ORIGINAL PAPER

# Estimating above-ground biomass of trees: comparing Bayesian calibration with regression technique

**Jürgen Zell · Bernhard Bösch · Gerald Kändler**

**Abstract** The commitment to report greenhouse gas emissions requires an estimation of biomass stocks and their changes in forests. When this was first done, representative biomass functions for most common tree species were very often not available. In Germany, an estimation method based on solid volume was developed (expansion procedure). It is easy to apply because the required information is available for nearly all relevant tree species. However, the distributions of neither parameters nor prediction intervals are available. In this study, two different methods to estimate above-ground biomass for *Norway spruce* (*Picea abies*), *European beech* (*Fagus sylvatica*), and *Scots pine* (*Pinus sylvestris*) are compared. First, an approach based on information from the literature was used to predict above-ground biomass. It is basically the same method used in greenhouse gas reporting in Germany and was applied with prior and posterior parameters. Second, equations for direct estimation of biomass with standard regression techniques were developed. A sample of above-ground biomass of trees was measured in campaigns conducted previously to the third National Forest Inventory in Germany (2012). The data permitted the application of Bayesian calibration (BC) to estimate posterior distribution of the parameters for the expansion procedure. Moreover, BC enables the calculation of prediction intervals which are necessary for error estimations required for reporting. The two methods are compared with regard to predictive accuracy via cross-validation, under varying sample sizes. Our findings show that BC of the expansion procedure performs better, especially when sample size is small. We therefore encourage the use of existing knowledge together with small samples of observed biomass (e.g., for rare tree species) to gain predictive accuracy in biomass estimation.

**Keywords** MCMC · Bayesian calibration · Error estimation · GHG-reporting · Biomass estimation · Cross-validation

## Introduction

Biomass estimation in forests is a topic of great interest, driven by at least two developments: (1) changes of biomass in forests correspond directly to changes in carbon absorbed or released to the atmosphere and are therefore the focus of global politically relevant mechanisms (Kyoto Protocol and resulting reporting commitments) and (2) it is important to calculate biomass because forests face an increased demand for wood energy. For both cases, the estimation of biomass should be accurate, efficient and provide an error estimate.

Estimation of biomass ($B$) in forests is usually based on biomass functions using standard measurements such as the diameter at breast height ($d_{1.3}$). Among a variety of biomass functions (Wirth et al. 2004; Zianis et al. 2005; Muukkonen 2007; Wutzler et al. 2008), one standard functional form is the use of the allometric relationship (Pretzsch 2001) with $d_{1.3}$ and/or height of trees ($h$) as explanatory variables. This functional form has shown to produce good estimates of biomass (e.g.: $B = \beta_0 d_{1.3}^{\beta_1} h^{\beta_2}$).

The other approach to estimate biomass of trees is based on appropriate use of existing knowledge. This procedure

J. Zell (✉)
WSL, Zürcherstrasse 111, 8903 Birmensdorf, Switzerland
e-mail: juergen.zell@wsl.ch

B. Bösch · G. Kändler
FVA, Wonnhaldestrasse 4, 79100 Freiburg, Germany

**Table 1** Ranges of tabulated values and number of observations used to construct the tables of Grundner and Schwappach (1952)

| | $d_{1.3}$ (cm) | | $h$ (m) | | $V_s$ (m$^3$) | | $V_t$ (m$^3$) | | $n$ |
|---|---|---|---|---|---|---|---|---|---|
| | Min | Max | Min | Max | Min | Max | Min | Max | |
| *Norway spruce* | 8 | 85 | 6 | 47 | 0.010 | 9.95 | 0.027 | 11.92 | 22,757 |
| *European beech* | 6 | 72 | 9 | 38 | 0.002 | 8.28 | 0.017 | 9.27 | 12,180 |
| *Scots pine* | 7 | 70 | 6 | 40 | 0.006 | 6.50 | 0.021 | 6.69 | 17,059 |

was used in official greenhouse-gas reports ("National Inventory Report") for the Kyoto Protocol in Germany (Umweltbundesamt 2009, p. 455, eq. 11). The estimation starts with the well-known solid volume ($V_s$), measured in m$^3$ with trees of at least 7 cm in diameter, and depends on diameter, length and stem-form. In recent years, it has been the focus of forest researchers to estimate solid volume precisely. The functions of Kublin (2003) are able to predict solid volume ($V_s$) of a tree using the variables $d_{1.3}$, $h$ and a further diameter at a second stem height (usually at 7 m, $d_7$), which defines a certain stem-form for a given diameter–height relationship.

Furthermore, based on the volume tables of Grundner and Schwappach (1952), it is possible to estimate total above-ground volume ($V_t$), given the solid volume. The tables rely on a broad data base and contain tabulated values on solid and total above-ground volume for common tree species, diameters and heights (see Table 1 and second row in Fig. 2). Based on the tables, it is possible to fit functions for the expansion from solid volume $V_s$ to total above-ground volume $V_t$. Finally, multiplying mean basic densities by total above-ground volume yields an estimate of total above-ground biomass. Hereafter, this approach will be called the "Expansion Procedure" (EP).

Zapata-Cuartas et al. (2012) introduced Bayesian standard techniques to estimate tree biomass with high precision and small sample sizes. They used the linearized form of the simplified biomass function and collected published parameter estimates in order to obtain prior information. Our approach is considerably different with regards to a variety of assumptions. One difference is the error term—we use an additive error with variance function—hence, we did not transform the data to a logarithmic scale. Further, our models are based on more than one predictor variable. The most important difference lies in the usage of prior information. We use old but also very large and common datasets. This requires a totally different formulation (namely, the EP), which can be seen as an indirect way of estimating biomass.

The aim of this article is twofold. Firstly, the currently used EP in official statistics for greenhouse gas reporting in Germany (Umweltbundesamt 2009) uses the EP with some parameter point estimates without knowledge of their distribution. The Bayesian calibration (BC) enables the construction of a sample of these distributions. They can be used to test for significance and to construct confidence intervals (in Bayesian notation, credibility intervals). Secondly, it seems clear that BC will be superior to any other estimate without prior knowledge if only small sample sizes are available. It remains unclear, however, how large the samples should be. Therefore, the effect of sample size was specifically analyzed, given the tree specific prior information. Therefore, the BC within the EP will be compared with an allometric regression approach for the tree species *Norway spruce*, *European beech* and *Scots pine*.

## Materials and methods

### Data

The above-ground biomass data were collected in several campaigns throughout Germany and in the State of Baden-Wurttemberg in preparation for the 3rd National Forest Inventory in 2013. Generally, since destructive measurement of trees is a demanding task and must be organized in collaboration with the forest service, a strictly randomized selection of samples is not feasible. We aimed at selecting sample trees covering the whole range of dimensions, and especially in the Germany-wide campaign, we collected the samples over a species-specific range after an analysis of species distribution based on National Forest Inventory data. It was assumed that stem-form is relevant in biomass estimation. Consequently, the data cover information on stem-forms such as diameter at breast height ($d_{1.3}$), total tree height and diameter at 7 m ($d_7$) or in the relative height of 30 % of total height. Sample size was 390, 218, and 127 for *N. spruce*, *E. beech*, and *S. pine*, respectively.

