

Article

Three-Stage Up-Scaling and Uncertainty Estimation in Forest Aboveground Biomass Based on Multi-Source Remote Sensing Data Considering Spatial Correlation

Xiangyuan Ding^{1,2,3}, Erxue Chen^{1,2,3,*}, Lei Zhao^{1,2,3}, Yaxiong Fan^{1,2,3}, Jian Wang^{1,2,3} and Yunmei Ma^{1,2,3}

- ¹ National Key Laboratory of Efficient Production of Forest Resources, Beijing 100091, China; dingxiangyuan@ifrit.ac.cn (X.D.); zhaolei@ifrit.ac.cn (L.Z.); fanyx@ifrit.ac.cn (Y.F.); wangjian@ifrit.ac.cn (J.W.); mayunmei@ifrit.ac.cn (Y.M.)
- ² Institute of Forest Resource Information Techniques, Chinese Academy of Forestry, Beijing 100091, China
- ³ Key Laboratory of Forestry Remote Sensing and Information System, National Forestry and Grassland Administration, Beijing 100091, China
- * Correspondence: chenerx@ifrit.ac.cn; Tel.: +86-010-62889164

Abstract: Airborne LiDAR (ALS) data have been extensively utilized for aboveground biomass (AGB) estimation; however, the high acquisition costs make it challenging to attain wall-to-wall estimation across large regions. Some studies have leveraged ALS data as intermediate variables to amplify sample sizes, thereby reducing costs and enhancing sample representativeness and model accuracy, but the cost issue remains in larger-scale estimations. Satellite LiDAR data, offering a broader dataset that can be acquired quickly with lower costs, can serve as an alternative intermediate variable for sample expansion. In this study, we employed a three-stage up-scaling approach to estimate forest AGB and introduced a method for quantifying estimation uncertainty. Based on the established three-stage general-hierarchical-model-based estimation inference (3sGHMB), an RK-3sGHMB inference method is proposed to make use of the regression-kriging (RK) method, and then it is compared with conventional model-based inference (CMB), general hierarchical model-based inference (GHMB), and improved general hierarchical model-based inference (RK-GHMB) to estimate forest AGB and uncertainty at both the pixel and forest farm levels. This study was carried out by integrating plot data, sampled ALS data, wall-to-wall Sentinel-2A data, and airborne P-SAR data. The results show that the accuracy of CMB ($R_{adj}^2 = 0.37$, $RMSE = 33.95$ t/ha, $EA = 63.28\%$) is lower than that of GHMB ($R_{adj}^2 = 0.38$, $RMSE = 33.72$ t/ha, $EA = 63.53\%$), while it is higher than that of 3sGHMB ($R_{adj}^2 = 0.27$, $RMSE = 36.58$ t/ha, $EA = 60.43\%$). Notably, RK-GHMB ($R_{adj}^2 = 0.60$, $RMSE = 27.07$ t/ha, $EA = 70.72\%$) and RK-3sGHMB ($R_{adj}^2 = 0.55$, $RMSE = 28.55$ t/ha, $EA = 69.13\%$) demonstrate significant accuracy enhancements compared to GHMB and 3sGHMB. For population AGB estimation, the precision of the proposed RK-3sGHMB ($p = 94.44\%$) is the highest, providing that there are sufficient sample sizes in the third stage, followed by RK-GHMB ($p = 93.32\%$) with sufficient sample sizes in the second stage, GHMB ($p = 90.88\%$), 3sGHMB ($p = 88.91\%$), and CMB ($p = 87.96\%$). Further analysis reveals that the three-stage model, considering spatial correlation at the third stage, can improve estimation accuracy, but the prerequisite is that the sample size in the third stage must be sufficient. For large-scale estimation, the RK-3sGHMB model proposed herein offers certain advantages.



Academic Editor: Eric Casella

Received: 6 January 2025

Revised: 11 February 2025

Accepted: 14 February 2025

Published: 16 February 2025

Citation: Ding, X.; Chen, E.; Zhao, L.; Fan, Y.; Wang, J.; Ma, Y. Three-Stage Up-Scaling and Uncertainty Estimation in Forest Aboveground Biomass Based on Multi-Source Remote Sensing Data Considering Spatial Correlation. *Remote Sens.* **2025**, *17*, 671. <https://doi.org/10.3390/rs17040671>

Copyright: © 2025 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Keywords: forest AGB; multi-source data; uncertainty; three-stage up-scaling; spatial correlation

1. Introduction

Forests are a pivotal component of the Earth's ecosystems, playing an essential role in ecological conservation, water source protection, air purification, and carbon cycling [1–3]. Aboveground biomass (AGB) is a critical metric in assessing forest health, and its rapid and precise estimation is imperative in the refined management of forest resources and research on carbon sequestration [4]. Traditionally, the acquisition of forest AGB has relied on design-based inference (DB), which uses plot-based data to estimate AGB at a regional level. However, this approach necessitates that the plot data adhere strictly to the principles of random or systematic sampling to ensure unbiased estimation. Therefore, this method faces significant challenges in terms of high labor and time costs when conducting large-scale surveys [5]. The advancement of remote sensing technology, particularly the utilization of multi-source remote sensing data, has opened up new possibilities in the large-scale inversion of forest AGB. Model-based inference has emerged as a widely adopted method of quantitative estimation in remote sensing, not only providing comprehensive and spatially continuous results but also significantly improving the population estimation precision under the same sample conditions [5,6]. Nonetheless, these methods are still contingent on sample size, and ensuring model accuracy requires a substantial number of samples, which perpetuates the issue of high labor and time costs. Multi-stage approaches have emerged as a research focus within multi-source remote sensing estimation technology, utilizing features strongly correlated with target variables to model and obtain a large number of relative true values, thereby expanding the original sample set, enhancing sample representativeness, and reducing the costs associated with obtaining measured samples [7–9].

In multi-stage modeling, intermediate variables play a crucial role in ensuring the precision of results. Light detection and ranging (LiDAR) data can accurately depict forest vertical structure information and have a strong correlation with forest resource parameters, such as average height, diameter at breast height (DBH), volume, and AGB, thus offering a distinct advantage in the estimation of forest resource parameters. Airborne LiDAR data (ALS), in particular, have been extensively applied in the quantitative inversion of forest resource parameters [10–12]. Due to the high cost of data acquisition, it is often challenging to achieve the quantitative estimation of forest resources across large areas using ALS data alone [13]. Therefore, ALS data are commonly employed as an intermediate variable in multi-stage modeling to expand the sample size [14–18]. Satellite LiDAR data, such as those from Global Ecosystem Dynamics Investigation (GEDI); Ice, Cloud, and Land Elevation Satellite-2 (ICESat-2); and Terrestrial Ecosystem Carbon Inventory Satellite (TECIS), are widely applied in multi-stage modeling, and they are also easier to obtain and involve relatively lower costs for end-users compared with ALS data [19–21]. Given the challenges involved in obtaining measured data for footprint locations due to positioning biases and inaccessibility, ALS data can be used as an intermediate variable to construct models that predict the relative true values of forest AGB at footprint locations [19,22–26]. This approach mitigates the impact of positioning errors and the difficulties associated with tracking footprints, meaning that the combination of sampled airborne LiDAR data and spatially discrete satellite waveform LiDAR data represents a common technical approach to large-area forest parameter estimation.

