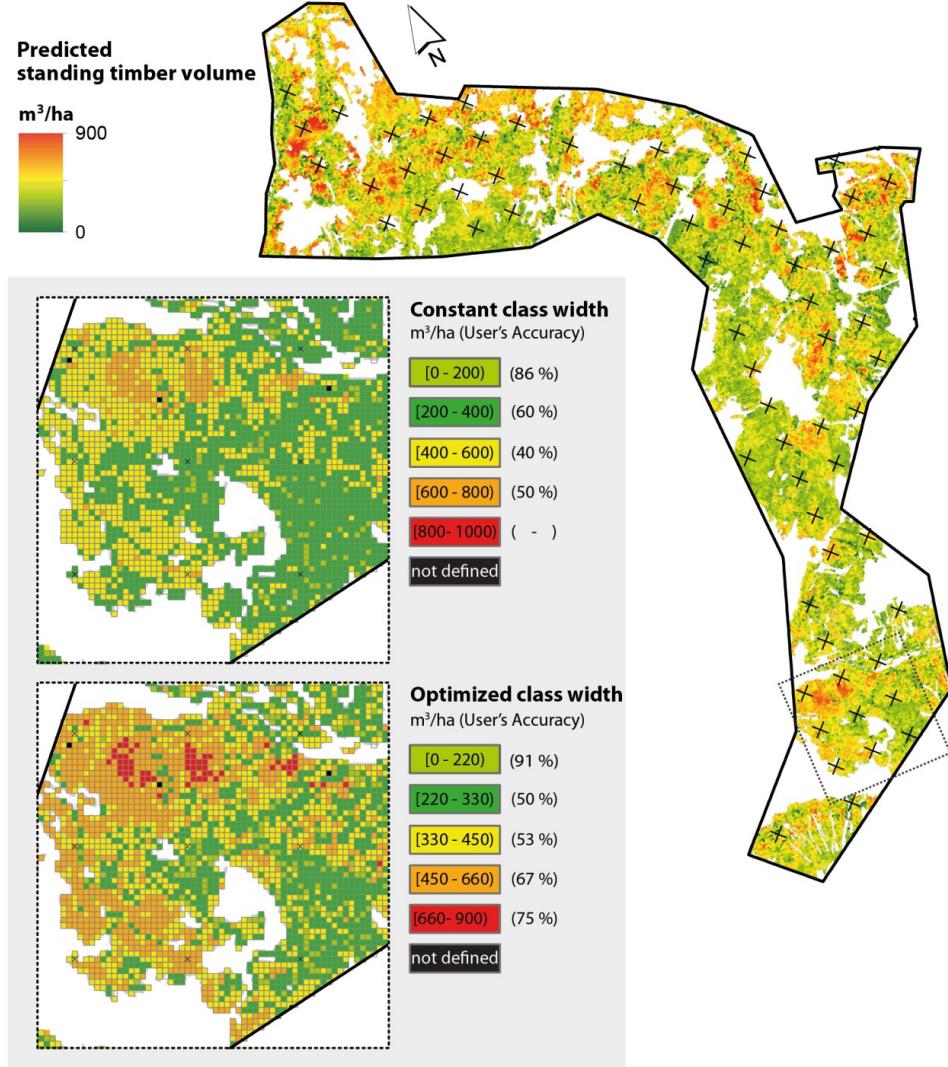


Figure 4.7: Volume map with model-predictions on a continuous scale for the entire study area covering 2000 hectares of forest with a spatial resolution of 25 meters; the subareas show the classified volume map using a constant class width of 200 m<sup>3</sup>/ha (upper) and locally-optimized class widths for five classes (lower).

#### 4.3.3 Calculation of the Timber Volume Map

Figure 4.7 shows the timber volume map for the entire study area estimated by the application of the regression model (Equation 4.8). An undesired property of the map was the occasional appearance of negative predictions, which are most likely caused by the negative signs of two of the regression coefficients (Equation 4.8). Additionally, the upper range of the predictions of the entire map (1002 m<sup>3</sup>/ha) turned out to exceed the upper valid model range (900 m<sup>3</sup>/ha). The values of these raster cells were consequently changed into 'Not Available' (NA). The map was additionally classified according to a constant class width of 200 m<sup>3</sup>/ha, as well as to the

corresponding optimized classification scheme presented in Section 4.3.2. As the classification schemes were also based on the valid model range, cells whose predictions did not cover any of the pre-defined classes were consequently marked as 'not defined' in the classified map. The number of these cases was, however, small compared to the total number of predicted cell values of the entire volume map (170 of 31,622 raster cells, i.e., 0.5%).