We applied two different procedures to estimate total above-ground biomass of the sample trees. With hardwoods, we applied randomized branch sampling (RBS) which is an efficient sampling technique, especially with large trees (Saborowski and Gaffrey 1999). Biomass of conifer species was assessed by measuring bole dimensions in 2-m sections to obtain volume and taking stem disks to estimate basic density (=oven dry mass/green volume) in order to convert volume to biomass. For the estimation of branch biomass, we used the following procedure. First, at each whorl, branches were counted. For a sample of
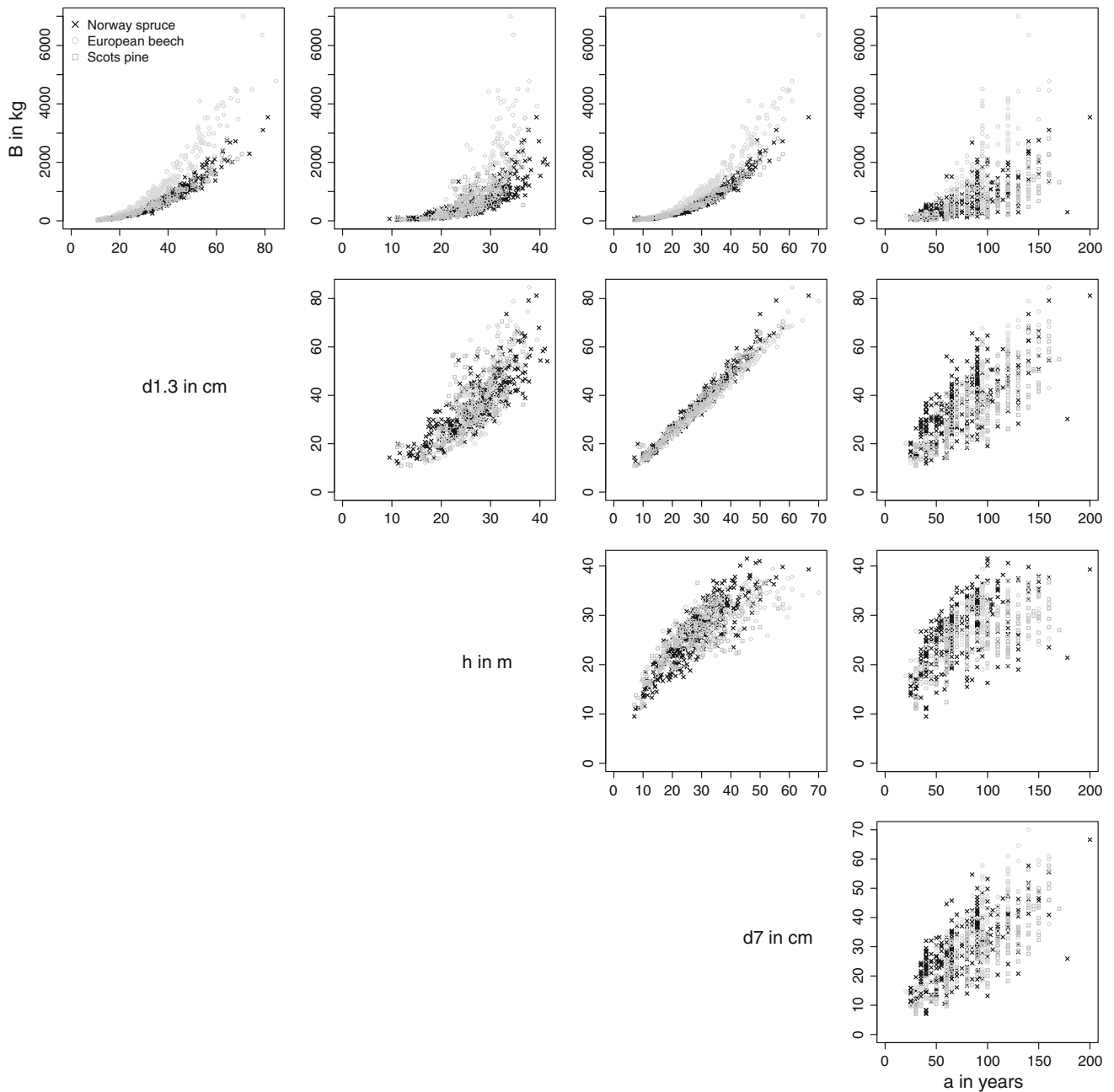
**Fig. 1** Sampled biomass data and explaining variables of the three tree species. Sample size was 390, 218 and 127 for *Norway spruce*, *European beech* and *Scots pine*. *a*: age in years

branches, base diameters were measured and dry matter was determined for a subsample. The biomass data were then pooled for all sampled trees of a particular species and an allometric model with branch biomass as response and base diameter as predictor was fitted. This was in turn used to estimate total branch biomass of a tree based on the base diameters of all branches. For branches without diameter measurements, diameters were imputed via a random draw from the diameter distribution of the distribution obtained for the individual tree.

In Fig. 1, all reasonable potential explanatory variables are shown as pairwise scatter plots. Although $d_{1.3}$ and $d_7$ are highly correlated, both variables remain in the regression model for *N. spruce* and *E. beech*.

Estimating parameters by regression

To estimate the allometric relationship with above-ground biomass, nonlinear regression analysis was used [function gnls in package nlme, Pinheiro and Bates (2000) in R, R