Currently, the data framework for one-stage inference, such as conventional model-based inference (CMB) [17], typically consists of a limited number of sample data and wall-to-wall data. The data framework for multi-stage inference includes the following three configurations: (1) a limited number of plot data combined with sampling ALS data and wall-to-wall data; (2) a limited number of plot data combined with satellite waveform LiDAR data and wall-to-wall data; (3) utilizing a limited number of plot data, sampling ALS

data, satellite waveform LiDAR data, and wall-to-wall data. In the multi-stage approach, except for the first stage, the dependent variables in subsequent stages are derived from the estimated values of the previous stage, which serve as relative true values. Consequently, error propagation can significantly influence the final result. Therefore, quantifying its uncertainty is crucial in practical applications. Scholars have conducted studies on this issue, employing models such as two-stage generalized hierarchical model-based inference (GHMB) [9] and three-stage generalized hierarchical model-based inference (3sGHMB) [8], but neither method accounts for the impact of spatial correlation. Geostatistical methods represent a widely recognized approach to mitigating the effects of spatial correlation and have been extensively employed in the estimation of forest resource parameters [27]. Several studies have integrated interpolation techniques into regression or machine learning algorithms, utilizing the spatial interpolation of residuals and then re-adding them to the original results to minimize the influence of spatial correlation on model outcomes. This integration has led to significant improvements in the accuracy of parameter estimations [28–31]. However, these methods are less commonly applied to multi-stage modeling, especially for uncertainty estimation. Zhao et al. [32] improved the GHMB model by incorporating spatial autocorrelation and introduced the RK-GHMB estimation method, which applies regression-kriging (RK) at the second stage of the modeling process. Their results demonstrate that integrating geostatistical methods can enhance the estimation accuracy of the model and make the uncertainty estimation more scientific and reasonable. To date, no studies have reported on the effectiveness of integrating geostatistical methods into three-stage inference. Additionally, there is a dearth of research comparing the advantages and disadvantages of three-stage inference relative to one- and two-stage inference [33].

Based on the aforementioned analysis, in this study, we introduce an RK model into the three-stage inference method 3sGHMB for forest AGB estimation, denoted as RK-3sGHMB, addressing the challenges arising from the fact that the current 3sGHMB does not account for spatial autocorrelation in the third stage. We then evaluate this method and compare it with CMB, GHMB, RK-GHMB, and 3sGHMB, analyzing their respective strengths and weaknesses, and assessing the effectiveness of the developed model. We also discuss the applicable scenarios of the developed model, with the goal of providing an efficient and high-precision multi-source remote sensing collaborative estimation technology for forest parameters, to support national forest resource annual monitoring services.

2. Materials and Methods

2.1. Study Area

The study area comprised the whole coverage area of two forest farms named Upper Yangge Forest Farm and Tidal Slag Forest Farm (Figure 1). The study area is located in Genhe City, Hulunbuir City, Inner Mongolia, where the terrain is undulating, with a relative elevation difference of between 100 m and 300 m. The average altitude is above 1000 m, characterizing it as a high-latitude, cold region. It has a cold temperate humid forest climate, with characteristics of a continental monsoon climate, featuring moist and cold air, long winters, short summers, and a rainy season from July to August each year. The vegetation is dominated by forest and grassland, with a forest coverage rate as high as 75%. Typical tree species include *Larix gmelinii*, *Betula platyphylla*, *Pinus sylvestris* var. *mongolica* Litv, and *Populus davidiana* [34].

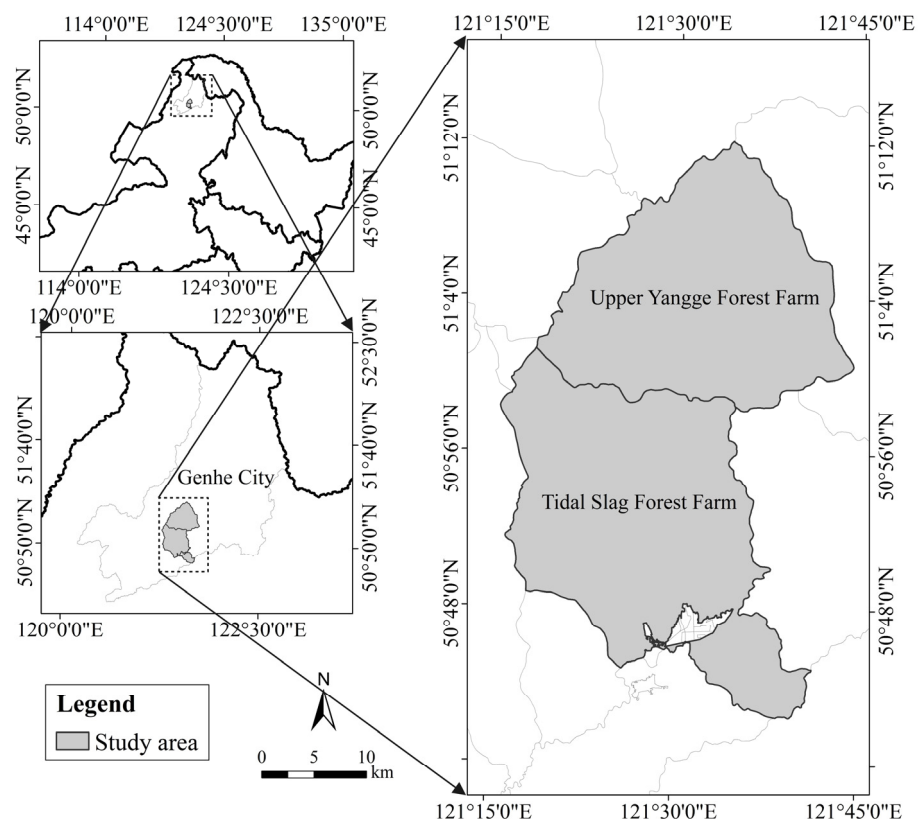


Figure 1. Location and coverage of the study area.

2.2. Data

2.2.1. Plot Data

The forest growth in the study area is relatively slow. To increase the sample size, data from two years, 2021 and 2022, were used for subsequent analysis. A total of 103 plots (25 m × 25 m) were obtained over the two years (Figure 2). Trees with a DBH greater than 5 cm within the plots were measured, and parameters such as DBH, tree height, east–west crown width, north–south crown width, and height to the base of the live crown were measured. The positions of the four corners and the center of each plot were precisely measured using a FindCM GPS-RTK produced by China Qianxun SI company with a measurement accuracy better than 0.15 m. Subsequently, the tree species-specific allometric growth equations proposed by Zhou et al. [35] were used to calculate the AGB of each tree within the plots. The total biomass and AGB density for each plot were calculated from the AGB of individual tree within the respective plots. The unit of total biomass is ton (t), and the unit of biomass density is tons per hectare (t/ha).