## 4.4 Discussion

### 4.4.1 Assessment of Map Accuracy

A core issue of this study was to use a class representation of continuous map predictions to estimate accuracy metrics of individual timber volume ranges (user's and producers accuracies). An evaluation for a large number of classification schemes revealed that the class accuracies can vary according to (i) which class width is used, (ii) how the class boundaries are chosen and (iii) a combination of both aspects. This variation in accuracy is not indicated by metrics such as the cross-validated RMSE<sub>cv</sub> of the prediction model, which can, however, be used to describe the overall prediction performance. When investigating the map for explicit timber volume ranges, purely relying on RMSE<sub>cv</sub> or R<sup>2</sup> metrics could, however, lead to over- or underestimation of the provided accuracy. Building the map according to classes and, consequently, estimating their individual accuracies is therefore considered a valuable step towards coping with the uncertainty of such maps.

### 4.4.2 Class Selection by Optimization Model

Regarding the choice of appropriate classes to represent the timber volume map, the proposed optimization model turned out to be of high value. While the overall accuracies between the optimized class width and constant class width approach did not differ significantly on a 5% significance level, the true benefit of the optimized class width approach were the properties of the confusion matrices thus obtained. The proposed classes revealed a more uniform distribution of references among the classes, ensuring that the estimated producer's and user's accuracies of each class were estimated by the highest possible number of references. The class width was thus successfully chosen smaller where a sufficient number of reference data were available, leading to a higher degree of detail and, with respect to an optimized adaptation of the class boundaries, without a loss of accuracy in those classes. In many of the evaluated classification schemes, the adaptation of the class boundaries led to higher producer's and user's accuracies compared to those obtained under constant class widths.

The optimization model has yet provided a method which allows for (i) finding an optimal classification scheme for a given number of classes and (ii) finding a number of classes, such that the class accuracies are acceptable for a user. If forest managers are interested in identifying stands with timber volumes greater than a selected minimum, perhaps for possible thinning treatments, one of the class boundaries can be fixed to that minimum volume.

While the maximization for the number of predictions and references included in the same class (Rule II, Section 4.2.2) gives a more realistic representation of the underlying regression model, one could argue that optimizing for this objective can lead to an 'overfitting' of the confusion matrix and, thus, to an increased generalization error of the accuracy metrics. This could not be investigated, since an appropriate method of bootstrapping and cross-validating for this kind of classification process has yet to be implemented. However, one could also slightly change the optimization model in order to optimize only for Rule I (i.e., to locally achieve smaller class widths and an even distribution of the reference data over the classes).

### 4.4.3 Regression Model

In the present study, the mean of the canopy height turned out to be the best predictor, providing an  $R^2$  of already 0.5. The fit was further improved to an adjusted  $R^2$  of 0.62 by extending the model using the predictor variables standard deviation, maximum height and the 75% quantile. This, however, came at the price of losing interpretability of the model (negative signs of some regression coefficients), as well as the occurrence of negative predictions in the timber volume map. The latter reassuringly happened only for 141 of 31,622 pixels (0.4%). To avoid negative predictions, we also computed the regression model using the log-response of the terrestrial timber volume and then calculated the map by back-transformation of the log-predictions into the original scale (Beauchamp & Olson, 1973). The log-response model, however, tended to predict unrealistically large timber volume values (up to 2400 m<sup>3</sup>/ha), exceeding the valid model range in 267 cases (0.16%) and, thus, not providing an improvement in prediction performance. The use of a log-response or non-linear regression model is also methodically sound in the design-based framework, with the restriction that only the variance under the external model assumption (Mandallaz, 2008b) can be calculated. This is because the design-based variance using the g-weight technique (Mandallaz et al., 2013; Mandallaz, 2013a) is only available within linear models on the original scale. Within the model-dependent framework, other alternatives to ensure non-negative predictions can be the use of nonlinear logistic regression models (McRoberts et al., 2013) or k-NN approaches. One could also consider the alternative use of external models (e.g., taken from literature) based on data that provide, for example, a larger model range or better coverage of a required timber volume range. It would be statistically sound to estimate the classification accuracies by applying an external model and consequently validating its classified predictions by available reference data of the study site. Using such models in the design-based framework of inventory would, however, again require calculating the variance under the external model approach.