**(a)** Mass density, based on Kollmann (1982)



**(b)** Extracted basic densities, Weibull Distributions



**(c)** Ratio $V_t$ to $V_s$ over $d_{1.3}$ (cm)



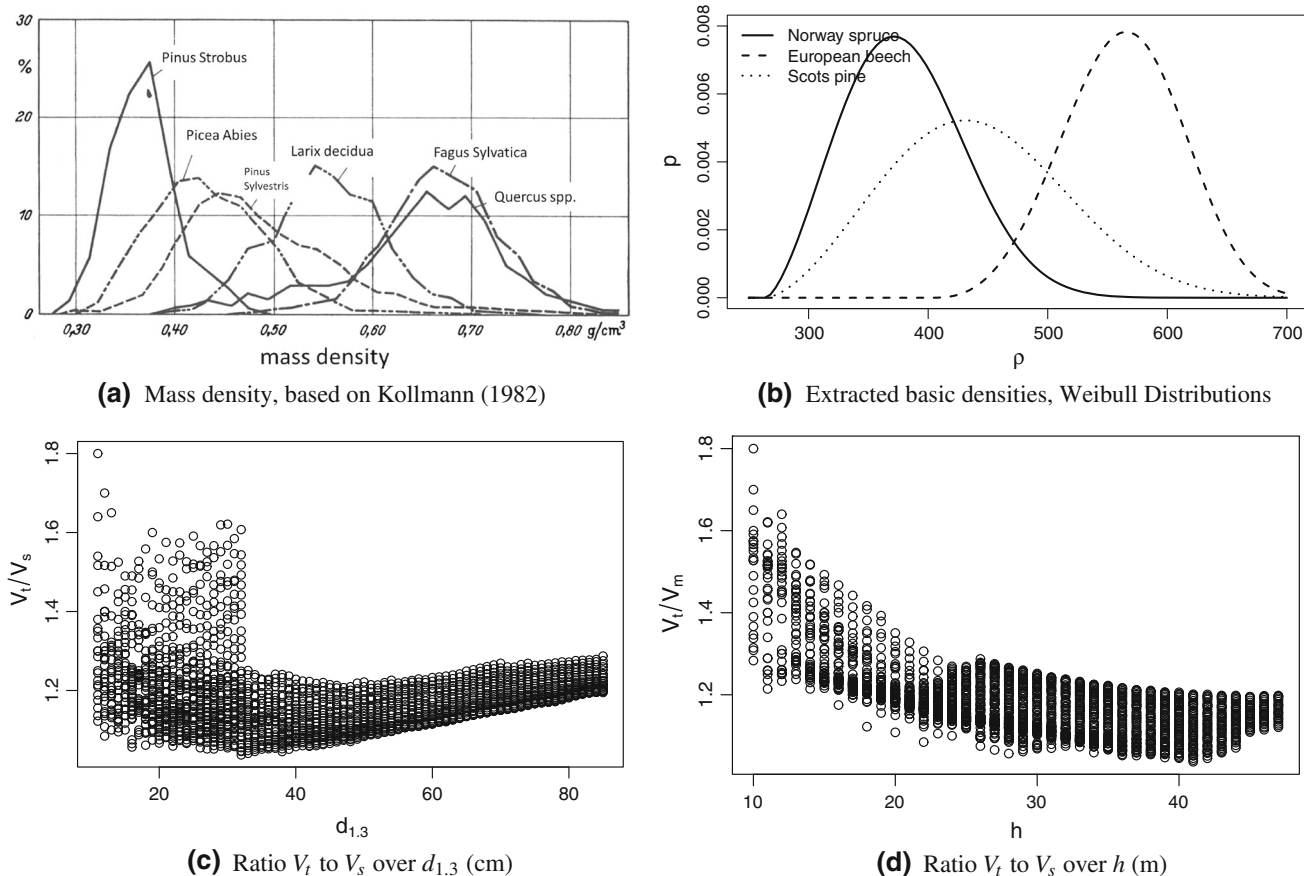**(d)** Ratio $V_t$ to $V_s$ over $h$ (m)

**Fig. 2** Example on the extraction of prior information for the expansion procedure. Second row **c** and **d** shows ratio of total to solid volume over $d_{1.3}$ and $h$. Values are the tables of Grundner and Schwappach (1952), the example is from *Norway spruce*

Development Core Team (2009)]. Above-ground biomass has a higher variability with increasing size. For model fitting, the data are often log transformed to linearize the equation and homogenize the variance. However, back-transformation introduces a bias in the expectation. While there are simple correction factors (Sprugel 1983), they also need assumptions. Therefore, we used the original scale of the data. Within the original scale, errors are additive. The increasing variance was modeled to depend directly on the estimated biomass ($\hat{B}_i$); hence, we assume the errors to be independent, but not identically distributed. Estimation started with a maximal model (Eq. 1), and nonsignificant terms were dropped until a final model was found. Only the allometric formulation was used in the analysis.

$$B_i = \beta_0 d_{1.3_i}{}^{\beta_1} h_i{}^{\beta_2} d_{7_i}{}^{\beta_3} a_i{}^{\beta_4} + \epsilon_i \tag{1}$$

Where $a$ is the age of the trees, $\beta_0$–$\beta_4$ are parameters to be estimated, and $\epsilon_i \sim N(0, \sigma^2 \hat{B}_i^{2\delta})$ has a non-constant variance. For $\epsilon_i$ different variance functions were tested:

1. Constant variance: $Var(\epsilon) = \sigma^2$
2. Exponential increasing variance: $Var(\epsilon) = \sigma^2 e^{2\delta \hat{B}}$

3. Variance increases with the power of estimated biomass: $Var(\epsilon) = \sigma^2 \hat{B}^{2\delta}$

Although type (2) and (3) require a further parameter δ, likelihood ratio test showed significant better results for type (3).

### Bayesian inference

#### Bayes rule

An introduction to the use of Bayes rule in the context of parameter estimation is given in Gelman et al. (2004). Denoting the parameters θ and the data $D$, gives the following formulation:

$$p(\theta \mid D) = \frac{p(D \mid \theta)p(\theta)}{\int p(\theta)p(D \mid \theta)d\theta} \tag{2}$$

In this formula $p(\theta)$ is called prior information. It comprises all is known about the parameter, before the data are measured. $p(D \mid \theta)$ is the likelihood of the data, given model output. $p(\theta \mid D)$ is called posterior probability

distribution and it contains all that is known about the parameters, after the data have been measured, inclusively the prior information. The denominator contains a multidimensional integral, which can be solved analytically by using appropriate conjugated prior distributions or it can be approximated by simulation (Gilks et al. 1995).

*Expansion procedure and prior information*

The expansion procedure contains several steps. First, the solid volume is estimated. This is an important step since this volume and the biomass are closely related. Then, an expansion to total above-ground volume takes place, which is of less importance, because it is in the range of a factor $\sim 1.2$. The multiplication by the basic densities is then again an important step, containing valuable new information. Solid volume $V_s = f(d_{1.3}, h, d_7)$ is expressed as a function of measured variables. It is predicted by the stem-form function from Kublin (2003), containing several cubic regression splines. The stem-form functions are based on section-wise (2 m) measurements of logs and are based on a nationwide data base. They are used by default in official statistics and reports, like the National Forest Inventory (BMELV 2009).

Total above-ground volume is then derived based on the tables from Grundner and Schwappach (1952). These tables were digitized, and the values are used as pseudo-observations. In Fig. 2c and d, the ratio of total to solid volume is plotted over $d_{1.3}$ and $h$ for *N. spruce*. This shows that the ratio is decreasing over diameter and height. Since total volume is always larger than solid volume, the logistic function was used such that any prediction will stay in a range of [0,1] for the expansion from solid to total volume. This function was fitted, and the estimated parameters $\theta$ are further used in EP.

$$V_t = V_s + V_s \frac{1}{1 + \exp(-(X\theta))} \tag{3}$$

Where $X$ is a design matrix containing the predictor variables $d_{1.3}$ and $h$. Since the regression is based on tabulated mean values, they have an unknown error. The tabulated values do not reflect the original uncertainty of measurements, nevertheless they still contain the relationship between solid and total volume in its dependency on diameter and height. Since there is no useful estimation of the standard errors of these parameters, a wide conservative standard error was assumed so that $\pm 2\sigma$ just overlaps zero. This corresponds to a weakly significant parameter, although the amount of underlying data (which are not available) would certainly result in smaller standard errors.