2.2.2. ALS Strip Data Scenario Simulation

Wall-to-wall ALS data were collected in August 2022 (Figure 2) using the CAF-LiCHy data acquisition system, with the LiDAR sensor model being Riegl LMS-Q680i produced by RIEGL that in Horn, Austria [36]. The LiDAR point cloud data underwent preprocessing, including attitude correction, noise point removal, coordinate transformation, flight line stitching, and system error correction. Subsequently, LASTOOLS software [37] was utilized for point cloud classification, denoising, and normalization to obtain normalized vegetation point cloud data. Fusion software [38] was employed to extract features such as the mean height, canopy density, height variance and standard deviation, coefficient of variation, interquartile distance, skewness, maximum height, and height percentiles for subsequent analysis [39,40]. Detailed information is presented in Table 1. The study area was divided

into 1,539,092 grids ($25\text{ m} \times 25\text{ m}$) denoted as N , and the mean of each LiDAR data feature within each grid was calculated for further analysis. In order to simulate real-world application scenarios, 13 strips were extracted to serve as a link to GEDI data (Figure 2b), with both the strip width and spacing set at 1.5 km. There were 80 plots located within the strips and 23 plots outside the strips. Plot data within the strips were used for model training, while plot data outside the strips were used for model validation.

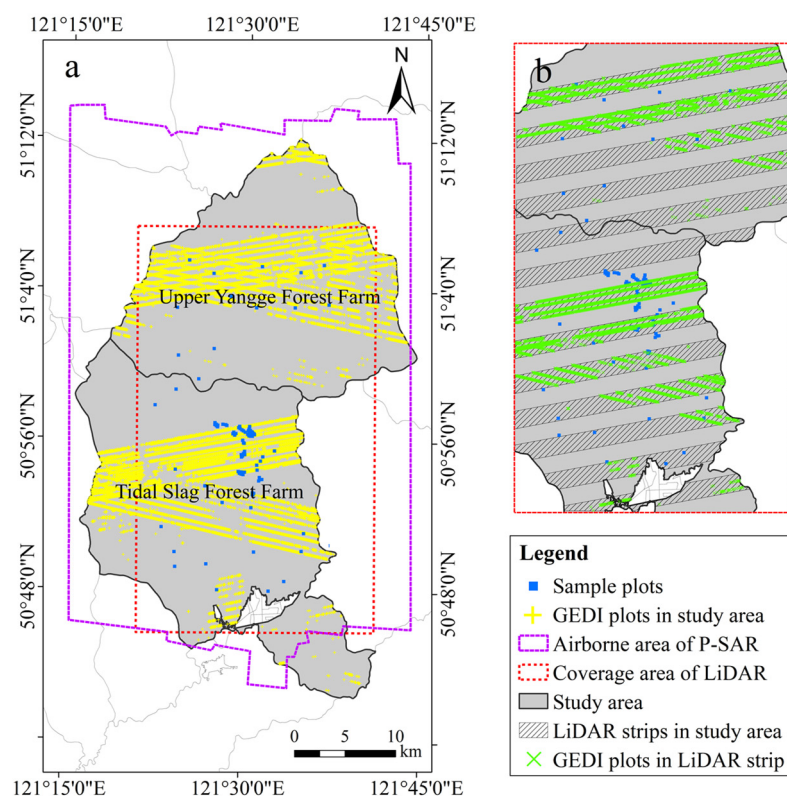


Figure 2. (a) Spatial distribution of sample plots, data from Sentinel-2A, P-SAR, and GEDI in the study area; and (b) the spatial distribution of ALS strips and GEDI data within the strips.

Table 1. Features of LiDAR.

Feature Name	Feature Symbol	Description
Mean height	Hmean	The mean height of the point cloud within a $25\text{ m} \times 25\text{ m}$ statistical unit
Forest canopy density	CD	The ratio of canopy backscatter points within a $25\text{ m} \times 25\text{ m}$ statistical unit to the total number of backscatter points
Variance and standard deviation	Hvar, Hstd	The variance and standard deviation of the point cloud within a $25\text{ m} \times 25\text{ m}$ statistical unit
Coefficient of variation	Hcv	Coefficient of variation in the point cloud within a $25\text{ m} \times 25\text{ m}$ statistical unit
Interquartile distance	Hint	Height interquartile distance within a $25\text{ m} \times 25\text{ m}$ statistical unit for point cloud data
Skewness	Hskw	The skewness of the points cloud within a $25\text{ m} \times 25\text{ m}$ statistical unit
Maximum and minimum	Hmax, Hmin	The maximum and minimum values of the point cloud within a $25\text{ m} \times 25\text{ m}$ statistical unit
Percentiles	H10, H20, H30, H40, H50, H60, ..., H95	Percentiles of the point cloud at different heights within a $25\text{ m} \times 25\text{ m}$ statistical unit

2.2.3. Airborne P-SAR

In this study, we obtained wall-to-wall Airborne P-SAR data for the study area, employing processing methodologies in accordance with those detailed by Fan et al. [41]. The primary steps encompass (1) multi-looking and refined Lee polarization filtering to generate multi-look complex (MLC) PolSAR data, followed by the application of a 5×5 moving window filter to the MLC results to further attenuate the impact of speckle noise; (2) the integration of SAR imaging orbit information into digital elevation model (DEM) data and the geocoding of the PolSAR data based on the radar–Doppler geolocation model to establish the correspondence between the SAR image slant range space and the geographic coordinate space while also obtaining relevant angular parameters that describe the local imaging geometry, such as projection angles and local incidence angles; (3) the implementation of a three-stage topographic–radiometric correction method [42] to address the polarization azimuth angle, effective scattering area, and angle effects, in order to compensate for the effects of terrain undulations on polarization characteristics; (4) conducting Yamaguchi decomposition [43] on the topographic–radiometric corrected polarization coherence matrix to extract various polarization decomposition components that, along with backscatter intensity, form the polarization feature set for each band (Table 2). The final processed multi-band polarimetric SAR data and features have a pixel resolution of $25 \text{ m} \times 25 \text{ m}$, ensuring spatial congruence with the sample plots.

Table 2. Polarimetric features extracted from P-SAR data in this study.

Feature Class	Feature Name	Feature Symbol
Backscatter intensity	HH polarization	PHH
	HV polarization	PHV
Yamaguchi decomposition	Surface scattering	Odd
	Dihedral scattering	Db1
	Volume scattering	Vol
	Helix scattering	Hlx

2.2.4. GEDI Data

The GEDI LiDAR system, launched aboard the International Space Station (ISS) in December 2018, comprises three lasers, with one laser divided into two coverage beams and the other two operating at full power. The coverage beams and full-power beams are able to detect the ground through 95% and 98% forest canopy cover, respectively [44]. GEDI creates spot trajectories on the Earth’s surface through beam jitter, with each trajectory consisting of circular footprints that are 60 m along-track and 25 m in diameter, and the distance between adjacent trajectories is approximately 600 m, with a scanning swath width of about 4.2 km [45,46]. These data have been widely applied to the estimation of AGB in forests [47,48]. The data used in this study are cover, canopy height, and height quantiles, i.e., rh10, rh20, rh30, rh40, rh50, rh60, rh70, rh75, rh80, rh85, rh90, rh95, and rh100 from the GEDI Level 2A and Level 2B data products. The data acquisition period was from 1 July 2022 to 30 September 2022, and some invalid spots were filtered based on the parameters inherent in the data [49]. The data were downloaded and preprocessed using the rGEDI package in R [50], with outliers removed, resulting in data features from 13,010 location spots, as shown in Figure 2a, of which 5715 GEDI spots are located within the strips (Figure 2b).

2.2.5. Sentinel-2A Data

Two scenes of sentinel-2A multispectral data were downloaded from the Copernicus Open Access Hub website (<https://dataspace.copernicus.eu/> (accessed on 10 May 2023)).

