**ARTICLE**

# 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

Alexander Massey, Daniel Mandallaz, and Adrian Lanz

**Abstract:** In 2009, the Swiss National Forest Inventory (NFI) turned from a periodic into an annual measurement design in which only one-ninth of the overall sample of permanent plots is measured every year. The reduction in sample size due to the implementation of the annual design results in an unacceptably large increase in variance when using the standard simple random sampling estimator. Thus, a flexible estimation procedure using two- and three-phase regression estimators is presented with a special focus on utilizing updating techniques to account for disturbances and growth and is applied to the second and third Swiss NFIs. The first phase consists of a dense sample of systematically distributed plots on a 500 m × 500 m grid for which auxiliary variables are obtained through the interpretation of aerial photographs. The second phase is an eightfold looser subgrid with terrestrial plot data collected from the past inventory, and the third and final phase consists of the three most recent annual subgrids with the current state of the target variable (stem volume). The proposed three-phase estimators reduce the increase in variance from 294% to 145% compared with the estimator based on the full periodic sample while remaining unbiased.

*Key words:* national forest inventories, three-phase regression, annual (paneled) inventory design, Swiss National Forest Inventory, remote sensing and previous measurement.

**Résumé :** L'inventaire Forestier National Suisse a changé en 2009 d'un système périodique (inventaire complet chaque 10 ans) à un relevé continu annuel par placettes permanentes sur un neuvième de l'échantillon total. La taille réduite de l'échantillon annuel conduit à une erreur inacceptablement grande de l'estimateur basé sur l'échantillonnage aléatoire simple. Pour remédier à ce problème nous proposons des estimateurs par régression utilisant un plan de sondage à deux ou trois phase et très flexibles. La variable cible considérée dans cet article est le volume sur pied. La première phase repose sur un réseau dense de photographies aériennes (500 m × 500 m), dont l'interprétation permet d'obtenir des variables explicatives sur la canopée (en particulier concernant les interventions), la deuxième phase consiste en placettes terrestres sur un sous-réseau (un huitième) de la première phase, dont les valeurs actualisées servent de variable explicative très performante, la troisième et dernière phase consiste des placettes avec les trois plus récentes données annuelles. La technique proposée permet d'obtenir des estimateurs sans biais et de réduire l'augmentation de la variance due au changement de système de 294 à 145 %.

*Mots-clés :* Inventaire forestier, estimation par régression avec plan d'échantillonnage à trois phases, Inventaire Forestier National Suisse avec relevé annuel, télédétection, actualisation de mesures antérieures.

## 1. Introduction

The fourth Swiss National Forest Inventory (NFI) is the first Swiss NFI to implement an annual measurement design in which one-ninth of its permanent remeasured terrestrial plots are assessed every year. Annual designs have the advantage of providing estimates continuously rather than periodically but have the major drawback of a dramatically reduced sample size if the simple random sampling (SRS) estimator is used for any given year. The loss in precision due to the reduction in sample size can be mitigated by grouping years together. For example, grouping three years of data, which seems reasonable for the Swiss NFI as this is approximately how long it took to conduct each of the previous three inventories, decreases the sample reduction from one-ninth to one-third but is still expected to inflate the variance by a factor of three on both the national and regional levels. Furthermore, grouping many years together deteriorates the concept of a coherent time point associated with the point estimate.

A better approach than grouping is to enhance the estimates with auxiliary data obtained from remote sensing data (i.e., satellite images, LiDAR, stereographic photos, etc.) or with information from past inventories. A plethora of such methods have been explored in the literature. In the context of incorporating past inventory data, Eskelson et al. (2009) and Johnson et al. (2003) compared various moving averages and weighted moving averages with the goal of achieving higher estimation precision by combining annual samples from multiple years. Another popular approach involves updating plots that were not in the annual sample using the current year's realization as a reference set. Such methods in the model-dependent case include combining inventory and remote sensing data using plot-level growth models (McRoberts 2001), spatiotemporal models (Houillier and Pierrat 1992), imputation techniques on the tree and plot levels (McRoberts 2001; Van Deusen 1997), tree-level growth models (Lessard et al. 2001), enhancing estimates over successive measurements using a Kalman filter (Dixon and Howitt 1979), and mixed estimation

**A. Massey and A. Lanz.** Swiss Federal Research Institute WSL, CH 8903 Birmensdorf, Switzerland.
**D. Mandallaz.** Chair of Land Use Engineering, ETH Zurich, CH 8092 Zurich, Switzerland.
**Corresponding author:** Alexander Massey (e-mail: alexander.massey@wsl.ch).

([Van Deusen 2002](#)). Further model-dependent methods, including regression updating, inverse weighting of new and old plots, and regression combined with growth-model projections, have been used to compare periodic and annual designs for both continuous forest inventory and sampling with partial replacement (SPR) ([Scott et al. 1999](#)). Literature concerning design-based methods for enhancing estimation under an annual design is currently scarce.

Here we present estimators in the model-assisted design-based framework specifically in the context of both two- and three-phase regression estimations under annual designs with no partial replacement. The theory draws heavily from the contributions of [Särndal et al. (2003)](#) and adaptations to the infinite population or Monte Carlo approach developed by [Mandallaz (2008)](#). While some authors have delved into similar model-assisted approaches using remote sensing in the two-stage estimator (see [Gregoire et al. 2011](#)), the three-phase setup in the specific context of regression within strata (see [Lüpke et al. 2012](#)), and two-phase sampling with partially exhaustive information ([Mandallaz et al. 2013](#)), this article presents three-phase regression estimation using the desirable g-weight technique ([Mandallaz 2013b](#)) and demonstrates its flexibility with data taken from permanent plots at the terrestrial level and remote sensing data derived from manually interpreted aerial photographs on which tree and forest stand data are observed. The proposed estimation procedure is flexible in the sense that it can easily be used to implement a variety of known variance reduction tools such as poststratification, regression estimation, and regression within stratification, as well as incorporate auxiliary information from multiple sources.

The objective here is to demonstrate how to integrate remote sensing data and previous inventory under an annual design in the Swiss NFI using regression estimation (REG). The main estimator of interest is a three-phase regression estimator (REG,3p) in which the first phase uses information gathered from remote sensing, the second phase is based on the previous plot measurements of the entire terrestrial sample, and the third phase is the current annual subsample for any given year or grouping of consecutive years. The classical two-phase regression estimator (REG,2p) will also be presented so that we can assess the efficiency of this scheme. Proofs are omitted but a basic overview of the estimators, as well as their associated asymptotic variance estimators, will be introduced and presented in such a way as to develop intuition and facilitate their implementation.

## 2. The sampling design

We have a well-defined population $P$ of trees $i \in 1, 2, …, N$ in forest $F$. For every tree, we have response variable $Y_i$. The parameter of interest is the spatial mean $\bar{Y} = \frac{1}{\lambda(F)} \sum_{i=1}^{N} Y_i$, where $\lambda(F)$ is the surface area of $F$. In this case, $Y_i$ is individual tree (stem) volume and $\bar{Y}$ is growing stock volume per unit area. The first phase is a large sample $s_1$ of $n_1$ points $x \in F$ distributed uniformly and assumed independently in forest $F$. Auxiliary information is taken in an arbitrary predefined region around each of these points that is ideally highly correlated to the response variable. It does not matter if it is qualitative, quantitative, or a combination of both. The auxiliary information is contained in row vector $\mathbf{Z}^t(x) = (\mathbf{Z}^{(1)t}(x), \mathbf{Z}^{(2)t}(x))$, where $\mathbf{Z}^{(1)t}(x)$ is the first-phase component known for a large sample size of points $x \in F$. $\mathbf{Z}^{(2)t}(x)$ is the second-phase component known only for points $x \in s_2$, where $s_2 \subset s_1$. The second-phase sample size is denoted $n_2$. The third-phase sample $s_3 \subset s_2$ draws $n_3$ points and corresponds directly to the observed ground truth measured by the field team for the desired year(s) under the annual design. Uniform and independent random sampling is used on an infinite population of points in a bounded domain of a plane for the selection of $s_1 \in F$. Simple random sampling without replacement is used in the selection of subsamples $s_2 \subset s_1$ and $s_3 \subset s_2$. Although most national forest inventories typically use systematic sampling schemes, the point and variance estimators can

still be considered acceptable (i.e., approximately unbiased) as long as the range of the spatial correlation of the target variable, primarily of the residuals, is small compared with the size of the area to be estimated, which is the case here (see chapter 7 of [Mandallaz 2008](#)).