### 4.4.4 Availability and Quality of Reference Data

In the present study, we used existing forest inventory data as reference and validation data, instead of acquiring these data in a special campaign. While this is in general both time and cost saving, the number of available reference data for calibrating the regression model and for validating the timber volume map was considerably limited, finally resulting in large confidence intervals of the accuracy metrics. However, in order to produce maps for operational forest management, probably over considerably large areas, one will realistically always depend on the use of existing inventory data due to limited financial resources. It should also be mentioned that the proposed methods assume the nominal coordinates of all field plots to be equal to the actual, true location of acquisition. This is in practice almost never the case, due to potential location errors, which can still be in the range of up to 10 meters, even under the use of GPS technique (Mauro et al., 2010; Steinmann et al., 2013). Although the remote sensing data can reveal positioning errors as well, they can be expected to be considerably smaller than the largest possible GPS location errors occasionally caused by dense vegetation and shielding of GPS reception on the ground. However, severe location errors should have appeared as outliers (or even leverage points) in the regression model. Since, in our case, neither the cross-validation nor an inspection of the regression model revealed any such outliers, the severity of location errors was assumed to be small. However, the  $R^2$  of the regression model is expected to be higher if the exact positions of the plot centers are known (Fuller, 2009), and the same can be assumed for the classification accuracies under the restriction that the locally-adaptive class width must be recomputed if the model is changed.

## 4.5 Conclusion

The methods proposed in this study provide better knowledge about the actual accuracy of a timber volume map. The classification of predictions into classes and the consequent computation of classification accuracy metrics improved the knowledge about the map accuracy, especially the accuracies of timber volume intervals. Additionally, considering the distribution of the reference data, as well as their corresponding predictions turned out to be key factors in choosing an appropriate classification scheme: the application of optimized locally-adaptive class widths ensured good statistical properties of the confusion matrices and also led to higher class accuracies and kappa coefficients compared to the approach of using constant class widths. The proposed methods, including the optimization of classification schemes, are thus not restricted to maps based on linear regression models, but can be applied to a larger class of prediction methods (e.g., k-NN estimation). Finally, the proposed design of the timber volume map was considered to be crucial for the reliability of the estimated accuracy metrics. Another advancement of the entire map set up is that the continuous map predictions can be directly used in the framework of design-based inventory methods. For example, several raster cells of the volume map could be merged, e.g., to 0.5 or one hectare, and the resulting larger cell could be used for true small area (synthetic) estimation (Mandallaz, 2013a). We expect this approach to also solve the problem of negative predictions. The design-based confidence intervals thus obtained for these estimates could then serve as a measure for local accuracy. The accuracy assessment of timber volume classes could also improve an automatic delineation of possible harvesting units by additionally considering the probability of each raster cell belonging to a certain class.

## Acknowledgements

We express our thanks to the reviewers for their comments and suggestions; to Professor Hans Rudolf Heinimann (Chair of Land Use Engineering, Swiss Federal Institute of Technology ETH Zurich) for his support; to the Forest Service of the canton of Grisons (Switzerland) for providing the forest inventory data; and to the Remote Sensing Laboratories (University of Zurich) for information on the LiDAR acquisition.



# Bibliography

- Akaike, H. (1992). Information theory and an extension of the maximum likelihood principle. In *Breakthroughs in statistics* (pp. 610–624). Springer.
- Akaike, H. (2011). *Akaike's Information Criterion*, (pp. 25–25). Springer Berlin Heidelberg: Berlin, Heidelberg.
- Baffetta, F., Corona, P., & Fattorini, L. (2010). Design-based diagnostics for k-nn estimators of forest resources this article is one of a selection of papers from extending forest inventory and monitoring over space and time. *Canadian Journal of Forest Research*, 41(1), 59–72.
- Baffetta, F., Fattorini, L., Franceschi, S., & Corona, P. (2009). Design-based approach to k-nearest neighbours technique for coupling field and remotely sensed data in forest surveys. *Remote Sensing of Environment*, 113(3), 463–475.
- Beauchamp, J. J. & Olson, J. S. (1973). Corrections for bias in regression estimates after logarithmic transformation. *Ecology*, 54(6), 1403–1407.
- Beaudoin, A., Bernier, P., Guindon, L., Villemaire, P., Guo, X., Stinson, G., Bergeron, T., Magnusson, S., & Hall, R. (2014). Mapping attributes of canada's forests at moderate resolution through k nn and modis imagery. *Canadian Journal of Forest Research*, 44(5), 521–532.
- Bitterlich, W. (1984). *The relascope idea. Relative measurements in forestry*. Commonwealth Agricultural Bureaux.
- Bohlin, J., Bohlin, I., Jonzén, J., & Nilsson, M. (2017). Mapping forest attributes using data from stereophotogrammetry of aerial images and field data from the national forest inventory. *SILVA FENNICA*, 51(2).
- Brassel, P. & Lischke, H. (2001). *Swiss National Forest Inventory: Methods and Models of the Second Assessment*. WSL Swiss Federal Research Institute, CH-8903 Birmensdorf. Technical report, ISBN 3-905620-99-5. URL <http://www.lfi.ch/publikationen/publ/methods/methods.pdf>.
- Breidenbach, J. (2015). **JoSAE**: Functions for some Unit-Level Small Area Estimators and their Variances. *R* package version 0.2.3.
- Breidenbach, J. & Astrup, R. (2012). Small area estimation of forest attributes in the norwegian national forest inventory. *European Journal of Forest Research*, 131(4), 1255–1267.
- Breidenbach, J., Kublin, E., McGaughey, R., Andersen, H.-E., & Reutebuch, S. E. (2008). Mixed-effects models for estimating stand volume by means of small footprint airborne laser scanner data. *Photogrammetric Journal of Finland*, 21(1), 4–15.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- Brosofske, K. D., Froese, R. E., Falkowski, M. J., & Banskota, A. (2014). A review of methods for mapping and prediction of inventory attributes for operational forest management. *Forest Science*, 60(4), 733–756.