Due to the influence of cloudy and wet weather, it was not possible to obtain imagery that is perfectly synchronized with the ground plot data acquisition time; therefore, only available imagery from the growing season with a close time match was acquired, dated 8 July 2022. The European Space Agency’s open-source SNAP software (<https://step.esa.int/main/download/snap-download/> (Version 8.0.0.0, accessed on 20 May 2023)) was utilized to perform radiometric calibration, atmospheric correction, and orthorectification on each scene. After these preprocessing steps had been performed, the pixel size of the Sentinel-2A data were set to 25 m × 25 m, consistent with the size of the sample plots; then, the two preprocessed scenes were mosaicked to obtain multispectral remote sensing data covering the entire study area. Common vegetation indices were selected as candidate features (Table 3), and an additional novel vegetation index, the kernel-normalized difference vegetation index (KNDVI), was included. This index is derived from the concept of kernel functions, and previous studies have shown that it can significantly improve the estimation accuracy of forest carbon stocks [51]. The corresponding kernel function vegetation indices, KNDVIre1, KNDVIre2, and KNDVIre3, were calculated with the methods shown in [51].

Table 3. Features of Sentinel-2A.

Feature	Symbol	Formula	Reference
Spectral reflectance	B2, B3, B4, B5, B6, B7, B8a, B11, B12	/	/
Normalized difference vegetation index	NDVI	$(B_8 - B_4) / (B_8 + B_4)$	[52]
Red-edge vegetation index	NDVIre1 NDVIre2 NDVIre3	$(B_8 - B_5) / (B_8 + B_5)$ $(B_8 - B_6) / (B_8 + B_6)$ $(B_8 - B_7) / (B_8 + B_7)$	[53]
Kernel-normalized difference vegetation index	KNDVI	$\tanh\left(\left(\frac{B_8 - B_4}{2\sigma_1}\right)^2\right)$	[51]
Kernel red-edge vegetation index	KNDVIre1	$\tanh\left(\left(\frac{B_8 - B_5}{2\sigma_2}\right)^2\right)$	
	KNDVIre2	$\tanh\left(\left(\frac{B_8 - B_6}{2\sigma_3}\right)^2\right)$	
	KNDVIre3	$\tanh\left(\left(\frac{B_8 - B_7}{2\sigma_4}\right)^2\right)$	
Difference vegetation index	DVI	$B_8 - B_4$	[54]
Red-edge difference vegetation index	DVIre1 DVIre2 DVIre3	$B_5 - B_4$ $B_6 - B_4$ $B_7 - B_4$	[55]
Enhanced vegetation index	EVI	$2.5 \times (B_8 - B_4) / (B_8 + 6 \times B_4 - 7.5 \times B_2 + 1)$	[56]
Red-edge enhanced vegetation index	EVIre1 EVIre2 EVIre3	$2.5 \times (B_5 - B_4) / (B_5 + 6 \times B_4 - 7.5 \times B_2 + 1)$ $2.5 \times (B_6 - B_4) / (B_6 + 6 \times B_4 - 7.5 \times B_2 + 1)$ $2.5 \times (B_7 - B_4) / (B_7 + 6 \times B_4 - 7.5 \times B_2 + 1)$	[55]
Simple ratio	RVI	B_8 / B_4	[57]
Soil-adjusted vegetation index	SAVI	$1.5 \times (B_8 - B_4) / (B_8 + B_4 + 0.5)$	[58]
Green normalized difference vegetation index	GNDVI	$(B_8 - B_3) / (B_8 + B_3)$	[59]
Narrow enhanced vegetation index	EVInirn	$2.5 \times (B_{8A} - B_4) / (B_{8A} + 6 \times B_4 - 7.5 \times B_2 + 1)$	[55]

$B_2, B_3, B_4, B_5, B_6, B_7, B_8, B_{8a}, B_{11},$ and B_{12} represent the bands corresponding to sentinel-2A; the value of $\sigma_1, \sigma_2, \sigma_3$ and σ_4 is $(B_8 + B_4)/2, (B_8 + B_5)/2, (B_8 + B_6)/2,$ and $(B_8 + B_7)/2,$ respectively.

2.3. Methods

2.3.1. Overview

We propose the application of RK at the third stage of a three-stage model, and in this study, we aim to confirm whether it can promote accuracy and reduce uncertainty in estimation results. We compare it with CMB, GHMB, RK-GHMB, and 3sGHMB to analyze their strengths and weaknesses in practical application scenarios. The flowchart for this study is depicted in Figure 3, and the data frameworks for the five methods are as follows:

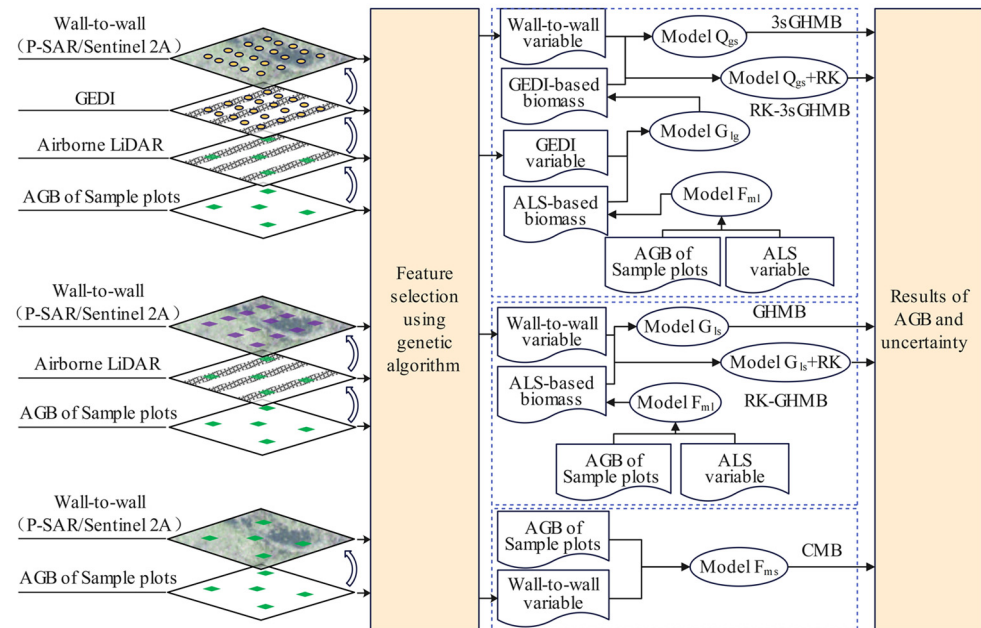


Figure 3. Overall workflow of the study; the green square blocks represent field sample plots; the purple square blocks represent ALS sample plots; the yellow circles represent the GEDI sample plots; the black bars represent the ALS strips; the elements enclosed within the blue boxes collectively constitute a unified whole.

Case A: CMB using plot data and wall-to-wall airborne P-SAR and Sentinel-2A data.

Case B: GHMB using plot data, sampled ALS strip data, and wall-to-wall airborne P-SAR and Sentinel-2A data.

Case C: RK-GHMB using the same data as in Case B.

Case D: 3sGHMB using plot data, sampled ALS strip data, GEDI data, and wall-to-wall airborne P-SAR and Sentinel-2A data.

Case E: RK-3sGHMB using the same data as in Case D.