For every $x \in s_3$, trees are selected from $P$ with inclusion probabilities $\pi_i$ equal to the area of an inclusion circle (possibly adjusted for the forest boundary) around the $i$th tree divided by $\lambda(F)$. For a given point $x$, the indicator variable for the $i$th tree is defined as follows:

$$(1) \qquad I_i(x): = \begin{cases} 1 \text{ if } i \in s_3(x) \\ 0 \text{ if } i \notin s_3(x) \end{cases}$$

$Y_i$ is evaluated for all trees $i$ such that $I_i(x) = 1$. Using the indicator variable and the inclusion probability, we can define the local density estimate.

$$(2) \qquad Y(x) := \frac{1}{\lambda(F)} \sum_{i=1}^{N} \frac{I_i(x)Y_i}{\pi_i}$$

Note that the local density estimate $Y(x)$ is assumed error-free given $x$ and the random selection $x$ is uniform in $F$. Thus, we are using the Monte Carlo approach in the design-based setup in which the relations $\mathbb{E}_x(Y(x)) = \frac{1}{\lambda(F)} \int_F Y(x) dx = \frac{1}{\lambda(F)} \sum_{i=1}^{N} Y_i = \bar{Y}$ hold for both concentric circles and variable radius (angle count) sampling of trees. This implies that $\frac{1}{n_3} \sum_{x \in s_3} Y(x)$ is an unbiased estimator of our parameter of interest $\bar{Y}$. However, this estimator does not utilize any information from the auxiliary vector $\mathbf{Z}^t(x)$, which is the motivation for using the model-assisted approach presented in this paper.

## 3. Model-assisted estimation

### 3.1. External models

We will start by presenting the difference estimator for three-phase sampling. It is highly general, intuitively simple, and leads naturally into the presentation of regression estimators. The difference estimator in three-phase sampling under the infinite population (Monte Carlo) approach is defined by [Mandallaz (2014)](#) based on ideas originating from [Särndal et al. (2003)](#) as

$$(3) \qquad \hat{Y}_{\text{DIFF,3p}} = \frac{1}{n_1} \sum_{x \in s_1} \hat{Y}_{M_1}(x) + \frac{1}{n_2} \sum_{x \in s_2} (\hat{Y}_{M_2}(x) - \hat{Y}_{M_1}(x)) + \frac{1}{n_3} \sum_{x \in s_3} (Y(x) - \hat{Y}_{M_2}(x))$$

where $\hat{Y}_{M_1}(x)$ and $\hat{Y}_{M_2}(x)$ are model predictions available for the first- and second-phase sampling points, respectively, that are error-free given $x$. We do not specify the type of model (e.g., linear, nonlinear, kNN, etc.) from which the predictions $\hat{Y}_{M_1}(x)$ and $\hat{Y}_{M_2}(x)$ are produced or that the models are unbiased in the model-dependent sense. We assume that the models are external (i.e., not fitted with the current sample of inventory data) so that the theoretical variance takes the form ([Mandallaz 2013b](#))

$$(4) \qquad \mathbb{V}(\hat{Y}_{\text{DIFF,3p}}) = \frac{1}{n_1} \mathbb{V}(Y(x)) + \left(1 - \frac{n_2}{n_1}\right) \frac{1}{n_2} \mathbb{V}(R_{M_1}(x)) + \left(1 - \frac{n_3}{n_2}\right) \frac{1}{n_3} \mathbb{V}(R_{M_2}(x))$$

with $R_{M_1}(x)) = Y(x) - \hat{Y}_{M_1}$ and $R_{M_2}(x)) = Y(x) - \hat{Y}_{M_2}$. It is easily checked using the total law of expectation that the overall expectation of the

point estimate is the parameter of interest (i.e., $\mathbb{E}_{1,2,3}\hat{Y}_{\text{DIFF,3p}} = \overline{Y}$) and $\hat{Y}_{\text{DIFF,3p}}$ is exactly design-unbiased. If the models are fit from the current sample of inventory data (i.e., dependent on the stochastic process of the sampling design), they are called internal, and the estimators given in section 3.2 can be used. If one disregards that a model is fit internally and applies eqs. 3 and 4, i.e., under an implicit external model assumption, then the point estimator is still approximately design-unbiased, but the variance will likely be slightly underestimated because the model fit is dependent on the sampling process. For linear models, the external model assumption is quite reasonable, especially given large sample size. The adequacy of eq. 4 is less clear for other models such as kNN.

## 3.2. Internal models

To incorporate the effect of the sampled data on the fit of the internal model and the overall precision of the estimators, one must first specify which type of model to be considered. Two nested linear models will be considered here: one to be fitted only with the reduced component, $\mathbf{Z}^{(1)t}(x)$, of the auxiliary vector available for all first-phase points in $s_1$; and a second model fitted with the full auxiliary vector $\mathbf{Z}^t(x) = (\mathbf{Z}^{(1)t}(x), \mathbf{Z}^{(2)t}(x))$ available for all second-phase points in $s_2$. The former will be referred to as $M_{\text{reduced}}$ ("reduced" refers to the reduced number of covariates in the nested framework) and the latter as $M_{\text{full}}$ ("full" indicates that all covariates in the auxiliary vector are used).

### 3.2.1. The reduced model

The model $M_{\text{reduced}}$

$$(5) \qquad Y(x) = \mathbf{Z}^{(1)t}(x)\boldsymbol{\alpha} + R^{(1)}(x)$$

has a vector of theoretical regression coefficients $\boldsymbol{\alpha}$ and theoretical residual term $R^{(1)}(x)$. In complete analogy to classical least squares in a finite universe, $\boldsymbol{\alpha}$ minimizes the integral of theoretical residuals squared over all possible points $x$, which is infinite in the Monte Carlo approach. Thus, $\boldsymbol{\alpha}$ minimizes $\int_F R^{(1)}(x)^2 dx = \int_F (Y(x) - \mathbf{Z}^{(1)t}(x)\boldsymbol{\alpha})^2 dx$ and leads to the normal equation $(\int_F \mathbf{Z}^{(1)}(x)\mathbf{Z}^{(1)t}(x)dx)\boldsymbol{\alpha} = \int_F Y(x)\mathbf{Z}^{(1)}(x)dx$. Of course, we have incomplete information to solve for $\boldsymbol{\alpha}$ directly for the infinite population of points $x$ because $Y(x)$ is observed only given the third-phase sample, whereas $\mathbf{Z}^{(1)t}$ is available only for $x \in s_1$. Thus, we estimate $\boldsymbol{\alpha}$ by sample copies of the normal equation leading to

$$(6) \qquad \begin{aligned} \hat{\boldsymbol{\alpha}} &= \left(\frac{1}{n_3}\sum_{x \in s_3}\mathbf{Z}^{(1)}(x)\mathbf{Z}^{(1)t}(x)\right)^{-1}\frac{1}{n_3}\sum_{x \in s_3}Y(x)\mathbf{Z}^{(1)}(x) \\ &= \left(\frac{1}{n_3}\mathbf{Z}^{(1)t}_{s_3}\mathbf{Z}^{(1)}_{s_3}\right)^{-1}\frac{1}{n_3}\sum_{x \in s_3}Y(x)\mathbf{Z}^{(1)}(x) \end{aligned}$$

where $\mathbf{Z}^{(1)}_{s_3}$ is the design matrix of explanatory variables constrained to the third-phase points, $x \in s_3$. In other words, $\hat{\boldsymbol{\alpha}}$ is simply the vector of regression coefficients obtained when fitting $M_{\text{reduced}}$ using ordinary least squares (OLS) to the third-phase sample.

### 3.2.2. The full model

In the same way as $M_{\text{reduced}}$, $M_{\text{full}}$ is defined

$$(7) \qquad \begin{aligned} Y(x) &= \mathbf{Z}^t(x)\boldsymbol{\beta} + R(x) \\ &= \mathbf{Z}^{(1)t}(x)\boldsymbol{\beta}^{(1)} + \mathbf{Z}^{(2)t}(x)\boldsymbol{\beta}^{(2)} + R(x) \end{aligned}$$

where $\boldsymbol{\beta}^t = (\boldsymbol{\beta}^{(1)t}, \boldsymbol{\beta}^{(2)t})$ is the vector of theoretical regression coefficients. Note that $\boldsymbol{\beta}^{(1)} = \boldsymbol{\alpha}$ only if $\mathbf{Z}^{(1)}(x)$ and $\mathbf{Z}^{(2)}(x)$ are orthogonal (i.e., independent), which is almost never the case in practice. The corresponding normal equation is $(\int_F \mathbf{Z}(x)\mathbf{Z}^t(x))\boldsymbol{\beta} = \int_F Y(x)\mathbf{Z}(x)$ and $\boldsymbol{\beta}$ is estimated as follows:

$$(8) \qquad \begin{aligned} \hat{\boldsymbol{\beta}} &= \left(\frac{1}{n_3}\sum_{x \in s_3}\mathbf{Z}(x)\mathbf{Z}^t(x)\right)^{-1}\frac{1}{n_3}\sum_{x \in s_3}Y(x)\mathbf{Z}(x) \\ &= \left(\frac{1}{n_3}\mathbf{Z}^t_{s_3}\mathbf{Z}_{s_3}\right)^{-1}\frac{1}{n_3}\sum_{x \in s_3}Y(x)\mathbf{Z}(x) \end{aligned}$$

where $\mathbf{Z}_{s_3}$ is the design matrix containing all explanatory variables present in the auxiliary vector but restricted to $x \in s_3$. $\hat{\boldsymbol{\beta}}$ is simply the vector of regression coefficients obtained when fitting $M_{\text{full}}$ using ordinary least squares (OLS) to the third-phase sample.