- Bundesministerium für Ernährung, L. u. V. (2011). Aufnahmeanweisung für die dritte Bundeswaldinventur BWI3 (2011 - 2012).
- Carroll, R. J., Ruppert, D., Stefanski, L. A., & Crainiceanu, C. M. (2006). *Measurement error in nonlinear models: a modern perspective*. CRC press.
- Chirici, G., Corona, P., Marchetti, M., Mastromandi, A., Maselli, F., Bottai, L., & Travaglini, D. (2012). K-nn forest: a software for the non-parametric prediction and mapping of environmental variables by the k-nearest neighbors algorithm. *European Journal of Remote Sensing*, 45(1), 433–442.
- Clementel, F., Colle, G., Farruggia, C., Floris, A., Scrinzi, G., & Torresan, C. (2012). Estimating forest timber volume by means of "low-cost" lidar data. *Italian Journal of Remote Sensing/Rivista Italiana di Telerilevamento*, 44(1).
- Cochran, W. G. (2007). *Sampling techniques*. John Wiley & Sons.
- Congalton, R. G. & Green, K. (2008). *Assessing the accuracy of remotely sensed data: principles and practices*. CRC press.
- Cullmann, A. D. (2016). **maSAE: Mandallaz' Model-Assisted Small Area Estimators**. R package version 0.1-5.
- Deo, R. K., Froese, R. E., Falkowski, M. J., & Hudak, A. T. (2016). Optimizing variable radius plot size and lidar resolution to model standing volume in conifer forests. *Canadian Journal of Remote Sensing*, 42(5), 428–442.
- Dowle, M. & Srinivasan, A. (2017). **data.table: Extension of ‘data.frame’**. R package version 1.10.4-1.
- Draper, N. R. & Smith, H. (2014). *Applied regression analysis*. John Wiley & Sons.
- Fahrmeir, L., Kneib, T., Lang, S., & Marx, B. (2013). *Regression: Models, Methods and Applications*. Springer Science & Business Media.
- Fehrman, L., Lehtonen, A., Kleinn, C., & Tomppo, E. (2008). Comparison of linear and mixed-effect regression models and ak-nearest neighbour approach for estimation of single-tree biomass. *Canadian Journal of Forest Research*, 38(1), 1–9.
- Franco-Lopez, H., Ek, A. R., & Bauer, M. E. (2001). Estimation and mapping of forest stand density, volume, and cover type using the k-nearest neighbors method. *Remote sensing of environment*, 77(3), 251–274.
- Fuller, W. A. (2009). *Measurement error models*, volume 305. John Wiley & Sons.
- Gauer, J. & Aldinger, E. (2005). Waldökologische Naturräume Deutschlands-Wuchsgebiete. *Mitteilungen des Vereins für Forstliche Standortskunde und Forstpflanzenzüchtung*, 43, 281–288.
- Gregoire, T. & Valentine Harry, T. (2008). Sampling strategies for natural resources and the environment chapman & hall/crc, boca raton. *Fla. London*.
- Gregoire, T. G. & Valentine, H. T. (2007). *Sampling Strategies for Natural Resources and the Environment*. CRC Press.
- Gustafson, P. (2003). *Measurement error and misclassification in statistics and epidemiology: impacts and Bayesian adjustments*. CRC Press.