Multiplying total volume by basic density gives an estimate of above-ground biomass. Given that mass density has been the focus of forest research over the past century,

not only the mean and standard deviation are known. Based on standard literature on forest technology (Kollmann 1982), the distribution of mass density for large samples is also known for common tree species in middle Europe (7,112 for *N. spruce*, 1,778 for *E. beech* and 2,418 for *S. pine*, see figure 2(a)).

Basic density in the EP is considered to be a parameter; hence, its prior distribution is of relevance. It was derived from an illustration in Kollmann (1982) and is presented in 2(a), since the original values are not available. The picture was overlaid with a digital polygon. Then subsamples of mass density were drawn proportional to the relative frequency. Based on the subsamples, a Weibull distribution was fitted, with an offset in the x-axis. The resulting prior distribution for basic density is given in Fig. 2(b). The offsets were estimated to be 265, 390 and 269 kg/m$^3$. Lower values are excluded from the posterior, since these offsets can be seen as reasonable lower physiological boundaries for basic densities (Hakkila 1972; Kollmann 1982).

*The likelihood and Metropolis Hastings algorithm*

Likelihood (see Eqs. 2 and 44) is the probability of observing the data given by the model output. It is calculated as the product of densities of normally distributed measurement errors, in this case the difference between observed and expected biomass $(B - \hat{B})$. As already discussed above in "Estimating parameters by regression" section, in the regression models in Eq. 1, variance in the biomass is not constant. Therefore, $\sigma^2$ was replaced by $\sigma^2 \hat{B}_i^{2\delta}$ resulting in the logarithmized likelihood ($\log(L)$), with $N$ being the total number of observations:

$$\log(L) = -\frac{1}{2} \sum_{i=1}^{N} \log\left(2\pi\sigma^2 \hat{B}_i^{2\delta}\right) - \sum_{i=1}^{N} \frac{\left(B_i - \hat{B}_i\right)^2}{2\sigma^2 \hat{B}_i^{2\delta}} \tag{4}$$

Since this likelihood is not a standard formulation of a normal model, we programmed a sampler based on Metropolis Hastings algorithm. This was originally described by Metropolis et al. (1953) and can be seen as a walk through the parameter space such that the visited points in the chain are a sample of the posterior distribution. In each step, a candidate parameter vector is generated randomly. For the generation of new proposal values, a covariance matrix is needed. In our application, the covariance matrix was adjusted manually, by testing short chains. Programming was done with R (R Development Core Team 2009), using library MASS (Venables and Ripley 2002) with function mvrnorm to produce proposals.

The expansion procedure and data were used to calculate the expected biomass $(\hat{B} = f(D, \theta))$. By using the Metropolis Hastings algorithm, a sample of the posterior distribution of the parameters was generated. The sampled

chain length was set to 400,000 and a thinning of 200 was applied to receive efficient estimates of the posterior probability distribution. Burn-in was discarded and set to 10 % of the chain length. Initially, different starting values were also tested to see whether the chain converges to the same level.

A new parameter vector can be accepted or rejected. Acceptance depends on the Metropolis ratio, which is the ratio of likelihood multiplied by the prior value for the candidate to the likelihood multiplied by the prior for the current parameter. If this ratio is above 1, the candidate will be accepted; if it is below one, it is accepted with probability equal to the Metropolis ratio. Proposal variances were changed to reach acceptance rates to lie between 20 and 30 % (Gelman et al. 2004, p. 307). The sampling was stopped, if the chains converged and the whole parameter space was visited. This was evaluated graphically by plotting the trace of the chains.

Final expansion procedure models

With the complete dataset, final models of both methods were also estimated. In the case of the regression models, nonsignificant parameters were dropped stepwise. The BC models needed re-parameterization because the high correlation of parameters resulted in unstable chains in the posterior. Re-parameterization was either done by transformation of variables or by dropping variables from the model. This resulted in BC models that contained only few parameters in the following forms:

– *N. spruce*

$$B = \left( \rho_0 - \theta_2 \log\left( 5 \frac{d_{1.3}}{a} \right) \right)$$
$$\left( V_s + V_s \frac{1}{1 + \exp(-(\theta_0 h + \theta_1 d_{1.3} h))} \right)$$

– *European beech*

$$B = \rho \left( V_s + V_s \frac{1}{1 + \exp(-(\theta_0 \log(a) + \theta_1 d_{1.3} h))} \right)$$

– *Scots pine*

$$B = \rho \left( V_s + V_s \frac{1}{1 + \exp(-(\theta_0 \log(h)))} \right)$$

With $\rho_0$ the basic density at (a theoretical) zero tree-ring width, otherwise $\rho$ is mean basic density of the above-ground tree.

For *N. spruce* four parameters were in the final BC model, two for the expansion, and two for basic density. For *S. pine* only two parameters were in the final model. Mass density shows for *N. spruce* a log-linear relationship

to mean tree-ring width ($10 d_{1.3}/(2a)$ in mm/y)[1], whereas diffuse-porous trees (like *E. beech*) show a weak linear relationship (Kollmann 1982). For *N. spruce,* this relation was also used in the BC of EP, whereas for *E. beech* and *S. pine*, the parameter for this relationship was extremely weak and was therefore omitted from the model.

Model comparison

Cross-validation was used to measure the predictive accuracy of the models. Consequently, random subsamples with fixed sizes were repeatedly chosen to partition a training and validation dataset. The fixed sizes of the training dataset were increased starting at 10 stepwise until 300 for *N. spruce*, 150 for *E. beech* and 100 for *S. pine*, based on the total sample size. With each training dataset size, 50 datasets were randomly partitioned. For each, a regression model (see Eq. 1) was fitted and the BC of the EP was also implemented.

To measure the predictive performance, a relative value was used, given that variance increases with size. Root Mean Square Percentage Error (RMSPE) was therefore used:

$$RMSPE = \sqrt{ \frac{1}{N} \sum_{i=1}^{N} \left( \frac{B_i - \hat{B}_i}{B_i} \right)^2 } \, 100 \% \qquad (5)$$

In addition, a one-sided sign test was applied to verify whether the null hypothesis "BC based on EP is not better than the predictions from allometric regression models" could be rejected. As a result, the absolute differences between observed and predicted values under both model types were thereby compared. If model one (regression) was closer to the data, the test statistic is negative, whereas if the converse is true (BC is closer), then positive:

$$d_i = |B_i - \hat{B}_{M1i}| - |B_i - \hat{B}_{M2i}|$$

The number of positive $d_i$ is binomially distributed and a one-sided binomial test was conducted on a 95 % confidence level (sign test). As mentioned above, this test was also repeated 50 times, i.e., under different data constellations. As a result, the proportion of rejected null hypotheses can be displayed over the size of training datasets.