### Zero mean property

Note that both $M_{\text{reduced}}$ and $M_{\text{full}}$ possess the same properties as in classical OLS regression such as zero mean empirical residuals (i.e., $\frac{1}{n_3}\sum_{x \in s_3}\hat{R}(x) = 0$ and $\frac{1}{n_3}\sum_{x \in s_3}\hat{R}^{(1)}(x) = 0$, where $\hat{R}^{(1)}(x) = Y(x) - \mathbf{Z}^{(1)t}(x)\hat{\boldsymbol{\alpha}}$ and $\hat{R}(x) = Y(x) - \mathbf{Z}^t(x)\hat{\boldsymbol{\beta}}$), and the orthogonality of the true residuals to the model space (i.e., $\mathbf{Z}^{(1)t}(x)\boldsymbol{\alpha}$ and $R^{(1)}(x)$ are orthogonal by construction and likewise for $M_{\text{full}}$). They are unbiased in the model-dependent sense but it should be emphasized that we will be using them in the design-based framework where we never need to make the assumption that the models hold true in the model-dependent sense.

### 3.3. The three-phase regression estimator

With the nested internal models $M_{\text{reduced}}$ and $M_{\text{full}}$ defined in the last section, the three-phase regression estimator is

$$(9) \qquad \begin{aligned} \hat{Y}_{\text{REG,3p}} = \frac{1}{n_1}\sum_{x \in s_1}\hat{Y}^{(1)}(x) &+ \frac{1}{n_2}\sum_{x \in s_2}(\hat{Y}(x) - \hat{Y}^{(1)}(x)) \\ &+ \frac{1}{n_3}\sum_{x \in s_3}(Y(x) - \hat{Y}(x)) \end{aligned}$$

where $\hat{Y}^{(1)}(x) = \mathbf{Z}^{(1)t}(x)\hat{\boldsymbol{\alpha}}$ and $\hat{Y}(x) = \mathbf{Z}^t(x)\hat{\boldsymbol{\beta}}$ are the predictions for the target variable $Y(x)$ at point $x$ based on the reduced and full sets of explanatory variables respectively.

The estimator can be rewritten in the more useful form

$$(10) \qquad \begin{aligned} \hat{Y}_{\text{REG,3p}} &= \left(\hat{\bar{\mathbf{Z}}}^{(1)}_{s_1} - \hat{\bar{\mathbf{Z}}}^{(1)}_{s_2}\right)^t\hat{\boldsymbol{\alpha}} + \left(\hat{\bar{\mathbf{Z}}}_{s_2} - \hat{\bar{\mathbf{Z}}}_{s_3}\right)^t\hat{\boldsymbol{\beta}} + \frac{1}{n_3}\sum_{x \in s_3}Y(x) \\ &= \left(\hat{\bar{\mathbf{Z}}}^{(1)}_{s_1} - \hat{\bar{\mathbf{Z}}}^{(1)}_{s_2}\right)^t\hat{\boldsymbol{\alpha}} + \hat{\bar{\mathbf{Z}}}^t_{s_2}\hat{\boldsymbol{\beta}} \end{aligned}$$

where $\hat{\bar{\mathbf{Z}}}^{(1)}_{s_1} = \frac{1}{n_1}\sum_{x \in s_1}\mathbf{Z}^{(1)}(x)$, $\hat{\bar{\mathbf{Z}}}^{(1)}_{s_2} = \frac{1}{n_2}\sum_{x \in s_2}\mathbf{Z}^{(1)}(x)$, $\hat{\bar{\mathbf{Z}}}_{s_2} = \frac{1}{n_2}\sum_{x \in s_2}\mathbf{Z}(x)$, and $\hat{\bar{\mathbf{Z}}}_{s_3} = \frac{1}{n_3}\sum_{x \in s_3}\mathbf{Z}(x)$.

As the true residuals are design-based uncorrelated with the predictions by construction in least squares, we get $\mathbb{V}(Y(x)) = \mathbb{V}(\hat{Y}^{(1)}(x)) + \mathbb{V}(R^{(1)}(x))$ for $M_{\text{reduced}}$ and $\mathbb{V}(Y(x)) = \mathbb{V}(\hat{Y}(x)) + \mathbb{V}(R(x))$ for $M_{\text{full}}$. Furthermore, if we accept the external model assumption (i.e., considering the internal models to be external), we can plug these decompositions into eq. 4 and the theoretical variance simplifies to

$$(11) \qquad \mathbb{V}(\hat{Y}_{\text{REG,3p}}) = \frac{1}{n_1}\mathbb{V}(\hat{Y}^{(1)}(x)) + \frac{1}{n_2}\mathbb{V}(R^{(1)}(x)) + \left(1 - \frac{n_3}{n_2}\right)\frac{1}{n_3}\mathbb{V}(R(x))$$

The sample could be used to estimate the terms in eq. 11.

However, here we are forgoing the external model assumption to develop a variance estimator that takes the effects of an internally fitted model into account (i.e., the sampling variability on $\hat{\alpha}$ and $\hat{\beta}$).

$\hat{Y}_{\text{REG,3p}}$ is still asymptotically design-unbiased and the variance estimator is written, for simplicity, in terms of g-weights (Mandallaz 2013a, 2014):

$$(12) \qquad g^{(1)}(x) := \hat{\bar{\mathbf{Z}}}_{s_1}^{(1)t}(\mathbf{A}_{s_2}^{(1)})^{-1}\mathbf{Z}^{(1)}(x)$$

$$(13) \qquad g(x) := \hat{\bar{\mathbf{Z}}}_{s_2}^{t}\mathbf{A}_{s_3}^{-1}\mathbf{Z}(x)$$

where

$$\mathbf{A}_{s_2}^{(1)} = \frac{1}{n_2}\sum_{x \in s_2}\mathbf{Z}^{(1)}(x)\mathbf{Z}^{(1)}(x)^t = \frac{1}{n_2}\mathbf{Z}_{s_2}^{(1)t}\mathbf{Z}_{s_2}^{(1)}$$

$$\mathbf{A}_{s_3} = \frac{1}{n_3}\sum_{x \in s_3}\mathbf{Z}(x)\mathbf{Z}(x)^t = \frac{1}{n_3}\mathbf{Z}_{s_3}^{t}\mathbf{Z}_{s_3}$$

with $\mathbf{Z}_{s_2}^{(1)}$ and $\mathbf{Z}_{s_3}$ being the design matrices of $M_{\text{reduced}}$ and $M_{\text{full}}$ defined over the entire second and third phases, respectively. The REG possesses the following useful calibration properties:

$$(14) \qquad \begin{aligned}\frac{1}{n_2}\sum_{x \in s_2}g^{(1)}(x)\mathbf{Z}^{(1)}(x) &= \hat{\bar{\mathbf{Z}}}_{s_1}^{(1)} \\ \frac{1}{n_2}\sum_{x \in s_2}g^{(1)}(x) &= 1\end{aligned}$$

$$(15) \qquad \begin{aligned}\frac{1}{n_3}\sum_{x \in s_3}g(x)\mathbf{Z}(x) &= \hat{\bar{\mathbf{Z}}}_{s_2} \\ \frac{1}{n_3}\sum_{x \in s_3}g(x) &= 1\end{aligned}$$

The $g^{(1)}(x)$-weighted mean of the auxiliary variables $\mathbf{Z}^{(1)}(x)$ over second-phase sample $s_2$ is equal to the mean over first-phase sample $s_1$ (i.e., $\hat{\bar{\mathbf{Z}}}_{s_1}^{(1)}$). Similarly, the $g(x)$-weighted mean of $\mathbf{Z}(x)$ over $s_3$ is equal to the mean of $\mathbf{Z}(x)$ over $s_2$ (i.e., $\hat{\bar{\mathbf{Z}}}_{s_2}$). The intuition follows that it is desirable to apply these weights to $Y(x)$ provided that $Y(x)$ is well correlated with the auxiliary vectors $\mathbf{Z}^{(1)}(x)$ and $\mathbf{Z}(x)$. The g-weights tend to 1 asymptotically and have a mean equal to 1 for $g^{(1)}(x)$ and $g(x)$. Furthermore, it is important to note that better statistical properties arise when the g-weights are used explicitly in the variance estimator than compared with a variance estimator derived by the external model assumption (Mandallaz 2013*a*, 2013*b*; Mandallaz et al. 2013; for a concrete example in the two-phase case, see p. 84 in Mandallaz 2008).