- Haara, A. & Kangas, A. (2012). Comparing k nearest neighbours methods and linear regression—is there reason to select one over the other? *Mathematical and Computational Forestry & Natural Resource Sciences*, 4(1), 50.
- Hill, A. (2013). Comparison of small area estimators in forest inventory using airborne laserscanning data. Master's thesis, ETH Zurich, University of Göttingen.
- Hill, A. (2017). *forestinventory*. GitHub repository.
- Hill, A., Breschan, J., & Mandallaz, D. (2014). Accuracy assessment of timber volume maps using forest inventory data and lidar canopy height models. *Forests*, 5(9), 2253–2275.
- Hill, A., Mandallaz, D., Buddenbaum, H., Stoffels, J., & Langshausen, J. (2017). Implementation of design-based small area estimations on forest district level in Rhineland-Palatinate by combining remote sensing data with data of the Third National German Inventory. Third International Workshop on Forest Inventory Statistics, Freiburg.
- Hoffmann, C. (1982). Die berechnung von tarifen für die waldinventur. *Forstwissenschaftliches Centralblatt*, 101(1), 24–36.
- Hollaus, M., Wagner, W., Maier, B., & Schadauer, K. (2007). Airborne laser scanning of forest stem volume in a mountainous environment. *Sensors*, 7(8), 1559–1577.
- Holmgren, J. (2004). Prediction of tree height, basal area and stem volume in forest stands using airborne laser scanning. *Scandinavian Journal of Forest Research*, 19(6), 543–553.
- Husmann, K., Rumpf, S., & Nagel, J. (2017). Biomass functions and nutrient contents of european beech, oak, sycamore maple and ash and their meaning for the biomass supply chain. *Journal of Cleaner Production*.
- Jakubowski, M. K., Guo, Q., & Kelly, M. (2013). Tradeoffs between lidar pulse density and forest measurement accuracy. *Remote Sensing of Environment*, 130, 245–253.
- Keller, M. (2011). Swiss national forest inventory. manual of the field survey 2004–2007. *Swiss Federal Research Institute WSL, Birmensdorf, CH*.
- Kirchhoefer, M., Schumacher, J., Adler, P., & Kändler, G. (2017). Considerations towards a novel approach for integrating angle-count sampling data in remote sensing based forest inventories. *Forests*, 8(7), 239.
- Kirkpatrick, S., Gelatt, C. D., Vecchi, M. P., et al. (1983). Optimization by simulated annealing. *science*, 220(4598), 671–680.
- Koch, B. (2010). Status and future of laser scanning, synthetic aperture radar and hyperspectral remote sensing data for forest biomass assessment. *ISPRS Journal of Photogrammetry and Remote Sensing*, 65(6), 581–590.
- Köhl, M., Magnussen, S. S., & Marchetti, M. (2006). *Sampling Methods, Remote Sensing and GIS Multiresource Forest Inventory*. Springer-Verlag Berlin Heidelberg.
- Kublin, E. (2003). Einheitliche beschreibung der schaftform–methoden und programme–bdatpro. *Forstwissenschaftliches Centralblatt*, 122(3), 183–200.
- Kublin, E., Breidenbach, J., & Kändler, G. (2013). A flexible stem taper and volume prediction method based on mixed-effects b-spline regression. *European journal of forest research*, 132(5–6), 983–997.