# Results

Expansion procedure

After biomass data were sampled, it was possible to compare them with predictions based on prior knowledge of the

---

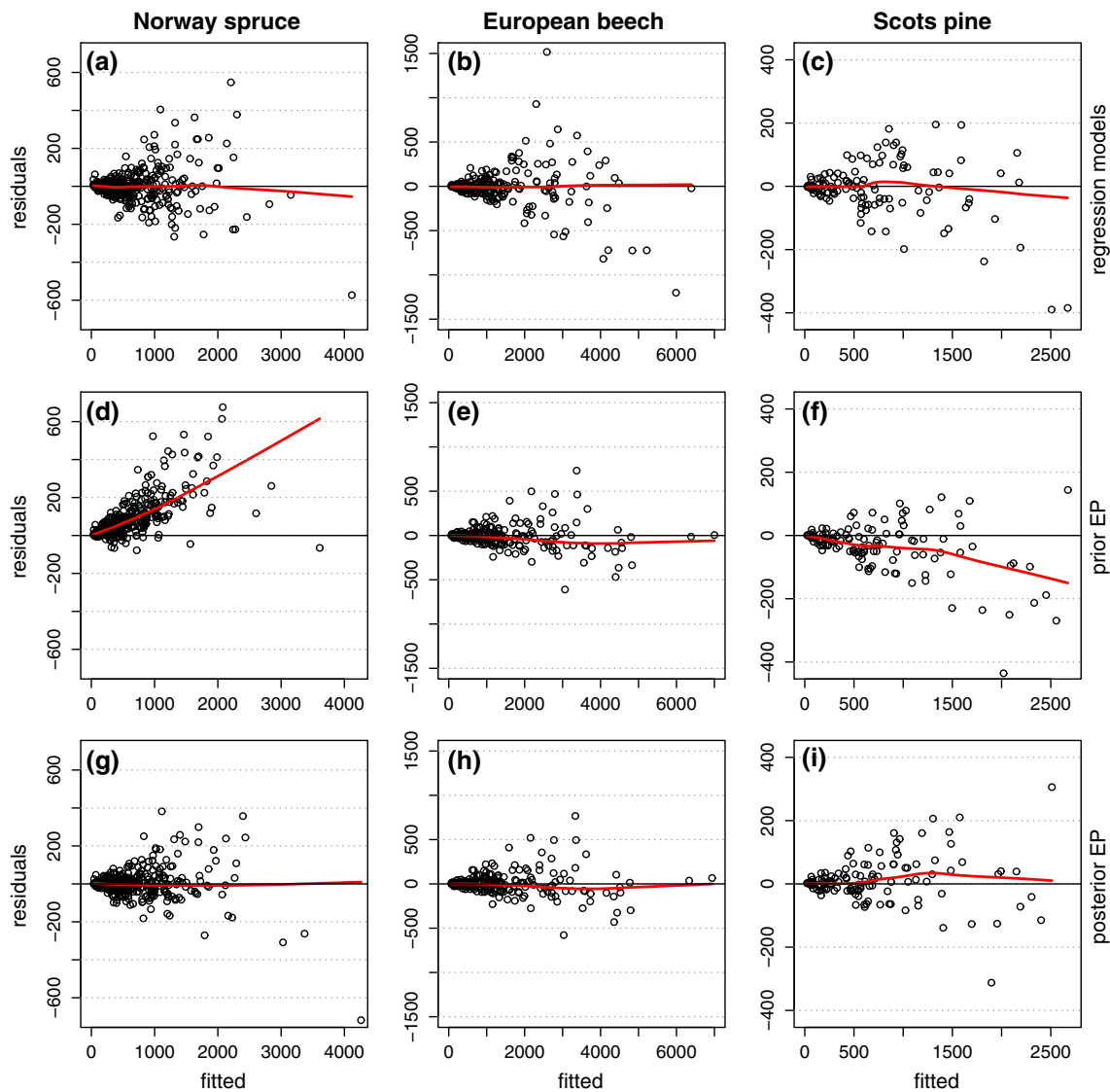[1] 10 because of a unit change from cm to mm and 1/2 because of the change from diameter to radius

**Fig. 3** Raw residuals $\hat{\epsilon} = B - \hat{B}$ for the tree species under different models. First row **a–c** residuals of the regression model. Second row **d–f** expansion procedure with prior information. Third row **g–i** expansion procedure with (mean) posterior of the parameters

EP. The comparison is shown as Tukey-Anscombe plots in Fig. 3 d–f and reveals clear and systematic underestimations for *N. spruce* and overestimation for *S. pine*. In contrast, *E. beech* prediction is in good accordance with the observed data and is visually indistinguishable from predictions based on the posterior distribution (after the data were observed). While further large residuals, especially in *E. beech*, are more common in the regression models, they are less common in the EP.

Figure 4 shows the 2.5 and 97.5 % quantiles and distribution of prior and posterior parameters. For *N. spruce*, the 95 % interval (based on the quantiles) drops from 337–528 in the prior, to 402–457 kg/m$^3$ in the posterior (see Fig. 4 d). In *E. beech*, the gain in precision of the

parameters is less obvious (see Fig. 4 e–g). For *S. pine*, a clear gain in precision of basic density and an abrupt drop at 450 kg/m$^3$ (Fig. 4 i can be observed.