There is no closed analytical formula for the theoretical variance under the internal model; however, an asymptotically consistent variance estimator derived by the g-weight technique is defined (Mandallaz 2013*b*) as

$$(16) \qquad \begin{aligned}\widehat{\mathbb{V}}(\hat{Y}_{\text{REG,3p}}) &= \frac{1}{n_1}\frac{\sum_{x \in s_1}(\hat{Y}^{(1)}(x) - \hat{\bar{Y}}_{s_1}^{(1)})^2}{n_1 - 1} \\ &+ \frac{1}{n_2}\frac{1}{n_3}\sum_{x \in s_3}(g^{(1)}(x)\hat{R}^{(1)}(x))^2 + \frac{1}{n_3^2}\left(1 - \frac{n_3}{n_2}\right)\sum_{x \in s_3}(g(x)\hat{R}(x))^2\end{aligned}$$

where $\hat{Y}^{(1)}(x) = \mathbf{Z}^{(1)t}(x)\hat{\boldsymbol{\alpha}}_2$ and $\hat{\bar{Y}}_{s_1}^{(1)} = \frac{1}{n_1}\sum_{x \in s_1}\hat{Y}^{(1)}(x)$. Notice that eq. 16 is asymptotically equivalent to the theoretical variance under the external model assumption given in eq. 11 because of the zero mean property of the empirical residuals and the fact that the g-weights tend to 1 asymptotically.

### 3.4. The two-phase regression estimator

We also want to quickly present the two-phase regression estimator (REG,2p) because of its usefulness in our study to assess the relative contribution of the individual phases. In the following,

the notational convention is to present the REG,2p in the notation of the REG,3p from section 3.3, but leaving out the first phase (i.e., $n_1 = n_2$). There is only the model $M_{\text{full}}$ (note that ignoring the second phase would be a similar process but with $M_{\text{reduced}}$ and changing the subscripts from 2 to 1). The REG,2p simplifies to (see chapter 6 in Mandallaz 2008)

$$(17) \qquad \begin{aligned}\hat{Y}_{\text{REG,2p}} &= \frac{1}{n_2}\sum_{x \in s_2}\hat{Y}(x) + \frac{1}{n_3}\sum_{x \in s_3}(Y(x) - \hat{Y}(x)) \\ &= \frac{1}{n_2}\sum_{x \in s_2}\mathbf{Z}^t(x)\hat{\boldsymbol{\beta}} + \frac{1}{n_3}\sum_{x \in s_3}(Y(x) - \mathbf{Z}^t(x)\hat{\boldsymbol{\beta}}) \\ &= (\hat{\bar{\mathbf{Z}}}_2 - \hat{\bar{\mathbf{Z}}}_3)^t\hat{\boldsymbol{\beta}} + \frac{1}{n_3}\sum_{x \in s_3}Y(x) \\ &= \hat{\bar{\mathbf{Z}}}_2^t\hat{\boldsymbol{\beta}}\end{aligned}$$

There is only one set of g-weights that is defined the same as in eq. 13. The approximately design-unbiased g-weight variance estimator is

$$(18) \qquad \begin{aligned}\widehat{\mathbb{V}}(\hat{Y}_{\text{REG,2p}}) &= \frac{1}{n_2}\frac{1}{n_3 - 1}\sum_{x \in s_3}(Y(x) - \bar{Y}_{s_3})^2 \\ &+ \frac{1}{n_3^2}\left(1 - \frac{n_3}{n_2}\right)\sum_{x \in s_3}(g(x)\hat{R}(x))^2\end{aligned}$$

The REG,2p has some special cases that are worth mentioning:

- double sampling for stratification — only categorical variables are present in the model (i.e., we have an ANOVA model); the strata weights are calculated based on the first phase;
- double sampling for regression — only continuous variables are present in the model; multiple linear regression is possible and a single intercept is possible;
- double sampling for regression within strata — a combination of categorical and continuous variables are present in the model; when an interaction term is not present between categorical and continuous variables, the regression lines within strata will have the same slope but different intercepts; when the interaction is present, the slope is allowed to be vary within strata, which is defined by the categorical variable; and the strata weights will be calculated based on the first phase.

These special cases tend to lose their concrete interpretations when extended to three-phase regression estimation. We emphasize that one should not restrict oneself to utilizing the REG for special cases, but merely that they are easily implemented. The goal intuitively in model selection is to minimize the residuals as they typically play the most prominent part in the variance estimation.

## 4. Case study

### 4.1. Study area

The Swiss NFI consists of some 160 000 systematically distributed and permanent sampling points for aerial photo interpretation and a subgrid of some 20 000 systematically distributed and permanently installed sampling points for terrestrial data collection (Lanz et al. (2010)). Approximately one-third of the points are located in forest and have been retained for this study. The restriction to a forest subsample is not straightforward and deserves some explanation. Firstly, no map of forest land exists that conforms to the NFI forest definition. Aerial photo interpretation and terrestrial inventory provide independent forest versus nonforest classifications. The classifications are highly correlated; less than 1% of the sampling points have a differing classification. In the standard estimation procedures of the NFI, the terrestrial classifi-

**Table 1.** Sample size for each phase by region.

| Region | Forested area (ha) | $n_1$ | $n_2$ | $n_3$ |
|---|---|---|---|---|
| Jura | 201 185 | 8 064 | 996 | 331 |
| Plateau | 230 014 | 9 294 | 1164 | 379 |
| Pre-Alps | 218 596 | 8 864 | 1092 | 347 |
| Alps | 370 842 | 15 082 | 1944 | 602 |
| South Alps | 151 570 | 6 466 | 821 | 240 |
| Switzerland | 1 172 207 | 47 770 | 6017 | 1899 |

**Note:** These sample sizes are based on a "working forest definition" derived solely from aerial stereophotography. In practice, some plots under this working definition are not actually measured terrestrially because they are deemed to be nonforest by the field crew. $n_3$ is adjusted to exclude plots for which the target variable was not measured and we assume that these responses are missing at random.

cation is the ground truth and the classification derived from aerial photo interpretation is considered as auxiliary information in a two-phase estimation procedure (double sampling for post-stratification).

For the purpose of this study, we retained first-phase sampling points that were classified as forest during aerial photo interpretation. From these points, the second phase is defined on a one-eighth subgrid. The third phase ideally should be based on a further subgrid based on three consecutive years taken from the annual subgrid implemented in current fourth Swiss NFI. However, in practice, 368 of the plots have no ground truth measured by the field crew because of a differing forest–nonforest classification. These plots are considered miss at random and were removed from the third phase post-factum. The sample sizes are given in Table 1.

The second and third Swiss NFIs are two-phase forest inventories in which permanent terrestrial plots are defined by systematically distributed points on interpenetrating square grids. The second Swiss NFI was conducted from 1993 to 1995, whereas the third Swiss NFI was conducted from 2004 to 2007. In both inventories, the first phase draws a sample from a regular 500 m × 500 m grid for auxiliary variables based on manually interpreted digital aerial photographs that include roughly 160 000 forest and nonforest plots. The second phase consists of roughly 20 000 plots taken on a $\sqrt{2}$ km by $\sqrt{2}$ km subgrid for target variables. Approximately one-third of these are forest plots and are evaluated based on terrestrial measurements. The survey partitions Switzerland (CH) into five production regions: Jura (JU), Plateau (SP), Pre-Alps (PA), Alps (AL), and South Alps (SA).

### 4.2. Ground data

Trees on each plot are re-identified and remeasured every inventory using calipers according to two concentric circles measuring 200 m² for trees with diameter at breast height (DBH) between 12 cm and 36 cm and 500 m² for trees with DBH greater than 36 cm. There are approximately 11 trees per plot on average. An approximation of timber volume (over bark bole volume) is determined by using DBH (one-way yield table stratified by species) on all eligible trees. For details, we refer to the reader to Brassel and Lischke (2001). For the purposes of this article, it is sufficient to assume that the local density estimate is a known error-free quantity.

The fourth Swiss NFI, which began in 2009, differs from the second and third NFIs due to the implementation of a new annual design in which one-ninth of the plots are measured every year using interpenetrating grids to allow for more up-to-date estimates to be available between inventory cycles. The drawback of this annual strategy is that these estimates are subject to a dramatically reduced sample size, which provides the motivation for exploring more efficient model-assisted methods.