- Lamprecht, S., Hill, A., Stoffels, J., & Udelhoven, T. (2017). A machine learning method for co-registration and individual tree matching of forest inventory and airborne laser scanning data. *Remote Sensing*, 9(5).
- Latifi, H., Nothdurft, A., & Koch, B. (2010). Non-parametric prediction and mapping of standing timber volume and biomass in a temperate forest: application of multiple optical/lidar-derived predictors. *Forestry*, 83(4), 395–407.
- Latifi, H., Nothdurft, A., Straub, C., & Koch, B. (2012). Modelling stratified forest attributes using optical/lidar features in a central european landscape. *International Journal of Digital Earth*, 5(2), 106–132.
- Lefsky, M. A., Cohen, W., Acker, S., Parker, G. G., Spies, T., & Harding, D. (1999). Lidar remote sensing of the canopy structure and biophysical properties of douglas-fir western hemlock forests. *Remote sensing of environment*, 70(3), 339–361.
- Liaw, A. & Wiener, M. (2002). Classification and regression by randomforest. *R News*, 2(3), 18–22.
- Lumley, T. (2016). *survey: Analysis of Complex Survey Samples*. R package version 3.32.
- Maack, J., Lingenfelder, M., Weinacker, H., & Koch, B. (2016). Modelling the standing timber volume of baden-württemberg-a large-scale approach using a fusion of landsat, airborne lidar and national forest inventory data. *International Journal of Applied Earth Observation and Geoinformation*, 49, 107–116.
- Magnussen, S., Eggermont, P., & LaRiccia, V. N. (1999). Recovering tree heights from airborne laser scanner data. *Forest science*, 45(3), 407–422.
- Magnussen, S., Mandallaz, D., Breidenbach, J., Lanz, A., & Ginzler, C. (2014). National forest inventories in the service of small area estimation of stem volume. *Canadian Journal of Forest Research*, 44(9), 1079–1090.
- Magnussen, S., Tomppo, E., & McRoberts, R. E. (2010). A model-assisted k-nearest neighbour approach to remove extrapolation bias. *Scandinavian Journal of Forest Research*, 25(2), 174–184.
- Mallows, C. L. (2000). Some comments on cp. *Technometrics*, 42(1), 87–94.
- Mandallaz, D. (2008a). *Sampling Techniques for Forest Inventories*. CRC Press.
- Mandallaz, D. (2008b). *Sampling techniques for forest inventories*. CRC Press.
- Mandallaz, D. (2013a). Design-based properties of some small-area estimators in forest inventory with two-phase sampling. *Canadian Journal of Forest Research*, 43(5), 441–449.
- Mandallaz, D. (2013b). *Regression Estimators in Forest Inventories with Three-Phase Sampling and Two Multivariate Components of Auxiliary Information*. Technical report, Department of Environmental Systems Science, ETH Zurich.
- Mandallaz, D. (2013c). *Regression Estimators in Forest Inventories with Two-Phase Sampling and Partially Exhaustive Information with Applications to Small-Area Estimation*. Technical report, Department of Environmental Systems Science, ETH Zurich.
- Mandallaz, D. (2013d). A three-phase sampling extension of the generalized regression estimator with partially exhaustive information. *Canadian Journal of Forest Research*, 44(4), 383–388.

- Mandallaz, D. (2015). *Mathematical Details of Two-Phase/Two-Stage and Three-Phase/Two-Stage Regression Estimators in Forest Inventories. Design-based Monte Carlo Approach*. Technical report, Department of Environmental Systems Science, ETH Zurich.
- Mandallaz, D., Breschan, J., & Hill, A. (2013). New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based monte carlo approach with applications to small-area estimation. *Canadian Journal of Forest Research*, 43(11), 1023–1031.
- Mandallaz, D., Hill, A., & Massey, A. (2016). *Design-Based Properties of Some Small-Area Estimators in Forest Inventory with Two-Phase Sampling - Revised Version*. Technical report, Department of Environmental Systems Science, ETH Zurich.
- Massey, A. & Mandallaz, D. (2015a). Comparison of classical, kernel-based, and nearest neighbors regression estimators using the design-based monte carlo approach for two-phase forest inventories. *Canadian Journal of Forest Research*, 45(11), 1480–1488.
- Massey, A. & Mandallaz, D. (2015b). Design-based regression estimation of net change for forest inventories. *Canadian Journal of Forest Research*, 45(12), 1775–1784.
- Massey, A., Mandallaz, D., & Lanz, A. (2014). Integrating remote sensing and past inventory data under the new annual design of the swiss national forest inventory using three-phase design-based regression estimation. *Canadian Journal of Forest Research*, 44(10), 1177–1186.
- Massey, A. F. (2015). *Multiphase Estimation Procedures for Forest Inventories under the Design-Based Monte Carlo Approach*. PhD thesis, ETH Zurich.
- Mathworks (2017). Matlab version 9.2.0.538062 (r2017a).
- Mauro, F., Valbuena, R., Manzanera, J., & García-Abril, A. (2010). Influence of global navigation satellite system errors in positioning inventory plots for tree-height distribution studies this article is one of a selection of papers from extending forest inventory and monitoring over space and time. *Canadian journal of forest research*, 41(1), 11–23.
- McRoberts, R. E., Næsset, E., & Gobakken, T. (2013). Inference for lidar-assisted estimation of forest growing stock volume. *Remote Sensing of Environment*, 128, 268–275.
- McRoberts, R. E., Næsset, E., Gobakken, T., & Bollandsås, O. M. (2015). Indirect and direct estimation of forest biomass change using forest inventory and airborne laser scanning data. *Remote Sensing of Environment*, 164, 36–42.
- McRoberts, R. E., Tomppo, E. O., Finley, A. O., & Heikkinen, J. (2007). Estimating areal means and variances of forest attributes using the k-nearest neighbors technique and satellite imagery. *Remote Sensing of Environment*, 111(4), 466–480.
- McRoberts, R. E., Tomppo, E. O., & Næsset, E. (2010). Advances and emerging issues in national forest inventories. *Scandinavian Journal of Forest Research*, 25(4), 368–381.
- Næsset, E. (1997). Estimating timber volume of forest stands using airborne laser scanner data. *Remote Sensing of Environment*, 61(2), 246 – 253.
- Næsset, E. (2002). Predicting forest stand characteristics with airborne scanning laser using a practical two-stage procedure and field data. *Remote sensing of environment*, 80(1), 88–99.