Allometric regression models

The resulting final regression models and parameters are given in Table 2. For *N. spruce* and *E. beech*, the complete set of predictors are highly significant ($p < 0.001$, see Table 2), only in *E. beech* age has a $p$ value of 0.0215. In the allometric biomass function for *S. pine*, the variables $a$ and $d_7$ are not significant and are therefore excluded from the final model. The Tukey-Anscombe plots in Fig. 3 a–c confirm the independence of the (raw) residuals and show
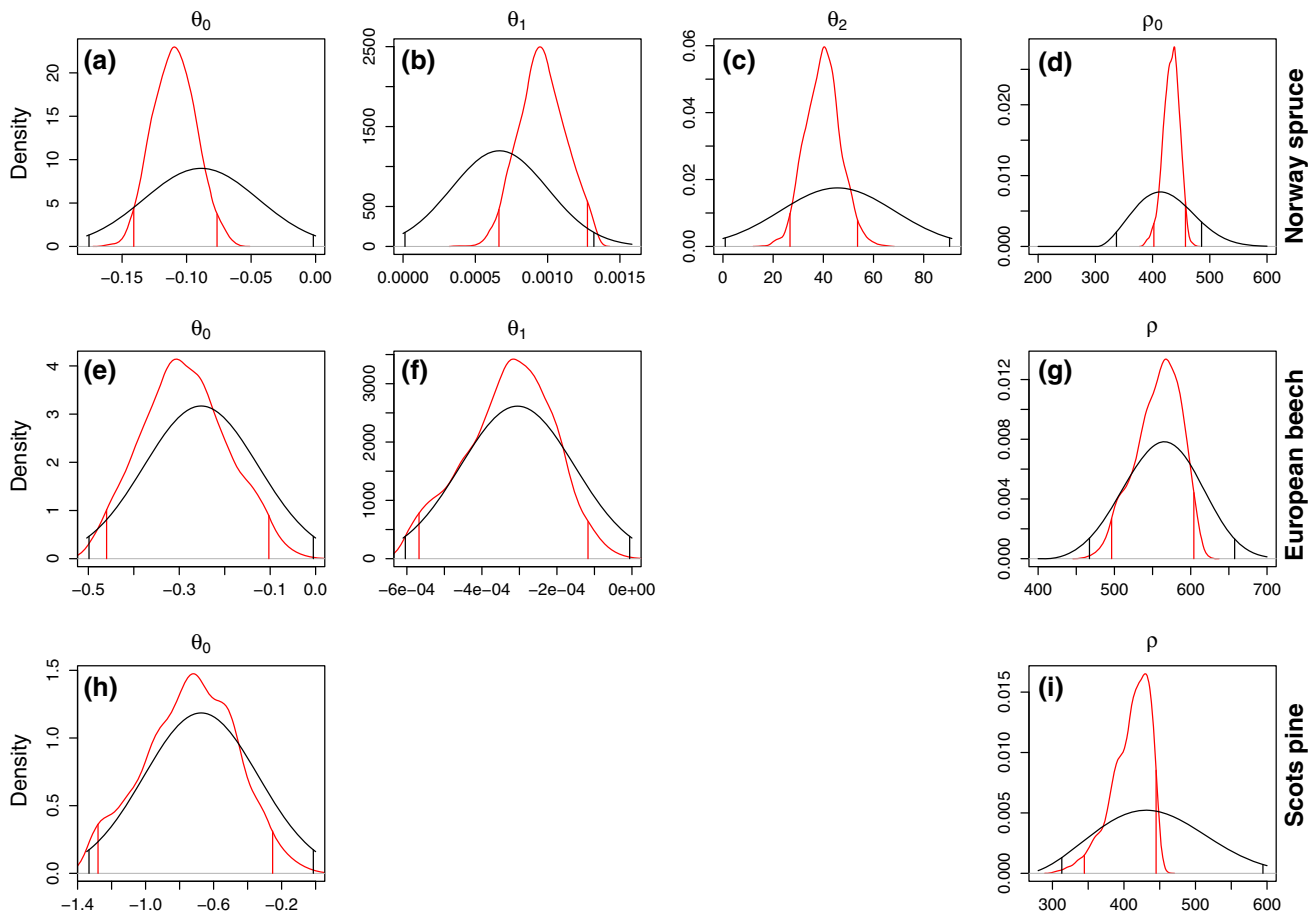
Fig. 4 Prior and posterior distribution of parameters of the BC models. $\rho_0$: basic density at zero tree-ring width, $\rho$: basic density, $\theta_0 - \theta_1$: expansion parameters, $\theta_2$: logarithmic decrease in basic density (kg/m$^3$) with mean tree-ring width (mm/year). 2.5 and 97.5 %-quantiles are presented as horizontal lines

**Table 2** Parameter estimates for *Norway spruce, European beech* and *Scots pine* for the allometric biomass functions. $B = \beta_0 d_{1.3}^{\beta_1} h^{\beta_2} d_7^{\beta_3} a^{\beta_4}$, *Scots pine* without $d_7$ and $a$

| Species | Param. | Value | se | t value | p value | $\hat{\sigma}$ | $\hat{\delta}$ |
|---|---|---|---|---|---|---|---|
| N. spruce | $\hat{\beta}_0$ | 0.074 | 0.0077 | 9.591 | <0.001 | 0.544 | 0.748 |
| | $\hat{\beta}_1$ | 0.993 | 0.0945 | 10.5089 | <0.001 | | |
| | $\hat{\beta}_2$ | 0.355 | 0.0534 | 6.6460 | <0.001 | | |
| | $\hat{\beta}_3$ | 1.055 | 0.1006 | 10.4887 | <0.001 | | |
| | $\hat{\beta}_4$ | 0.157 | 0.0215 | 7.3119 | <0.001 | | |
| E. beech | $\hat{\beta}_0$ | 0.0752 | 0.0134 | 5.6013 | <0.001 | 0.038 | 1.157 |
| | $\hat{\beta}_1$ | 0.8312 | 0.1461 | 5.6908 | <0.001 | | |
| | $\hat{\beta}_2$ | 0.6778 | 0.0736 | 9.2090 | <0.001 | | |
| | $\hat{\beta}_3$ | 1.3356 | 0.1503 | 8.8874 | <0.001 | | |
| | $\hat{\beta}_4$ | −0.0702 | 0.0303 | −2.3164 | 0.0215 | | |
| S. pine | $\hat{\beta}_0$ | 0.0235 | 0.0047 | 5.0073 | <0.001 | 0.1588 | 0.934 |
| | $\hat{\beta}_1$ | 2.2392 | 0.1005 | 22.2919 | <0.001 | | |
| | $\hat{\beta}_2$ | 0.6399 | 0.1170 | 5.4677 | <0.001 | | |

**Table 3** Parameter estimates of simplified models for *Norway spruce*, *European beech* and *Scots pine*. Two model types, $B = \beta_0 d_{1.3}^{\beta_1}$, and $B = \beta_0 d_{1.3}^{\beta_1} h^{\beta_2}$, are used. The model including $h$ for *Scots pine* is already presented in Table 2