In this study, the N grid of the study area is denoted as population U . The grids for Tidal Slag Forest Farm constitute the subpopulation U_1 , with a total number of 821,909, denoted as N_1 . The grids for Upper Yangge Forest Farm constitute the subpopulation U_2 , with a total number of 717,183, denoted as N_2 . For the sake of clarity in subsequent discussions, the model within the CMB is designated as F_{ms} ; the first-stage models for the GHMB, RK-GHMB, 3sGHMB, and RK-3sGHMBs are represented by F_{ml} ; the second-stage models for GHMB and RK-GHMB are denoted by G_{ls} and G_R (in Figure 3, G_R is denoted as “model $G_{ls} + RK$ ”), respectively; the second-stage models for 3sGHMB and RK-3sGHMB are denoted by Q_{lg} ; and their third-stage models are labeled as Q_{gs} and Q_R (in Figure 3, Q_R is denoted as “model $Q_{gs} + RK$ ”). The dataset constituted by the sample plots is represented by S_I , which is located in the airborne LiDAR strips and includes the measured forest AGB, ALS features, and wall-to-wall features. The sample size n for S_I is 80. The sample plots outside the ALS strips are used for validation, with a count of 23. From the GEDI spots within the ALS strips, a random selection is made to construct the S_{II} dataset, which includes ALS features, GEDI features, and wall-to-wall features, and the sample size k is 5000. Since there are no measured data at the GEDI spot locations, data for an additional GEDI spot for which LiDAR-based biomass is available are randomly selected to form a validation dataset to verify the G_{lg} model; the sample size of the dataset is 1000. From all GEDI spots, random selection is performed to create the S_{III} dataset for model training, which includes GEDI features and wall-to-wall features. The sample size of S_{III} is $h = 10,000$.

2.3.2. Feature Selection

We employed a genetic algorithm for feature selection, which is a heuristic method that simulates natural selection to solve optimization problems and has a wide range of uses in the field of remote sensing [60]. The R^2 and $RMSE$ of internal and external validation samples were used as evaluation metrics. We implemented genetic algorithm feature selection for each model, utilizing the `gafs` function within the `caret` package in `R` [61].

2.3.3. Case A: The CMB

This case follows the well-established model-based inference [9,17]. The general form of the model is as follows:

$$F_{ms} : \mathbf{y} = f(\mathbf{X}_{S_I}, \boldsymbol{\alpha}_{S_I}) + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Omega}) \quad (1)$$

where \mathbf{y} is the vector of measured forest AGB data, \mathbf{X}_{S_I} is an $n \times q$ dimension matrix of wall-to-wall features with respect to S_I , $\boldsymbol{\alpha}_{S_I}$ is the corresponding parameter vector, and $\boldsymbol{\varepsilon}$ is the random error following $N(0, \boldsymbol{\Omega})$. f represents the model form for fitting. In this study, we used logistic regression for model fitting and estimation due to its ability to ensure non-negative values and provide a certain physical significance [15]. The same model form was also used in Case B, Case C, Case D, and Case E. The pixel uncertainty is calculated based on the methods proposed by Saarela et al. [9] and McRoberts et al. [62], using the following formula:

$$RMSE(\hat{y}_{F_{ms_i}}) = \sqrt{\tilde{\mathbf{X}}_i^T Cov(\hat{\boldsymbol{\alpha}}_{S_I}) \tilde{\mathbf{X}}_i + V(\varepsilon_i)} \quad (2)$$

$$\tilde{\mathbf{X}}_i = \frac{\partial F_{ms}(\hat{\boldsymbol{\alpha}}_{S_I}, \mathbf{X}_i)}{\partial \hat{\boldsymbol{\alpha}}_{S_I}} \quad (3)$$

where $RMSE(\hat{y}_{F_{ms_i}})$ represents the uncertainty of the model F_{ms} at map unit i ; $\hat{y}_{F_{ms_i}}$ is the predicted AGB using model F_{ms} for map unit i ; $\tilde{\mathbf{X}}_i$ is a $(q + 1)$ -length vector composed of the partial derivatives of the model F_{ms} at unit i with respect to $\hat{\boldsymbol{\alpha}}_{S_I}$, where q is the number of parameters. $Cov(\boldsymbol{\alpha}_{S_I}) = (\tilde{\mathbf{X}}_{S_I}^T \boldsymbol{\Omega} \tilde{\mathbf{X}}_{S_I})^{-1}$, and $\tilde{\mathbf{X}}_{S_I}$ is an $n \times q$ matrix of partial derivatives of the model F_{ms} based on the dataset S_I with respect to $\hat{\boldsymbol{\alpha}}_{S_I}$, and $\boldsymbol{\Omega}$ is the variance-covariance matrix of the residuals. Without considering spatial autocorrelation, the diagonal elements are constituted by $V(\varepsilon_i)$, which is the variance of ε_i . It can be calculated by the method proposed by McRoberts [62], with the formula $V(\varepsilon_i) = a \bar{y}_{S_I}^b + \delta$, where a and b are the model parameters, δ is the residuals term, \bar{y}_{S_I} and is the mean of the predicted forest AGB for each group after stratification.

The population mean is as follows:

$$\mu_{F_{ms}} = \frac{1}{N} \sum_{i=1}^N y_{F_{ms_i}} \quad (4)$$

with the variance of $\mu_{F_{ms}}$ being

$$Var(\mu_{F_{ms}}) = \mathbf{I}_U^T \tilde{\mathbf{X}}_U Cov(\hat{\boldsymbol{\alpha}}_{S_I}) \tilde{\mathbf{X}}_U^T \mathbf{I}_U \quad (5)$$

where \mathbf{I}_U is an N -length vector with each element being $1/N$, and $\tilde{\mathbf{X}}_U$ is a partial derivative matrix of F_{ms} with respect to the population U . When estimating the population mean and variance for the two sub-regions, U and N are replaced by U_1, N_1 and U_2, N_2 .

2.3.4. Case B: The GHMB

In the case of GHMB, the forest AGB values estimated by sampling ALS data serve as extended samples for the estimation of the full-coverage forest AGB [9]. Similarly to Equation (1), the general form of the model for the first stage is as follows:

$$F_{ml} : \mathbf{y} = f(\mathbf{P}_{S_I}, \boldsymbol{\beta}_{S_I}) + \mathbf{o}, \mathbf{o} \sim N(0, \boldsymbol{\Omega}^*) \tag{6}$$

where \mathbf{P}_{S_I} is an $n \times p$ dimensional matrix of ALS features with respect to S_I ; $\boldsymbol{\beta}_{S_I}$ is the corresponding model parameter vector; and \mathbf{o} is the random error following $N(0, \boldsymbol{\Omega})$.

The second model general form is as follows:

$$G_{Is} : \mathbf{y}_{F_{ml}} = f(\mathbf{X}_{S_{II}}^*, \boldsymbol{\alpha}_{S_{II}}) + \mathbf{v}, \mathbf{v} \sim N(0, \mathbf{C}) \tag{7}$$

where $\mathbf{y}_{F_{ml}}$ is the predicted AGB estimated by the model F_{ml} , $\mathbf{X}_{S_{II}}^*$ is a $k \times g$ matrix of wall-to-wall data features with respect to S_{II} ; $\boldsymbol{\alpha}_{S_{II}}$ is a vector of the corresponding model parameters; and \mathbf{v} is a vector of random error following $N(0, \mathbf{C})$.