- Næsset, E. (2014). Area-based inventory in norway - from innovation to an operational reality. In *Forest Applications of Airborne Laser Scanning - Concepts and Case Studies* chapter 11, (pp. 216–240). Springer.
- Nink, S., Hill, J., Buddenbaum, H., Stoffels, J., Sachtleber, T., & Langshausen, J. (2015). Assessing the suitability of future multi-and hyperspectral satellite systems for mapping the spatial distribution of norway spruce timber volume. *Remote Sensing*, 7(9), 12009–12040.
- Nothdurft, A., Saborowski, J., & Breidenbach, J. (2009). Spatial prediction of forest stand variables. *European Journal of Forest Research*, 128(3), 241–251.
- Ott, E., Frehner, M., Frey, H.-U., & LÄ¶scher, P. (1997). *GebirgsnadelwÄlder: Ein praxisorientierter Leitfaden fÄr eine standortgerechte Waldbehandlung*. Haupt..
- Pinheiro, J. & Bates, D. (2000). *Mixed-Effects Models in S and S-PLUS*. Springer-Verlag.
- Pirlot, M. (1996). General local search methods. *European journal of operational research*, 92(3), 493–511.
- R Core Team (2017). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rayward-Smith, V. J., Osman, C., Reeves, C. R., & Smith, G. D. (1996). *Modern heuristic search methods*. John Wiley.
- Richards, J. A. & Richards, J. (1999). *Remote sensing digital image analysis. An introduction*, volume 3. Springer.
- Saborowski, J., Marx, A., Nagel, J., & Böckmann, T. (2010). Double sampling for stratification in periodic inventories-infinite population approach. *Forest ecology and management*, 260(10), 1886–1895.
- SAS Institute Inc. (2015). *SAS/STAT Software, Version 9.4*. Cary, NC.
- Schreuder, H. T., Gregoire, T. G., & Wood, G. B. (1993). *Sampling Methods for Multiresource Forest Inventory*. John Wiley & Sons.
- Srivastava, A. K., Srivastava, V. K., & Ullah, A. (1995). The coefficient of determination and its adjusted version in linear regression models. *Econometric reviews*, 14(2), 229–240.
- Steinmann, K., Mandallaz, D., Ginzler, C., & Lanz, A. (2013). Small area estimations of proportion of forest and timber volume combining lidar data and stereo aerial images with terrestrial data. *Scandinavian journal of forest research*, 28(4), 373–385.
- Stoffels, J., Hill, J., Sachtleber, T., Mader, S., Buddenbaum, H., Stern, O., Langshausen, J., Dietz, J., & Ontrup, G. (2015). Satellite-based derivation of high-resolution forest information layers for operational forest management. *Forests*, 6(6), 1982–2013.
- Straub, C., Dees, M., Weinacker, H., & Koch, B. (2009). Using airborne laser scanner data and cir orthophotos to estimate the stem volume of forest stands. *Photogrammetrie-Fernerkundung-Geoinformation*, 2009(3), 277–287.
- Tomppo, E. (2006). The finnish multi-source national forest inventory—small area estimation and map production. *Forest inventory—methodology and applications*. Springer, Dordrecht, NL, (pp. 195–224).

- Tonolli, S., Dalponte, M., Vescovo, L., Rodeghiero, M., Bruzzone, L., & Gianelle, D. (2011). Mapping and modeling forest tree volume using forest inventory and airborne laser scanning. *European Journal of Forest Research*, 130(4), 569–577.
- Van Aardt, J. A., Wynne, R. H., & Scrivani, J. A. (2008). Lidar-based mapping of forest volume and biomass by taxonomic group using structurally homogenous segments. *Photogrammetric Engineering & Remote Sensing*, 74(8), 1033–1044.
- von Lüpke, N. (2013). *Approaches for the Optimisation of Double Sampling for Stratification in Repeated Forest Inventories*. PhD thesis, University of Göttingen.
- von Lüpke, N., Hansen, J., & Saborowski, J. (2012). A three-phase sampling procedure for continuous forest inventory with partial re-measurement and updating of terrestrial sample plots. *European Journal of Forest Research*, 131(6), 1979–1990.
- von Lüpke, N. & Saborowski, J. (2014). Combining double sampling for stratification and cluster sampling to a three-level sampling design for continuous forest inventories. *European journal of forest research*, 133(1), 89–100.
- von Thünen-Institut (2014). Dritte Bundeswaldinventur 2012. Accessed: 2017-02-03.
- White, J. C., Coops, N. C., Wulder, M. A., Vastaranta, M., Hilker, T., & Tompalski, P. (2016). Remote sensing technologies for enhancing forest inventories: A review. *Canadian Journal of Remote Sensing*, 42(5), 619–641.
- Wickham, H. (2009). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- Zianis, D., Muukkonen, P., Mäkipää, R., Mencuccini, M., et al. (2005). Biomass and stem volume equations for tree species in Europe. *Silva Fennica*, Monographs 4.