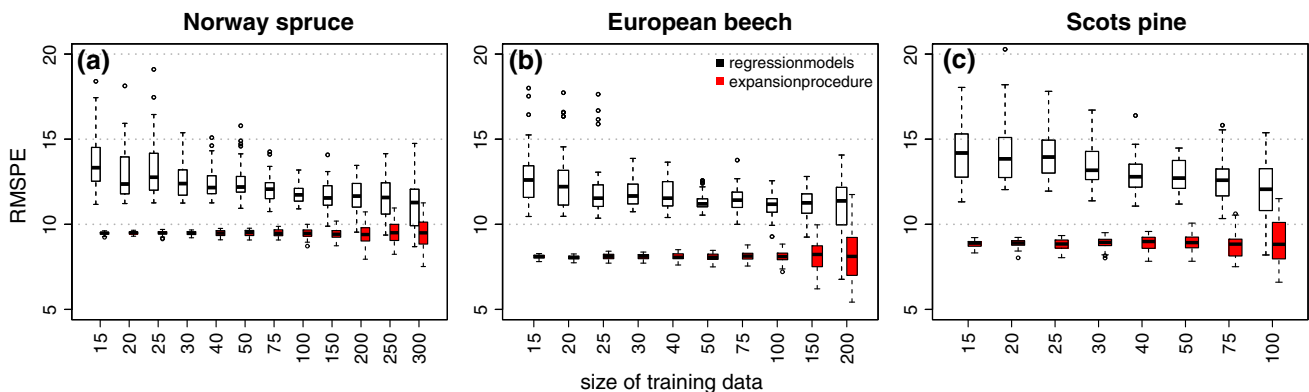| Species | Model | Param. | Value | se | $t$ value | $p$ value | $\hat{\sigma}$ | $\hat{\delta}$ |
|---------|-------|--------|-------|-----|-----------|-----------|------|------|
| N. spruce | $B = \beta_0 d_{1.3}^{\beta_1}$ | $\hat{\beta}_0$ | 0.1010 | 0.0077 | 13.0724 | <0.001 | 0.185 | 0.968 |
|  |  | $\hat{\beta}_1$ | 2.4134 | 0.0216 | 111.7073 | <0.001 |  |  |
| E. beech | $B = \beta_0 d_{1.3}^{\beta_1}$ | $\hat{\beta}_0$ | 0.1527 | 0.0149 | 10.2836 | <0.001 | 0.099 | 1.073 |
|  |  | $\hat{\beta}_1$ | 2.4511 | 0.0275 | 89.1028 | <0.001 |  |  |
| S. pine | $B = \beta_0 d_{1.3}^{\beta_1}$ | $\hat{\beta}_0$ | 0.0398 | 0.0093 | 4.2956 | <0.001 | 0.238 | 0.909 |
|  |  | $\hat{\beta}_1$ | 2.6966 | 0.0754 | 35.7461 | <0.001 |  |  |
| N. spruce | $B = \beta_0 d_{1.3}^{\beta_1} h^{\beta_2}$ | $\hat{\beta}_0$ | 0.0493 | 0.0044 | 11.1155 | <0.001 | 0.181 | 0.947 |
|  |  | $\hat{\beta}_1$ | 2.0319 | 0.0369 | 55.0796 | <0.001 |  |  |
|  |  | $\hat{\beta}_2$ | 0.6307 | 0.0529 | 11.9261 | <0.001 |  |  |
| E. beech | $B = \beta_0 d_{1.3}^{\beta_1} h^{\beta_2}$ | $\hat{\beta}_0$ | 0.0253 | 0.0038 | 6.6726 | <0.001 | 0.044 | 1.156 |
|  |  | $\hat{\beta}_1$ | 2.0559 | 0.0372 | 55.1969 | <0.001 |  |  |
|  |  | $\hat{\beta}_2$ | 0.9670 | 0.0735 | 13.1631 | <0.001 |  |  |



**Fig. 5** Root Mean Squared Percentage Error (RMSPE) for each 50 models based on different sizes of training datasets

that they are centered around zero. Further, for practical usage in standard inventories, a simplified version with predictors based only on $d_{1.3}$ and $d_{1.3}$ together with $h$ are presented in the "Appendix" in Table 3.

Prediction accuracy

Based on the 50 randomly chosen training datasets, the Root Mean Squared Percentage Error (RMSPE, Eq. 5) shows better prediction accuracy compared with the regression models (Fig. 5). The range of RMSPE values increases in BC with increasing sample size.

The sign test in Fig. 6 shows that BC is significantly better in prediction than regression models for sample sizes of <50. This predominance can be observed for nearly all data constellations and all tree species. The advantage of BC slightly falls with increasing sample size, but is still superior even when nearly all observations are used to build the models.
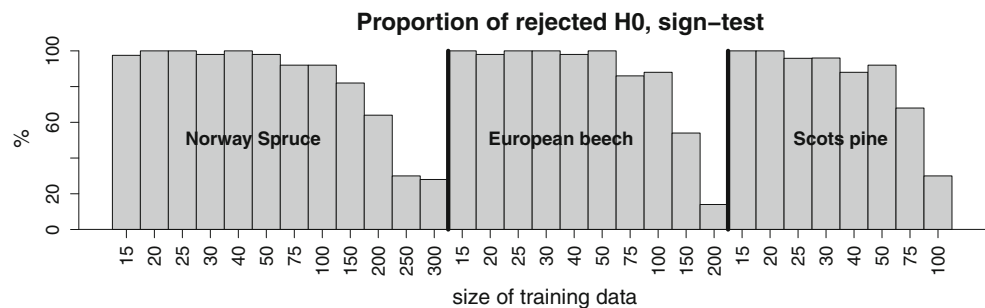
## Discussion

Expansion procedure and allometric regression models

Zapata-Cuartas et al. (2012) used prior information for two parameters of a simplified allometric biomass function ($log(B) = log(\beta_0) + \beta_1 \; log(d_{1.3})$). They were able to shown that sample sizes can be greatly reduced without loss of precision in RMSE. This is possible because Bayes theorem enables the use of prior knowledge in the process of parameter estimation, compared with classical fitting by least squares that can only make use of observed data. We found the same effects, although our analysis is different with respect to:

- *Tree species* our models are species-specific, since the amount of information depends heavily on the tree species. Further, we assume that allometry is species-dependent (see, e.g., Table 2). Since the EP uses wood density, it is by definition species-specific. However, in

**Fig. 6** Proportion of rejected sign tests for 50 randomly partitioned datasets over increasing training dataset size



the context of tropical forests (Zapata-Cuartas et al. 2012), with a diversity of tree species and little knowledge regarding a specific species, the use of general scaling rules is appropriate, and hence, the usage of general priors for any tree species is sensible.

- *Linearization* Zapata-Cuartas et al. (2012) used a linearized form in the allometry, enabling the usage of standard regression techniques and also implemented Bayesian techniques in statistical packages (e.g., MCMCglmm in R). Linearization has a large advantage, since heteroscedasticity disappears. However, it comes with the cost of a bias in back-transformation to the original scale (Jensens Inequality), although a correction can be estimated (Sprugel 1983). We did not linearize, since there are no good grounds for the use of a multiplicative error term in the original scale and the heteroscedasticity can be well handled by the variance function, as presented in "Estimating parameters by regression" section. By comparing the quantile–quantile plots of the residuals in both cases (transformed normally distributed and untransformed with variance function), the latter appears to be slightly better. The disadvantage of this uncommon variance function is its slightly more complicated likelihood function (see Eq. 4), which contains a non-constant first summand ($\hat{B}_i$). This is not directly tractable by the standard MCMC-sampler. Therefore, a sampler based on the Metropolis-Hastings algorithm was written and implemented in R.

- *Allometry versus EP* Zapata-Cuartas et al. (2012) priors are used directly for the two parameters in the allometric relationship. We used a more complicated formulation of the biomass, namely the EP. This was for two reasons. First, more information is available for the construction of EP, since it relies on standard measures, sampled and described in various books and articles (e.g., estimation of total volume by taper functions or distributions of basic densities). In total, Zapata-Cuartas et al. (2012) found 134 biomass functions, whereas ,e.g., only independent samples of basic densities for single trees number more than 10,000 (see Expansion procedure and prior information section).