### 4.3. Digital aerial stereophotography

As previously mentioned, the first phase of the Swiss NFI consists of a large sample of plots taken from a 500 m × 500 m grid using digital aerial stereophotography. The main purpose of these photographs is to make a forest–nonforest decision so that field crews do not have to visit nonforest plots. In the second and third NFIs, other continuous landscape variables such as canopy height information were also assessed at these plots. The basic sampling unit consists of a 50 m × 50 m square interpretation area containing 25 equally spaced raster points arranged in a 5 × 5 design. The analogue true color photos were flown at a scale of 1 : 30 000 and scanned at a 14 μm resolution. Once digitized, the photos had an aerial ground resolution of 0.42 m and a RMS error of <1 m after triangulation.

Each lattice point was assigned one thematic cover class by a photo interpreter using a three-dimensional stereo softcopy station (SOCET SET 5.0, BAE Systems). Cover classes included tree vegetation (>3 m), shrub and herb vegetation (<3 m), soil and sand and gravel, rock, nonnatural surfaces, and open water. Canopy height information was calculated for each lattice point by taking the difference between the photogrammetically measured surface elevation by the interpreter and a bilinearly interpolated 25 m spaced terrain model provided by the Swiss Federal Office of Topography. For cases in which the forest border was predefined, a forest boundary line was digitized on screen.

The forest–nonforest classification was made based on the continuous landscape described in such a way as to mirror the NFI final forest definition that will be made terrestrially by the field crew. In the Swiss NFI, the main requirements are that the tree canopy height, excluding burned, cut, damaged, afforested, or regeneration forests, is greater than 3 m, there is a minimum tree canopy cover of 20% within the forest boundary line, and there is a minimum stand width of 25 m (for full mathematical details, see Keller's section on p. 51 in Brassel and Lischke 2001).

### 4.4. Sampling frame

The sampling frame is selected using the forest–nonforest decision from the aerial stereophotographs associated with the third NFI as a working definition for forest boundary. A small percentage of plots (about 0.9%) were designated "no decision" because of the difficulty of interpretation due to external factors such as weather. Only plots designated by the aerial photos as "forest" were included in the sampling frame. A forest–nonforest definition taken on the ground by the field crew is also possible if a nonforest stratum is allowed (note that misclassification of forested plots in the nonforest stratum is allowable), but for the purposes of evaluating the efficiency of the models in the forest area, it is more appropriate to limit the frame to the forested area.

### 4.5. Target variable

The target variable of interest is the growing stock volume (VOL) in cubic metres per hectare (m³·ha⁻¹) of living trees excluding shrub forest in the accessible forest. There is a national threshold of 12 cm DBH for trees to be included in the estimates.

### 4.6. Auxiliary variables

The considered first-phase auxiliary variables included the mean canopy height of the raster points identified in forest (MCH), the variance of these points (VARCH), the proportion of these raster points lying on a coniferous tree (CONPROP), the proportion over all raster points in the forest area (INFORESTPROP), the final stratification variable defined by disturbance class (defined in the next section) and region (REGION:DIST), and all possible interactions between the aforementioned variables. Five quartiles of canopy heights were also considered.

The main variable of interest in the second phase is the plot's previous measurement for growing stock volume (PREV_VOLUME) obtained in the second NFI. PREV_VOLUME is usually one of the

strongest predictors of its current state assuming that the time to remeasurement is not unreasonably long. In theory, its predictive accuracy can be improved by accounting for two types of forest changes — loss due to mortality and disturbances due to natural or unnatural causes— and growth, including both physical tree growth and in-growth (i.e., trees surpassing the selection threshold). In the design-based setup, it is assumed that the predictors are known and error-free given $x$ so we can directly input predictions derived from any external tree or plot-level model into the proposed REGs as auxiliary variables. Here we will account for loss using a stratification variable (DIST) based on remote-sensing data acquired in the first phase. Although any model that provides good predictions can be used, for this case, study growth will be accounted for in each plot by aggregating predictions derived from tree-level linear growth models applied to previous inventory data. This variable, denoted PREV_VOLUME_UPDATED, does not account for volume caused by in-growth, but this is thought to be acceptable as in-growth is not expected to account for much of the overall plot volume.

### 4.6.1. Accounting for disturbance

The stratification variable, DIST, is intended to account for potential loss of growing stock at each plot due to natural and human causes. DIST is derived using canopy height information from two stereophotos taken from two different time frames. The first photo was taken just before the second NFI (flight years occurred in 1987–1993), and the second was the most recent taken before the next remeasurement in the third NFI (flight years occurred in 1998–2005). The key variable used to aid in disturbance prediction is the mean canopy height ratio (MCHR), defined as the mean canopy height of the raster points in the latter photo divided by the mean at the same raster points from the former photo. Unreasonable values for the raster heights (e.g., below 0 or above 55) are considered missing. The idea behind the MCHR is that if a disturbance occurred (i.e., harvesting, wind throw, etc.) in the window between the two photos, then this ratio is likely to be small.

The MCHR is prone to considerable noise that may arise from the preciseness of the locations of the raster points from photo to photo and the time of year that the photo was taken. Because no linear model was found in which MCHR was statistically significant as a continuous explanatory variable to predict current plot volume, it was decided to break it into classes for use as a stratification variable. Five classes were selected: $(-\infty, 0.4]$, $(0.4, 0.8]$, $(0.8, 1.2]$, $(1.2, 1.8]$, and $(1.8, \infty)$.

The harvesting intensities vary by region in Switzerland, so the effectiveness of these stratification classes were assessed visually. We employ the intuition behind "regression within strata" in which we have a categorical variable in the first and second phases and a continuous variable (PREV_VOLUME) in the second phase. Conceptually, the effectiveness of regression within strata comes when the slope of the regression line varies between strata. The slope of the regression lines within strata are clearly displayed in Fig. 1. As expected, the larger slopes correspond with larger values of MCHR. The correlations between past and current measurement are generally strong, indicating that the stratum is suitable to fit a regression line in. When the slopes between strata were too similar, the strata were collapsed, as was also the case if

a stratum contained too few samples corresponding to the third phase (recall that coefficients of the regression must be fit to the final phase where the response is known). Table 2 contains the third-phase sample size within strata for plots designated as forest plots in both the previous and current inventories as defined by the aerial photos.

The final disturbance class stratification (DIST) selected was $(-\infty, 0.8]$, $(0.8, 1.2]$, and $(12, \infty)$ within Jura, $(-\infty, 0.8]$, $(0.8, 1.2]$, and $(12, \infty)$ within Plateau, $(-\infty, 0.8]$ and $(08, \infty)$ within the Pre-Alps, and $(-\infty, 0.8]$ and $(08, \infty)$ within the Alps; the South Alps was treated as its own stratum without any disturbance classes (see Table 2). The final stratification by region (denoted REGION:DIST) also includes a stratum for all new forest plots that were not defined as forest in the second NFI by the aerial photography decision.
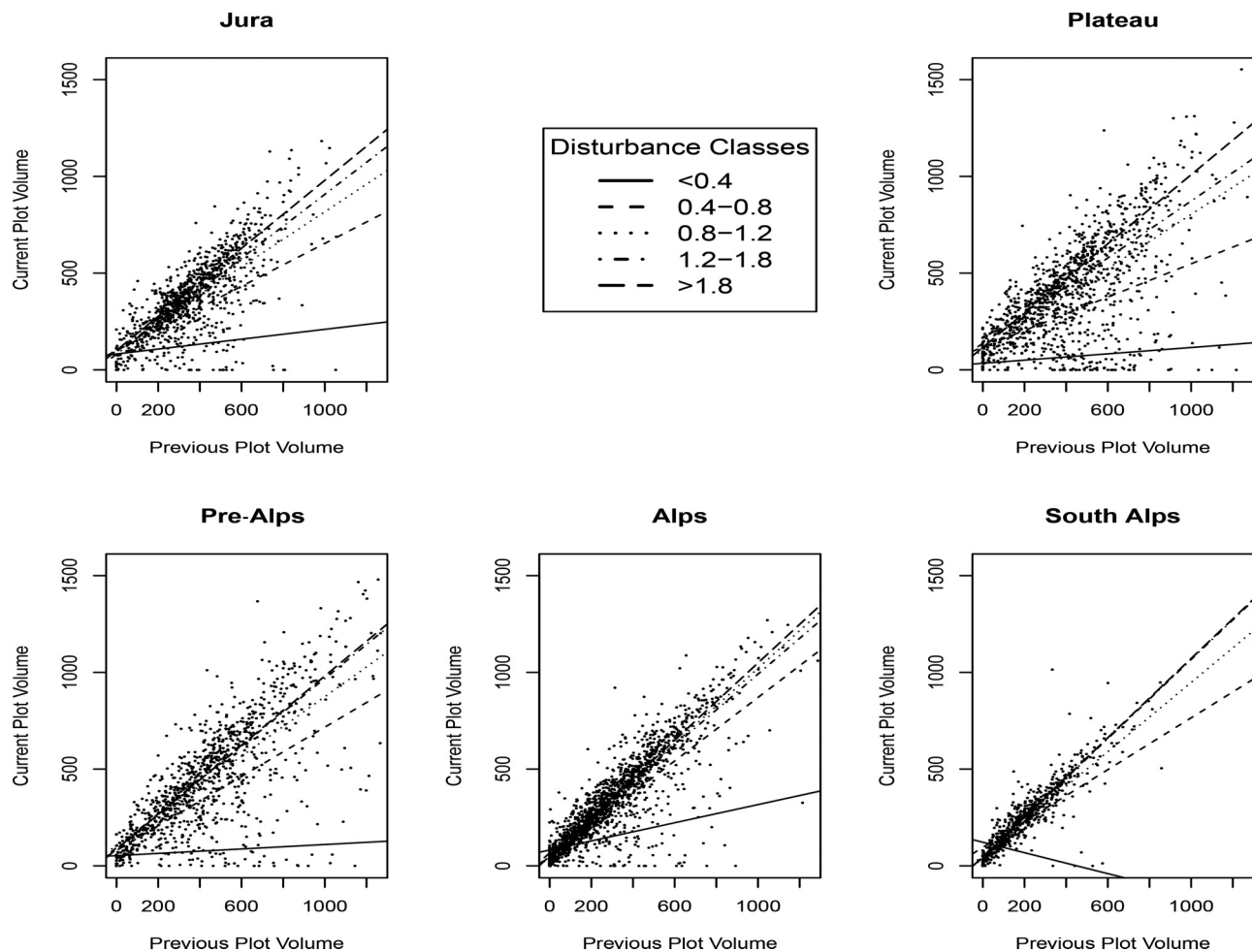
### 4.6.2. Accounting for growth on the tree level