# Curriculum Vitae

## Personal Details

---

Name	Andreas Christian Hill
Date of birth	15.01.1986
Place of birth	Trier (Germany)

## Education

---

02/2014–03/2018	Ph.D. Student ETH Zurich, Department of Environmental Systems Science, Chair of Forest Engineering
04/2015–05/2017	Diploma of Advanced Studies ETH ETH Zurich, Applied Statistics
04/2010–08/2013	Master of Science (M.Sc.) University of Göttingen, Forest Sciences and Forest Ecology with study focus on Forest Ecosystem Analysis and Information Processing
10/2006–03/2010	Bachelor of Science (B.Sc.) University of Göttingen, Forest Sciences and Forest Ecology

## Professional Experience

---

Since 10/2013	Researcher Department of Environmental System Science, Chair of Forest Engineering, ETH Zurich
10/2011–04/2012	Research Assistant Department of Environmental System Science, Chair of Forest Engineering, ETH Zurich
09/2007–03/2011	Research Assistant Department of Forest Inventory and Remote Sensing, Faculty of Forest Sciences, University of Göttingen

## Publications in Conference Proceedings and Workshops

---

- [1] **Hill, A.**, Mandallaz, D., Buddenbaum, H., Stoffels, J., Langshausen, J. (2017): Implementation of design-based small area estimations on forest district level in Rhineland-Palatinate by combining remote sensing data with data of the Third German National Forest Inventory. In *3rd International Workshop on Forest Inventory Statistics*. Freiburg, Baden-Württemberg Germany.
- [2] **Hill, A.**, Stoffels, J., Langshausen, J. (2016): . In *CARISMA-workshop on large-scale mapping and estimation of forest resources*. Ås, Norway
- [3] **Hill, A.**, Breschan, J. (2014): Automatic Design of Efficient Harvesting Units using Remote Sensing and Field Data. In *24th IUFRO World Congress*. Salt Lake City, Utah, USA
- [4] Breschan, J., **Hill, A.** (2014): Validation of timber volume maps derived from remote sensing data. In *24th IUFRO World Congress*. Salt Lake City, Utah, USA

## Publications in Scientific Journals

---

- [1] Lamprecht, S., **Hill, A.**, Stoffels, J., Udelhoven, T.(2017): A Machine Learning Method for Co-Registration and Individual Tree Matching of Forest Inventory and Airborne Laser Scanning Data. *Remote Sensing*, 9 (5). doi: 10.3390/rs9050505
- [2] **Hill, A.**, Breschan, J., Mandallaz, D. (2014): Accuracy Assessment of Timber Volume Maps using Forest Inventory Data and LiDAR Canopy Height Models. *Forests*, 5 (9). 2253-2275. doi: 10.3390/f5092253
- [3] Mandallaz, D., Breschan, J., **Hill, A.** (2013): New Regression Estimators in Forest Inventories with Two-Phase Sampling and Partially Exhaustive Information: a Design-Based Monte Carlo Approach with Applications to Small-Area Estimation. *Canadian Journal of Forest Research*, 43 (11). 1023-1031. doi: 10.1139/cjfr-2013-0181

## Other Publications

---

- [1] Mandallaz, D., **Hill, A.**, Massey, A. (2016): Design-based properties of some small-area estimators in forest inventory with two-phase sampling - revised version. *Technical Report*, Department of Environmental Systems Science, ETH Zurich. doi: 10.3929/ethz-a-010579388
- [2] **Hill, A.**, Massey, A., Mandallaz D. (2016): forestinventory: Design-Based Global and Small-Area Estimations for Multiphase Forest Inventories. R package version 0.1.0 *CRAN Repository* url: <https://CRAN.R-project.org/package=forestinventory>