The uncertainty is calculated using the same method as Equation (2):

$$RMSE(\hat{\mathbf{y}}_{G_{Is_i}}) = \sqrt{\tilde{\mathbf{X}}_i^{*T} Cov(\hat{\boldsymbol{\alpha}}_{S_{II}}) \tilde{\mathbf{X}}_i^* + V(v_i)} \tag{8}$$

$$\tilde{\mathbf{X}}_i^* = \frac{\partial G_{Is}(\hat{\boldsymbol{\alpha}}_{S_{II}}, \mathbf{X}_i^*)}{\partial \hat{\boldsymbol{\alpha}}_{S_{II}}} \tag{9}$$

where $\hat{\mathbf{y}}_{G_{Is_i}}$ is the predicted AGB using G_{Is} at the unit i ; $\tilde{\mathbf{X}}_i^*$ is a $(g + 1)$ -length vector of partial derivatives of the G_{Is} model with respect to $\hat{\boldsymbol{\alpha}}_{S_{II}}$; $Cov(\hat{\boldsymbol{\alpha}}_{S_{II}})$ is the estimated covariance matrix for $\hat{\boldsymbol{\alpha}}_{S_{II}}$; $V(v_i)$ is the variance of random error v_i at unit i , with the calculation approach being identical to that of $V(\varepsilon_i)$ and the model form being $V(v_i) = a\hat{\mathbf{y}}_{S_{II}}^2 + b\hat{\mathbf{y}}_{S_{II}} + c + \delta_2$ with a , b , and c as corresponding parameters; and δ_2 is the random error. According to Saarela et al. [9], $Cov(\hat{\boldsymbol{\alpha}}_{S_{II}})$ in the two-stage method can be expressed as follows:

$$Cov(\hat{\boldsymbol{\alpha}}_{S_{II}}) = (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} + \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} Cov(\mathbf{y}_{F_{mlS_{II}}}) \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^* (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} \tag{10}$$

where $\tilde{\mathbf{X}}_{S_{II}}^*$ is the $k \times g$ matrix of partial derivatives of the model F_{ms} based on dataset S_{II} with respect to $\hat{\boldsymbol{\alpha}}_{S_{II}}$. $Cov(\mathbf{y}_{F_{mlS_{II}}}) = \tilde{\mathbf{P}}_{S_{II}} Cov(\hat{\boldsymbol{\beta}}_{S_I}) \tilde{\mathbf{P}}_{S_{II}}^T$, $\mathbf{y}_{F_{mlS_{II}}}$ is the predicted AGB using the F_{ml} model with respect to S_{II} ; $\tilde{\mathbf{P}}_{S_{II}}$ is the $k \times p$ matrix of partial derivatives of the model F_{ml} based on dataset S_{II} with respect to $\hat{\boldsymbol{\beta}}_{S_I}$; $Cov(\hat{\boldsymbol{\beta}}_{S_I}) = (\tilde{\mathbf{P}}_{S_I}^T \boldsymbol{\Omega}^* \tilde{\mathbf{P}}_{S_I})^{-1}$, with $\tilde{\mathbf{P}}_{S_I}$ being the $n \times p$ matrix of partial derivatives of the model F_{ml} based on dataset S_I with respect to $\hat{\boldsymbol{\beta}}_{S_I}$.

The population mean estimated by GHMB can be written as follows:

$$\mu_{G_{Is}} = \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{y}}_{G_{Is_i}} \tag{11}$$

The corresponding variance of $\mu_{G_{Is}}$ is

$$Var(\mu_{G_{Is}}) = \mathbf{l}_U^T \tilde{\mathbf{X}}_U^* [(\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} + \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} Cov(\mathbf{y}_{F_{mlS_{II}}}) \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^* (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1}] \tilde{\mathbf{X}}_U^T \mathbf{l}_U \tag{12}$$

where $\tilde{\mathbf{X}}_U^*$ is a partial derivative matrix of G_{Is} with respect to the population U . When estimating the population mean and variance for the two sub-regions, U and N are replaced with U_1, N_1 and U_2, N_2 .

2.3.5. Case C: The RK-GHMB

Zhao et al. [27] proposed an improved GHMB model named RK-GHMB, which uses an RK model in the second stage of GHMB, considering the spatial correlation of residual errors. The general form of the first stage is identical to Equation (6). Assuming that the residual, which cannot be predicted using the trend model of (1), is second-order stationarity, the model for the second stage is as follows:

$$G_R : \hat{y}_{G_{ls}-RK_i} = \hat{y}_{G_{ls_i}} + V_{kriging} \tag{13}$$

where $\hat{y}_{G_{ls}-RK_i}$ is the predicted AGB of the model G_R , and $V_{kriging}$ is the residual result interpolated by ordinary kriging, which can be calculated as follows:

$$V_{kriging} = \sum_{j=1}^k d_j v_j, \quad \sum_{j=1}^k d_j = 1 \tag{14}$$

where d_j is the kriging weight of the residual v_j at the j th index; the total weights must add up to one. The d_j can be calculated as follows:

$$\mathbf{d} = \mathbf{C}^{*-1} \cdot \mathbf{c}_{G_{ls_i}} \tag{15}$$

$$\mathbf{c}_{G_{ls_i}} = \begin{bmatrix} Cov(v_i, v_1) \\ Cov(v_i, v_2) \\ \vdots \\ Cov(v_i, v_k) \end{bmatrix} \tag{16}$$

where \mathbf{d} is an m -length vector of d_j kriging weights; \mathbf{C}^* is a $k \times k$ covariance matrix of residuals for S_{II} , with its diagonal elements constituted by $V(v_i)$ and the off-diagonal elements calculated using $V(v_i)$, $V(v_j)$, and the spatial correlation ρ_{ij} ; i.e., $Cov(v_i, v_j) = \sqrt{V(v_i) \cdot V(v_j)} \cdot \rho_{ij}$. ρ_{ij} is expressed as $\rho_{ij} = 1 - \frac{\gamma(d_{ij})}{C_0 + C_1}$, where $\gamma(d_{ij})$ is a semi-variogram, and its form is either exponential, gaussian, spherical, or linear. C_0 and C_1 are the nugget and partial sill, respectively. $\mathbf{c}_{G_{ls_i}}$ is a k -length vector of covariance for the unit map i .

Thus, the uncertainty of RK-GHMB can be written as follows:

$$RMSE(y_{G_{ls}-RK_i}) = \sqrt{\begin{bmatrix} (\tilde{\mathbf{X}}_i^* - \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \mathbf{c}_{G_{ls_i}})^T \\ \left\{ (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} + \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} Cov(\mathbf{y}_{F_{mlS_{II}}}) \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^* (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} \right\} \\ (\tilde{\mathbf{X}}_i^* - \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \mathbf{c}_{G_{ls_i}}) + V(v_i) - \mathbf{c}_{G_{ls_i}}^T \mathbf{C}^{*-1} \mathbf{c}_{G_{ls_i}} \end{bmatrix}} \tag{17}$$

The population mean estimated by RK-GHMB is

$$\mu_{G_{ls}-RK_i} = \frac{1}{N} \sum_{i=1}^N \hat{y}_{G_{ls}-RK_i} \tag{18}$$

with the variance estimator as

$$Var(\mu_{G_{ls}-RK_i}) = \left\{ \begin{bmatrix} l_U^T \left\{ (\tilde{\mathbf{X}}_U^* - \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \mathbf{c}_{G_{ls_U}}) \right. \\ \left. \left\{ (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} + \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} Cov(\mathbf{y}_{F_{mlS_{II}}}) \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^* (\tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \tilde{\mathbf{X}}_{S_{II}}^*)^{-1} \right\} \right. \\ \left. \left. (\tilde{\mathbf{X}}_U^* - \tilde{\mathbf{X}}_{S_{II}}^{*T} \mathbf{C}^{*-1} \mathbf{c}_{G_{ls_U}}) \right\} l_U \right\} \tag{19}$$

where the parameters are the same as in Equation (12).