Here we explain the derivation of PREV_VOLUME_UPDATED by examining two linear tree-level growth models with variables that were selected with a predictive criteria such as the Akaike information criterion (AIC). The models expand on the idea of fully utilizing all of the previous measurements of the Swiss NFI data (note that this time, we are on the tree level, not the plot level). A logically strong predictor of the future DBH would be the past DBH combined with the previous DBH increment obtained during the previous remeasurement cycle (i.e., trees must have been measured twice before to acquire the previous increment). Of course, this is not possible for plots that entered the Swiss NFI sample for the first time in the previous inventory. Thus, two separate linear models are proposed to account for these two cases. DBH is used as the target response variable because it can be directly input into the one-way yield table (refer to section 4.2).

The models were trained on all available terrestrial Swiss NFI data. After assembling a comprehensive list of plausible predictive variables (e.g., past DBH, time between measurements, categorical species–region groupings, etc.) accompanied with the recommended variance stabilizing transformations (including a log-transformed response variable), a backwards selection algorithm comparing AIC was implemented. AIC has good statistical properties in variable selection because it protects against overfitting. However, given the very large sample sizes of the training set (there were 55 248 trees for which the previous DBH growth increment exists and 120 297 trees for which only last known DBH existed), there were many variables that were strongly statistically significant but not necessarily relevant. Thus, the effect on the response of the maximum value of each remaining explanatory variable was considered. In the presence of a logged response variable $y$, the interpretation of the coefficients is multiplicative rather than additive because $\log(\hat{y}) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \ldots + \hat{\beta}_p x_p \Leftrightarrow \hat{y} = \exp \hat{\beta}_0 \exp \hat{\beta}_1 x_1 \ldots \exp \hat{\beta}_p x_p$. So the potential effect on $y$ of the maximal value of any given explanatory variable $x_q$ is $|1 - \exp \hat{\beta}_q \max x_q|$. To demonstrate the interpretation, if $|1 - \exp \hat{\beta}_q \max x_q| = 0.0034$, then given that all the other explanatory variables stay the same, the maximal effect of $x_q$ on $y$ is only 0.34%, implying weak predictive relevance in the model. Any variable with a maximal effect less than 7% was dropped. Note that the same idea can be applied to the effect of the median and minimum values of $x_q$. The following models were selected:

$$(19) \quad \log(\text{DBH}) = \gamma_0 + \gamma_1 \log(\text{PREV}_{\text{DBH}}) + \gamma_2 \log(\text{PREV}_{\text{DBH}})^2 + \gamma_3 \log(\text{PREV}_{\text{INC}} + 10) + \gamma_4 \text{VEGUNTIL} + \gamma_5 \log(\text{BASFPH} + 1) + \gamma_{6:9} \text{SCHICHT} + \gamma_{10} \log(\text{Z25}) + \gamma_{11:39} \text{TARNR} + \epsilon$$

$$(20) \quad \log(\text{DBH}) = \eta_0 + \eta_1 \log(\text{PREV}_{\text{DBH}}) + \eta_2 \log(\text{PREV}_{\text{DBH}})^2 + \eta_3 \log(\text{DDOM}) + \eta_4 \text{VEGUNTIL} + \eta_5 \log(\text{BASFPH} + 1) + \eta_{6:9} \text{SCHICHT} + \eta_{10} \log(\text{Z25}) + \eta_{11:39} \text{TARNR:VEGUNTIL} + \eta_{40} \log(\text{GWL3}) + \eta_{41:69} \text{TARNR} + \epsilon$$

**Fig. 1.** Regression lines for individual disturbance classes based on mean canopy height ratio (MCHR) by region with plot volumes in m³·ha⁻¹. The points represent all known plot volumes for the second Swiss National Forest Inventory (NFI) (previous plot volume) and the third Swiss NFI (current plot volume). Each regression line depicted is fit only on a subset of these points, which belong to a particular disturbance class definable for the entire first phase.



**Table 2.** Third-phase sample sizes in forest within disturbance classes.

| Disturbance class | Jura | Plateau | Pre-Alps | Alps | South Alps | Switzerland |
|---|---|---|---|---|---|---|
| <0.4 | 2a | 18d | 11g | 26i | 0k | 57 |
| 0.4–0.8 | 33a | 44d | 37g | 89i | 19k | 222 |
| 0.8–1.2 | 122b | 134e | 143h | 169j | 59k | 627 |
| 1.2–1.8 | 124c | 14f | 115h | 198j | 87k | 673 |
| >1.8 | 52c | 37f | 47h | 112j | 67k | 315 |

**Note:** a, b, c, …, k represent final stratification membership after collapsing disturbance classes with similar regression line slopes and (or) low third-phase sample size.

where DBH = current observed DBH, $PREV_{DBH}$ = (centered) diameter at breast height in the previous inventory, $PREV_{INC}$ = previous observed change in DBH standardized by the number of vegetative periods, VEGUNTIL = number of vegetative periods, BASFPH = basal area of standing living trees (in m²·ha⁻¹), SCHICHT = a categorical description of the story (i.e., upper, middle, lower, no story, story, undistinguishable), Z25 = altitude of above sea level, TARNR = a categorical variable defining region and species, GWL3 = a measure of site quality in terms of maximum mean annual biomass increment (in kg·ha⁻¹·year⁻¹) (Keller 1978), and DDOM = an index describing the mean diameter of the 100 thickest trees per hectare. DDOM is calculated on the plot as the mean DBH of the largest trees. Trees over 36 cm have a weight of 20 and

trees with DBH between 12 cm and 36 cm have a weight of 50. The thickest trees are selected such that the sum of their weights is greater than or equal to 100. More technical definitions of these variables can be found in Brassel and Lischke (2001).

Equation 19 has an $R^2$ = 0.9815 and is applied when the $PREV_{INC}$ is available, whereas eq. 20 has $R^2$ = 0.9676. The $R^2$ is high because of the strong influence of predicting the future DBH knowing the past one. Both models were fitted using data collected in all four Swiss NFIs excluding information on the response variable collected during the third NFI (i.e., the growth update is externally fit). PREV_VOLUME_UPDATED is derived by applying the growth models to all individual trees present on the plot in the previous measurement and summing their predicted volumes together.

### 4.7. Model selection

The $M_{reduced}$ was selected based on the first-phase remote-sensing variables from the third Swiss NFI using a forward variable selection procedure based on AIC and then an additional stepwise selection procedure with the predictors of the resulting model. The procedure was used to individually fit linear models on all regions, as well as Switzerland as a whole, using VOL as the response variable. The resulting models selected all contained the variables MCH, CONPROP, INFORESTPROP, and REGION:DIST but differed slightly in which sets of interaction terms were included.