The second reason is that the EP was used to estimate biomass for official statistics; hence, it is a common way to estimate biomass. Our analysis shows the underlying distribution of the parameters in use and enables the construction of prediction intervals. Both of these possibilities are new and became possible through the use of Bayesian calibration.

Large residuals of *E. beech* (also in *N. spruce*) are visible in the regression models, compared with the EP. Regression models are only based on observations, viz., there is no prior information. Since the EP already contains the solid volume, it stabilizes the predictions because it contains most of the total above-ground volume and hence biomass, resulting in less extreme residuals.

The abrupt change in basic density in *S. pine* may be an effect of altered silvicultural systems. *S. pine*—like *N. spruce*—has a decreasing basic density over mean tree-ring width (Kollmann 1982). It may be that the older *S. pine* trees used in Kollmann (1982) had smaller tree-ring widths, which was an effect of a different silvicultural treatment. Unfortunately, it was not possible to include this effect in the *S. pine* BC model because the parameters had excessively high correlations and the parameters where not significant (based on the 2.5 and 97.5 %-quantiles).

The effect of increasing RMSPE measurement for BC models can be explained directly by the method: the EP can produce estimates even without any directly observed biomass. Having a few observations, the prior information dominates the prediction. In the case of *E. beech*, where the prior information is already close to observed biomass data, the RMSPE measurement (Fig. 5 b is also lower in BC compared with the other tree species. The RMSPE is therefore smaller when prior information dominates the prediction. With increasing sample sizes, RMSPE shows similar ranges to those from the regression models.

## Conclusion and outlook

Biomass estimation in forests has gained importance in the recent years. Especially, the reporting commitment

within the scope of the Kyoto Protocol requires estimates of biomass in forests together with an error estimate. Due to the lack of representative biomass functions for common tree species (except those from Wirth et al. (2004) and Wutzler et al. (2008)), Germany has decided so far to use existing knowledge (basically the solid volume of trees) to predict biomass in forests (Umweltbundesamt 2009). Essentially, the old procedure is the same as the EP presented here. It is reasonable to use this source of information to get a best possible estimate of biomass, although the accuracy and bias in the prediction are unknown.

We found that prior information in the EP can result in excellent biomass predictions, as the example in Fig. 3 shows for *E. beech*. However, this could be a coincidence, because *N. spruce* and *S. pine* estimates show biases in large trees. A bias in large trees may have strong implications for the estimation of sink and sources in the carbon budget of forests. If, hypothetically, the underlying distribution of trees in forests changes to larger trees, then a negative bias in the applied functions can give a result of a decrease in carbon stock even though the opposite is in fact true (and vice versa).

Sampling biomass for large trees is expensive and destructive. Here we present a method, whereby small samples sizes can be efficiently used to construct biomass functions. This method predictive accuracy is highly competitive compared with conventional biomass functions. We therefore encourage the use of small biomass investigations of rare tree species and to join these data with prior information. For instance, stem-form functions are available for 36 different species. The tables of Grundner and Schwappach (1952) contain the nine most import species in Germany, and Kollmann (1982) has collected basic densities for nearly all species in Germany, at least with an estimate of means and ranges. Even prediction intervals can be generated based on the posterior distribution of parameters. The applicability of the proposed method may easily be expanded by simpler assumptions such as:

- stem volume functions based solely on diameter and/or height
- expansion from solid to total volume, e.g., based only on diameter
- basic density could easily be assumed to be normal; mean and standard error are given in many textbooks

A drawback of BC is the use of subjective prior information. Although we refer to previously published values as much as possible, we are aware that it may be seen as an influential way to predict biomass. However, the obvious advantages of the method far outweigh this limitation.

## Appendix

For practical usage in standard inventories, simplified models are presented in Table 3. An additive error term with a variance function as presented in "Estimating parameters by regression" section was used.

## References

BMELV (2009) Bundeswaldinventur—Alle Ergebnisse und Berichte. http://www.bundeswaldinventur.de/ 01.05.2013

Gelman A, Carlin JB, Stern HS, Rubin DB (2004) Bayesian data analysis. Chapman and Hall, Boca Raton

Gilks W, Richardson S, Spiegelhalter D (eds) (1995) Markov chain Monte Carlo in practice. Springer, Berlin

Grundner S (1952) Massentafeln zur Bestimmung des Holzgehaltes stehender Waldbäume und Waldbestände, 10th edn. Verlag Paul Parey, Berlin

Hakkila P (1972) Utilisation of residual forest biomass. Springer, Berlin

Kollmann F (1982) Technologie des Holzes und der Holzwerkstoffe, vol 1. Springer, Berlin

Kublin E (2003) Einheitliche Beschreibung der Schaftform—Methoden und Programme—BDATPro. Forstwissenschaftl Centralbl 122(3):183–200

Metropolis N, Rosenbluth A, Rosenbluth M, Teller A, Teller E (1953) Equations of state calculations by fast computing machines. J Chem Phys 21:1087–1092

Muukkonen P (2007) Generalized allometric volume and biomass equations for some tree species in europe. Eur J For Res 126:157–166

Pinheiro JC, Bates DM (2000) Mixed-effects models in S and S-Plus. Springer, Berlin

Pretzsch H (2001) Modellierung des Waldwachstums. Parey, Berlin

R Development Core Team (2009) R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna. http://www.R-project.org, ISBN 3-900051-07-0

Saborowski J, Gaffrey D (1999) RBS, Ein mehrstufiges Inventurverfahren zur Schätzung von Baummerkmalen ii. Modifiziertes RBS-verfahren. AFJZ 170(12):223–227

Sprugel D (1983) Correcting for bias in log-transformed allometric equations. Ecology 64:209–210

Umweltbundesamt (2009) Submission under the United Nations Framework Convention on Climate Change and the Kyoto Protocol 2011. German greenhouse gas inventory 1990 - 2009. Tech. rep., Federal Environmental Agency. http://www.umweltbundesamt.de/publikationen/submission-under-united-nations-framework-0

Venables W, Ripley B (2002) Modern applied statistics with S. Statistics and Computing, Springer, Berlin

Wirth C, Schumacher J, Schulze ED (2004) Generic biomass functions for *Norway spruce* in central europe—a meta-analysis

approach towards prediction and uncertainty estimation. Tree Physiol 24:121–139

Wutzler T, Wirth C, Schumacher J (2008) Generic biomass functions for common beech (*European beech*) in central europe: predictions and components of uncertainty. Can J For Res 38(6):1661–1675

Zapata-Cuartas M, Sierra Ca, Alleman L (2012) Probability distribution of allometric coefficients and Bayesian estimation of aboveground tree biomass. For Ecol Manage 277:173–179 doi:10.1016/j.foreco.2012.04.030

Zianis D, Muukkonen P, Mäkipää R, Mencuccini M (2005) Biomass and stem volume equations for tree species in europe. Silva Fennica 4:63