2.3.6. Case D: The 3sGHMB

The 3sGHMB model includes three stages [8]. The general form of the first stage is identical to the GHMB, and the form of the second stage is as follows:

$$Q_{lg} : y_{F_{ml}} = f(\mathbf{Z}_{S_{II}}, \boldsymbol{\eta}_{S_{II}}) + \mathbf{e}, \mathbf{e} \sim N(0, \boldsymbol{\Sigma}) \tag{20}$$

where $\mathbf{Z}_{S_{II}}$ is a $k \times l$ dimension matrix of GEDI features with respect to S_{II} ; $\boldsymbol{\eta}_{S_{II}}$ is the corresponding model parameter vector; and \mathbf{e} is the random error following $N(0, \boldsymbol{\Sigma})$. The general form of the third stage can be written as follows:

$$Q_{gs} : y_{G_{ls}} = f(\mathbf{X}_{S_{III}}^{**}, \boldsymbol{\alpha}_{S_{III}}) + \boldsymbol{\theta}, \boldsymbol{\theta} \sim N(0, \boldsymbol{\Lambda}) \tag{21}$$

where $\mathbf{X}_{S_{III}}^{**}$ is an $h \times r$ dimension matrix of wall-to-wall features with respect to S_{III} ; $\hat{\boldsymbol{\alpha}}_{S_{III}}$ is the corresponding model parameter vector; and $\boldsymbol{\theta}$ is the random error following $N(0, \boldsymbol{\Lambda})$.

The estimation of uncertainty is analogous to Equation (8):

$$RMSE(y_{Q_{gs_i}}) = \sqrt{\tilde{\mathbf{X}}_i^{**T} Cov(\hat{\boldsymbol{\alpha}}_{S_{III}}) \tilde{\mathbf{X}}_i^{**} + V(\theta_i)} \tag{22}$$

$$\tilde{\mathbf{X}}_i^{**} = \frac{\partial Q(\hat{\boldsymbol{\alpha}}_{S_{III}}, \mathbf{X}_i^{**})}{\partial \hat{\boldsymbol{\alpha}}_{S_{III}}} \tag{23}$$

where $y_{Q_{gs_i}}$ is the predicted AGB using Q_{gs} at the map unit i ; $\tilde{\mathbf{X}}_i^{**}$ is a $(r + 1)$ -length vector of partial derivatives of the Q_{gs} model with respect to $\hat{\boldsymbol{\alpha}}_{S_{III}}$; $Cov(\hat{\boldsymbol{\alpha}}_{S_{III}})$ is the estimated covariance matrix for $\hat{\boldsymbol{\alpha}}_{S_{III}}$; $V(\theta_i)$ is the variance of random error θ_i at the map unit i , with a calculation approach identical to that of $V(\varepsilon_i)$ and the model form $V(\theta_i) = a\tilde{y}_{S_{IIIi}}^b + \delta_4$ with a and b as corresponding parameters; and δ_4 is the random error. Analogous to $Cov(\hat{\boldsymbol{\alpha}}_{S_{II}})$, $Cov(\hat{\boldsymbol{\alpha}}_{S_{III}})$ can be expressed as follows:

$$Cov(\hat{\boldsymbol{\alpha}}_{S_{III}}) = (\tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} + \tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} Cov(\mathbf{y}_{Q_{lgS_{III}}}) \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} \tag{24}$$

where $\tilde{\mathbf{X}}_{S_{III}}^{**}$ is an $h \times r$ matrix of partial derivatives of the model Q_{gs} based on dataset S_{III} with respect to $\hat{\boldsymbol{\alpha}}_{S_{III}}$. $\boldsymbol{\Lambda}$ is a diagonal matrix; each element is estimated using the model form presented as $V(\theta_i)$. The covariance of $\mathbf{y}_{Q_{lgS_{III}}}$ can be expressed using $\tilde{\mathbf{Z}}_{S_{III}}$ and $\hat{\boldsymbol{\eta}}_{S_{II}}$; i.e., $Cov(\mathbf{y}_{Q_{lgS_{III}}}) = \tilde{\mathbf{Z}}_{S_{III}} Cov(\hat{\boldsymbol{\eta}}_{S_{II}}) \tilde{\mathbf{Z}}_{S_{III}}^T$. $\tilde{\mathbf{Z}}_{S_{III}}$ is an $h \times r$ matrix of partial derivatives of the model Q_{lg} based on dataset S_{III} with respect to $\hat{\boldsymbol{\eta}}_{S_{II}}$. The covariance matrix of $\hat{\boldsymbol{\eta}}_{S_{II}}$ was estimated as follows:

$$Cov(\hat{\boldsymbol{\eta}}_{S_{II}}) = (\tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} + \tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} Cov(\mathbf{y}_{F_{mlS_{II}}}) \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}} (\tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tag{25}$$

where $\tilde{\mathbf{Z}}_{S_{II}}$ is a $k \times r$ matrix of partial derivatives of the model Q_{lg} based on dataset S_{II} with respect to $\hat{\boldsymbol{\eta}}_{S_{II}}$; the calculation method employed in $Cov(\mathbf{y}_{F_{mlS_{II}}})$ is consistent with the approach utilized in Equation (10). By replacing $Cov(\mathbf{y}_{F_{mlS_{II}}})$ and $Cov(\hat{\boldsymbol{\eta}}_{S_{II}})$ with its estimators, the covariance matrix estimator of $\hat{\boldsymbol{\alpha}}_{S_{III}}$ can be expressed as

$$Cov(\hat{\boldsymbol{\alpha}}_{S_{III}}) = \begin{aligned} & (\tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} + \\ & \tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{Z}}_{S_{III}} (\tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} + \\ & \tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{Z}}_{S_{III}} \tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \hat{\mathbf{Z}}_{S_{II}} (\hat{\mathbf{P}}_{S_{II}}^T \boldsymbol{\Omega} \hat{\mathbf{P}}_{S_{II}})^{-1} \hat{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{III}} \\ & (\tilde{\mathbf{Z}}_{S_{II}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \boldsymbol{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} \end{aligned} \tag{26}$$

The population mean estimated by 3sGHMB is

$$\mu_{Q_{gs}} = \frac{1}{N} \sum_{i=1}^N \hat{y}_{Q_{gs}} \tag{27}$$

with the variance estimator being

$$\begin{aligned} Var(\mu_{Q_{gs}}) = & \tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{Z}}_{S_{III}} (\tilde{\mathbf{Z}}_{S_{II}}^T \mathbf{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} + \\ & \tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{Z}}_{S_{III}} (\tilde{\mathbf{Z}}_{S_{II}}^T \mathbf{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} + \\ & \tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{Z}}_{S_{III}} (\tilde{\mathbf{Z}}_{S_{II}}^T \mathbf{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} \tilde{\mathbf{Z}}_{S_{II}}^T \mathbf{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}} \\ & (\tilde{\mathbf{Z}}_{S_{II}}^T \mathbf{\Sigma}^{-1} \tilde{\mathbf{Z}}_{S_{II}})^{-1} \tilde{\mathbf{Z}}_{S_{III}}^T \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**} (\tilde{\mathbf{X}}_{S_{III}}^{**T} \mathbf{\Lambda}^{-1} \tilde{\mathbf{X}}_{S_{III}}^{**})^{-1} \tilde{\mathbf{X}}_U^{**T} \mathbf{l}_U \end{aligned} \tag{28}$$

where $\tilde{\mathbf{X}}_U^{**}$ is a partial derivative matrix of Q_{gs} with respect for the population U .