The following model was chosen to represent Switzerland and its regions:

$$(21) \quad \text{VOL} = \alpha_0 + \alpha_1\text{MCH} + \alpha_2\text{CONPROP} + \alpha_3\text{INFORESTPROP} \\ + \alpha_{4:15}\text{REGION:DIST} + \alpha_{16}\text{MCH:INFORESTPROP} \\ + \alpha_{17:28}\text{MCH:REGION:DIST} + \epsilon$$

The model represented by eq. 21 was selected so that the interpretation of the coefficients could be intuitively inferred. For parsimony, the interaction between REGION:DIST and INFORESTPROP was dropped as it was difficult to interpret and only negligibly increased the adjusted $R^2$ when applied to all of Switzerland.

The inclusion of CONPROP as opposed to VARCH is somewhat surprising and is likely due to a correction of underestimated MCH on plots where many coniferous trees with pointier crowns are present. The raster point is less likely to be aligned with the highest part of the crown, thus explaining the positive coefficient for CONPROP in all regions. The interaction MCH:INFORESTPROP accounts for plots that are likely to be close to the forest boundary, and MCH:REGION:DIST is the disturbance-corrected prediction using MCH.

The second-phase model, $M_{\text{full}}$, incorporating PREV_VOLUME and all of the variables present in eq. 21, is

$$(22) \quad \text{VOL} = \beta_0 + \beta_1\text{MCH} + \beta_2\text{CONPROP} + \beta_3\text{INFORESTPROP} \\ + \beta_{4:15}\text{REGION:DIST} + \beta_{16}\text{MCH:INFORESTPROP} \\ + \beta_{17:28}\text{MCH:REGION:DIST} + \beta_{29}\text{PREV\_VOLUME} \\ + \beta_{30:41}\text{PREV\_VOLUME:REGION:DIST} + \epsilon$$

The presence of the MCH:PRODEG:DIST interaction term obscures the interpretation of the PREV_VOLUME:REGION:DIST coefficients. When MCH:REGION:DIST is dropped from both models, the coefficients have the expected interpretation of being smaller in disturbance classes corresponding to higher expectation of harvesting. However, the goal of the model selection is prediction, so both interactions were left in. It should be noted that there is a new forest strata imbedded in the REGION:DIST variable that corresponds to plots that do not contain any information about previous measurement. If zero is imputed for PREV_VOLUME in these strata, the design matrix will be singular and the model will be unable to be calculated. A simple workaround to circumvent this computational issue is to arbitrarily impute a value of 1 if the ID number of the plot is even and 0 if it is odd. This negligibly affects the residuals of the model and allows the model to be computed.

It should also be noted that the growth update models unfortunately did not significantly improve the correlations within region between ground truth observed in the third Swiss NFI and the update from the second Swiss NFI. This is because the time between remeasurements were all between 9 and 14 vegetative cycles. When the phase-in of the annual design is complete, the potential vegetative cycles range from 1 to 9. However, during the phase-in, the maximum number of cycles is potentially as much as 15. Thus, the available Swiss NFI data were not adequate to test the efficacy of the growth models in these situations. The REG automatically adjusts for average growth (both tree growth and in-growth) when the regression coefficient for PREV_VOLUME is more than 1. As a result, only one of the estimators presented here contains PREV_VOLUME_UPDATED.

### 4.8. Estimators considered

Nine estimators were considered to assess the efficiency of the three-phase REG and its update components. Estimates 1 and 2 are based on simple random sampling without any model (i.e., mean

**Table 3.** Estimates and models considered (results in Table 4).

| Estimate | $\mathbf{Z}^{(1)t}(x)$ | $\mathbf{Z}^{(2)t}(x)$ |
|---|---|---|
| 1* | REGION | — |
| 2* | REGION | — |
| 3 | MCH CONPROP INFORESTPROP REGION:DIST MCH:INFORESTPROP MCH:REGION:DIST | — |
| 4 | — | PREV_VOLUME PREV_VOLUME:REGION |
| 5 | — | PREV_VOLUME_UPDATED PREV_VOLUME_UPDATED:REGION |
| 6 | MCH CONPROP INFORESTPROP REGION:DIST MCH:INFORESTPROP MCH:REGION:DIST | PREV_VOLUME PREV_VOLUME:REGION:DIST |
| 7 | MCH CONPROP INFORESTPROP REGION:DIST MCH:INFORESTPROP MCH:REGION:DIST | PREV_VOLUME_UPDATED PREV_VOLUME_UPDATED:REGION:DIST |
| 8 | MCH CONPROP INFORESTPROP STRATUM MCH:INFORESTPROP MCH:REGION | PREV_VOLUME PREV_VOLUME:REGION |
| 9 | MCH CONPROP INFORESTPROP STRATUM MCH:INFORESTPROP MCH:REGION | PREV_VOLUME_UPDATED PREV_VOLUME_UPDATED:REGION |

**Note:** A colon ":" indicates a variable interaction term. For estimates 6–9 in the Pre-Alps (PA), there was a computational singularity in the model fit due to multicollinearity associated with the interaction term between the previous measurement and the "new forest" stratum. Thus, for estimates in the Pre-Alps, the interaction terms are dropped in $\mathbf{Z}^{(2)t}(x)$.

*Within each region, a standard estimate assuming simple random sampling was used, except for Switzerland (CH), for which double sampling for stratification is used with the production region as strata (this matches more closely to the current estimation procedure for the Swiss NFI that uses a stratification by region).

of the plot volumes) for each region. Double sampling for post-stratification using REGION as the stratification variable was used for the estimate across entire Switzerland, because it more closely corresponds to the current estimation procedure implemented in the Swiss NFI, which stratifies by region. Estimate 1 is the estimator used in the Swiss NFI under the old periodic design and takes into account all forested plots in the entire terrestrial sample of the third NFI. Estimate 2 only considers forested plots corresponding to the first 3 years of annual subsamples in the fourth NFI. Comparing estimate 2 with estimate 1 gives insight into the loss in precision associated with changing from a periodic to a continuous design when the standard SRS estimator is used to produce an estimate for a 3-year time period (note that the periodic Swiss NFI took 3 years to evaluate).

Estimates 3, 4, and 5 utilize the two-phase REG. Estimate 3 removes the second phase altogether and considers a model using only remote sensing variables in the auxiliary vector. This is the effect of ignoring the previous measurement altogether. Estimate 4 removes the first phase and only considers the previous measurement and

**Table 4.** Point estimates and standard errors (in parentheses) for timber volume (in m³·ha⁻¹), and adjusted $R^2$ values (in brackets) for $M_{reduced}$ and $M_{full}$, respectively, based on estimates and models from Table 3 for the Third Swiss National Forest Inventory.

| Estimate | Jura | Plateau | Pre-Alps | Alps | South Alps | Switzerland |
|---|---|---|---|---|---|---|
| 1 | 364.35 (6.51) [—] | 392.30 (7.81) [—] | 442.36 (9.12) [—] | 308.54 (5.41) [—] | 228.61 (5.99) [—] | 348.95 (3.32) [0.03] |
| 2 | 368.10 (11.30) [—] | 397.39 (13.64) [—] | 440.01 (15.33) [—] | 309.56 (9.21) [—] | 235.32 (10.88) [—] | 351.41 (5.69) [0.03] |
| 3 | 368.54 (9.64) [0.28] | 397.47 (11.29) [0.34] | 437.21 (13.51) [0.28] | 309.32 (7.85) [0.34] | 231.97 (8.89) [0.38] | 350.05 (4.66) [0.36] |
| 4 | 367.16 (9.17) [0.53] | 392.81 (11.68) [0.39] | 440.10 (12.80) [0.49] | 310.83 (6.87) [0.72] | 227.82 (7.08) [0.86] | 349.74 (4.43) [0.59] |
| 5 | 365.28 (9.08) [0.54] | 392.57 (11.56) [0.41] | 440.72 (12.80) [0.49] | 311.08 (6.85) [0.72] | 229.13 (6.97) [0.89] | 349.73 (4.41) [0.60] |
| 6 | 369.23 (8.11) [0.28; 0.60] | 394.03 (10.23) [0.34; 0.52] | 439.41 (12.14) [0.28; 0.52] | 311.62 (6.25) [0.34; 0.75] | 227.95 (6.51) [0.38; 0.87] | 350.49 (4.01) [0.36; 0.65] |
| 7 | 367.85 (7.98) [0.28; 0.62] | 395.24 (10.16) [0.34; 0.53] | 440.48 (12.20) [0.28; 0.52] | 311.95 (6.24) [0.34; 0.75] | 229.35 (6.41) [0.38; 0.89] | 350.92 (4.00) [0.36; 0.65] |
| 8 | 367.06 (8.33) [0.28; 0.58] | 393.20 (10.49) [0.34; 0.47] | 439.93 (12.17) [0.28; 0.51] | 310.71 (6.40) [0.33; 0.74] | 227.95 (6.51) [0.38; 0.87] | 349.91 (4.10) [0.34; 0.62] |
| 9 | 365.46 (8.22) [0.28; 0.59] | 392.85 (10.41) [0.34; 0.49] | 440.84 (12.24) [0.28; 0.51] | 310.95 (6.39) [0.33; 0.74] | 229.35 (6.41) [0.38; 0.89] | 349.91 (4.09) [0.34; 0.63] |

region information at the second phase without the growth update. Estimate 5 does the same as estimate 4 but incorporates the growth update. Estimates 4 and 5 demonstrate the effect of not ignoring the remote sensing data.

Estimates 6, 7, 8, and 9 utilize the three-phase REG. Estimates 6 and 7 incorporate the selected model, $M_{reduced}$, and $M_{full}$ in which variables PREV_VOLUME and PREV_VOLUME_UPDATED are used, respectively. Estimate 7 shows the effect of the growth update in the proposed three-phase REG setup. Finally, estimates 8 and 9 are the same as estimates 6 and 7, respectively, except the disturbance classes are removed. However, the "new forest" stratum is still included and is represented by the dummy variable STRATUM. These estimates give insight into the efficacy of incorporating the proposed disturbance update. All estimators are presented in detail in Table 3. The sample sizes for each phase by region are found in Table 1.

### 4.9. Results

The estimates with their corresponding variance estimates are found in Table 4. The estimates are all well within 2 standard deviations of the full third NFI estimate that contains the full sample, which is expected given their unbiasedness. As anticipated in large samples, the differences between the two variance estimates under the external model assumption (not shown) and the variances derived using the internal model are small for each region and empirically nonexistent for Switzerland as a whole (CH). In simulations, the external variance was usually smaller; however, here it was slightly larger for Jura (JU) and Plateau (SP).

The improvements to the standard errors should be compared with the standard estimate (i.e., estimate 2) that only contains the first 3 years of annual sample to assess the gain in efficiency. The two-phase REG that uses only remote sensing data (estimate 3)

shows a substantial improvement in variance across all regions and for CH, but not as much of an improvement as using only the previous measurement (estimates 4 and 5) (note that for CH, this also includes a simple stratification by region). The only exception was in the SP, where the adjusted $R^2$ for estimate 4 was the lowest (0.39) among all regions. The likely explanation is that despite the fact that the adjusted $R^2$ was higher than in estimate 3 (0.34), the effect of the larger first-phase sample size in estimate 3 was great enough to produce a lower standard error. The growth update offered only a slight improvement in estimate 5 over estimate 4.

Estimates 6, 7, 8, and 9 using the three-phase REG show a further decrease in standard error compared with estimates 3, 4, and 5. When the disturbance correction was removed in estimates 8 and 9, there was a slight increase in variance (note that there is no disturbance correction in South Alps, so it always remains exactly the same). Disturbance effects account for only a small proportion of the overall growing stock, so this is to be expected. As expected, there is a greater reduction in standard error in regions with heavier harvesting such as SP and JU. The Alps (AL) and South Alps (SA) have lighter harvesting, which is demonstrated by the negligible effect of their disturbance updates. Regions with fewer disturbances showed greater reductions in standard error compared with estimate 2 when the previous measurement was included in the model.

The growth update's effect when included in the three-phase REG was very modest, as evidenced by only slight decreases in standard error in most cases. The only exception was in the Pre-Alps, where the update had a slightly adverse effect; however, in the Pre-Alps, the interaction terms were dropped from $\mathbf{Z}^{(2)t}(x)$ for estimates 6–9 due to computational singularities in the model fit. In most cases, growth updating using only tree-level growth mod-

els provided a slight improvement during the estimation process. We note that the regression estimator automatically internally accounts for part of the overall growth, when the coefficient for the previous measurement variable is greater than one. In this case, it seems that the tree-level growth update did not bring much improvement. Further improvement might be possible if in-growth was accounted for.

## 5. Conclusions

The two- and three-phase REGs show clear usefulness in integrating remote sensing and past inventory data under an annual design. Furthermore, their implementation into statistical software is simplified by their close connection to classical linear regression. The estimation procedure in which remote sensing is used in the first phase and the previous inventory data are included in the second phase offers substantial reduction in the variance compared with using only the current annual subsample, with the added benefit of virtually no trade-off in bias. The updates for both growth and disturbance of the previous measurement variable provided only a modest improvement in standard error. Regardless, these estimators and this estimation scheme can clearly be recommended for use in regional and national forest estimation for growing stock volume. Their generality and flexibility allow for implementation with a wide variety of data sources, including LiDAR, satellite images, and high-resolution digital aerial photographs. The inclusion of a large first-phase sample as opposed to wall-to-wall data provides for great gains in computational efficiency and thus is an attractive alternative for large-scale national forest inventories.

## References

Brassel, P., and Lischke, H. 2001. Swiss National Forest Inventory: methods and models of the second assessment. Technical report, WSL Swiss Federal Research Institute, Birmensdorf.

Dixon, B., and Howitt, R. 1979. Continuous forest inventory using a linear filter. For. Sci. **20**(4): 675–689.

Eskelson, B.N.I., Temesgen, H., and Barrett, T.M. 2009. Estimating current forest attributes from paneled inventory data using plot-level imputation: a study from the Pacific Northwest. For. Sci. **55**: 64–71.

Gregoire, T., Ståhl, G., Næsset, E., Gobakken, T., Nelson, R., and Holm, S. 2011. Model-assisted estimation of biomass in a LiDAR sample survey in Hedmark County, Norway. Can. J. For. Res. **41**(1): 83–95. doi:10.1139/X10-195.

Houillier, F., and Pierrat, J.-C. 1992. Application des modèles statistiques spatio-temporels aux échantillonnages forestiers successifs. Can. J. For. Res. **20**(12): 1988–1995. doi:10.1139/x92-259.

Johnson, D., Williams, M., and Czaplewski, R. 2003. Comparison of estimators for rolling samples using forest inventory and analysis data. For. Sci. **49**: 50–63.

Keller, W. 1978. Einfacher ertragskundlicher Bonitätsschlüssel für Waldbestände in der Schweiz. Technical Report 54, Eidgenössische Anstalt für das forstliche Versuchswesen.

Lanz, A., Brändli, U.-B., Brassel, P., Ginzler, C., Kaufmann, E., and Thürig, E. 2010. Switzerland. *In* National forest inventories: pathways for common reporting. Chapter 36. *Edited by* E. Tomppo, T. Gschwantner, M. Lawrence, R. McRoberts. Springer. pp. 555–565.

Lessard, V., McRoberts, R., and Holdaway, M. 2001. Imputation on model-based updating techniques for annual forest inventories. For. Sci. **47**: 301–310.

Lüpke, N., Hansen, J., and Saborowski, J. 2012. A three-phase sampling procedure for continuous forest inventory with partial re-measurement and updating of terrestrial sample plots. Eur. J. For. Res. **131**: 1979–1990. doi:10.1007/s10342-012-0648-z.

Mandallaz, D. 2008. Sampling techniques for forest inventories. Chapman and Hall, Boca Raton, Florida.

Mandallaz, D. 2013*a*. Design-based properties of some small-area estimators in forest inventory with two-phase sampling. Can. J. For. Res. **43**(5): 441–449. doi:10.1139/cjfr-2012-0381.

Mandallaz, D. 2013*b*. Regression estimators in forest inventories within three-phase sampling and two multivariate components of auxiliary information. Technical report, ETH Zurich, Department of Environmental Systems Science. Available from http://e-collection.library.ethz.ch.

Mandallaz, D. 2014. A three-phase sampling extension of the generalized regression estimator with partially exhaustive information. Can. J. For. Res. **44**(4): 383–388. doi:10.1139/cjfr-2013-0449.

Mandallaz, D., Breschan, J., and 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. Can. J. For. Res. **43**(11): 1023–1031. doi:10.1139/cjfr-2013-0181.

McRoberts, R. 2001. Imputation on model-based updating techniques for annual forest inventories. For. Sci. **47**: 322–330.

Särndal, C., Swenson, B., and Wretman, J. 2003. Model assisted survey sampling. Springer Series in Statistics, New York.

Scott, C., Köhl, M., and Schnellbächer, H.J. 1999. A comparison of periodic and annual forest surveys. For. Sci. **45**(3): 433–451.

Van Deusen, P. 1997. Annual first inventory statistical concepts with emphasis on multiple imputation. Can. J. For. Res. **27**(3): 379–384. doi:10.1139/x96-211.

Van Deusen, P. 2002. Comparison of some annual forest inventory estimators. Can. J. For. Res. **32**(11): 1992–1995. doi:10.1139/x02-115.