2.3.7. Case E: The Proposed RK-3sGHMB

Similarly to Equation (13), we replaced the third general form Q_{gs} as follows:

$$Q_R : y_{Q_{gs}-RK_i} = y_{Q_{gs_i}} + V_{kriging}^* \tag{29}$$

where $y_{Q_{gs_i}}$ is the predicted AGB of the model Q_{gs} , and $V_{kriging}^*$ is the residual result interpolated by ordinary kriging, which can be calculated as follows:

$$V_{kriging}^* = \sum_{j=1}^h w_j \theta_j, \sum_{j=1}^h w_j = 1 \tag{30}$$

where w_j is the kriging weight of the residual θ_j at the j th index; the total weights must add up to one. The w_j can be calculated as follows:

$$\mathbf{w} = \mathbf{\Lambda}^{*-1} \cdot \mathbf{c}_{Q_{gs_i}} \tag{31}$$

$$\mathbf{c}_{Q_{gs_i}} = \begin{bmatrix} Cov(\theta_i, \theta_1) \\ Cov(\theta_i, \theta_2) \\ \vdots \\ Cov(\theta_i, \theta_h) \end{bmatrix} \tag{32}$$

where \mathbf{w} is an h -length vector of the w_j kriging weights; $\mathbf{\Lambda}^*$ is an $h \times h$ covariance matrix of residuals for S_{III} , with its diagonal elements constituted by $V(\theta_i)$, its off-diagonal elements calculated using $V(\theta_i)$, $V(\theta_j)$, and the spatial correlation ρ_{ij}^* ; i.e., $Cov(\theta_i, \theta_j) = \sqrt{V(\theta_i) \cdot V(\theta_j)} \cdot \rho_{ij}^*$. ρ_{ij}^* is expressed as $\rho_{ij}^* = 1 - \frac{\gamma(w_{ij})}{C_0^* + C_1^*}$, where $\gamma(w_{ij})$ is a semi-variogram in which the form is determined in the same manner as in $\gamma(d_{ij})$. C_0^* and C_1^* are the nugget and partial sill, respectively. $\mathbf{c}_{Q_{gs_i}}$ is an h -length vector of covariance for the unit map i .

Thus, the uncertainty of RK-3sGHMB for the map unit i can be written as follows:

$$RMSE(y_{Q_{gs}-RK_i}) = \sqrt{\left\{ \begin{aligned} & (\tilde{X}_i^{**} - \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} c_{Q_{gs_i}})^T \\ & (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} + \\ & \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{Z}_{S_{III}} (\tilde{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}})^{-1} \tilde{Z}_{S_{III}}^T \Lambda^{*-1} \tilde{X}_{S_{III}} (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} + \\ & \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{Z}_{S_{III}} \tilde{Z}_{S_{II}}^T \Sigma^{-1} \hat{Z}_{S_{II}} (\tilde{P}_{S_I}^T \Omega \tilde{P}_{S_I})^{-1} \hat{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}} \\ & (\tilde{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}})^{-1} \tilde{Z}_{S_{III}}^T \Lambda^{*-1} \tilde{X}_{S_{III}}^{**} (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} \\ & (\tilde{X}_i^{**} - \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} c_{Q_{gs_i}}) + V(\theta_i) - c_{Q_{gs_i}}^T \Lambda^{*-1} c_{Q_{gs_i}} \end{aligned} \right\}} \quad (33)$$

The population mean estimated by RK-3sGHMB can be expressed as

$$\mu_{Q_{gs}-RK} = \frac{1}{N} \sum_{i=1}^N \hat{y}_{Q_{gs}-RK_i} \quad (34)$$

with the variance estimator being

$$Var(\mu_{Q_{gs}-RK_i}) = \left\{ \begin{aligned} & I_U^T \left\{ (\tilde{X}_U^{**} - \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} c_{Q_{gs_U}}) \right. \\ & (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} + \\ & \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{Z}_{S_{III}} (\tilde{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}})^{-1} \tilde{Z}_{S_{III}}^T \Lambda^{*-1} \tilde{X}_{S_{III}} (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} + \\ & \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{Z}_{S_{III}} \tilde{Z}_{S_{II}}^T \Sigma^{-1} \hat{Z}_{S_{II}} (\tilde{P}_{S_I}^T \Omega \tilde{P}_{S_I})^{-1} \hat{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}} \\ & (\tilde{Z}_{S_{II}}^T \Sigma^{-1} \tilde{Z}_{S_{II}})^{-1} \tilde{Z}_{S_{III}}^T \Lambda^{*-1} \tilde{X}_{S_{III}}^{**} (\tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} \tilde{X}_{S_{III}}^{**})^{-1} \\ & \left. (\tilde{X}_U^{**} - \tilde{X}_{S_{III}}^{**T} \Lambda^{*-1} c_{Q_{gs_U}}) \right\} I_U \end{aligned} \right\} \quad (35)$$

where the parameters are the same as in Equation (28). If there is no spatial correlation in the dataset S_{III} , then $c_{Q_{gs_i}}$ and $c_{Q_{gs_U}}$ are approximately equal to the zero vector, and then Equations (34) and (35) are approximately equal to Equations (22), (27) and (28).

2.3.8. Evaluation Criteria for Modeling and Up-Scaling

We selected the adjusted coefficient of determination (R_{adj}^2), root mean square error (RMSE), and estimation accuracy (EA) as the evaluation indicators for the model fitting. The calculation formulas are given as (36), (37), (38), and (39).

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (36)$$

$$R_{adj}^2 = 1 - \left(1 - R^2\right) \frac{n - 1}{n - 1 - P} \quad (37)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - 1}} \quad (38)$$

$$EA = \left(1 - \frac{RMSE}{\bar{y}}\right) \times 100\% \quad (39)$$

where y_i represents the reference data; \hat{y}_i represents the predicted data; \bar{y} refers to the mean value of the reference data; and n is the number of reference data.

The precision of the population prediction for the five methods was evaluated using standard error (SE) and estimation precision (P) indicators, as shown in Formulas (40) and (41) [63].

$$SE = \sqrt{V(\bar{\mu}_{model})} \quad (40)$$

$$P = \left(1 - \frac{t \cdot SE}{\bar{y}}\right) \times 100\% \quad (41)$$

where $V(\bar{\mu}_{model})$ is the variance estimation of mean value by model; t is a reliability indicator, with a 95% confidence interval being $t = 1.96$.

3. Results

3.1. Feature Selection Results

Genetic algorithms were employed to perform feature selection for each model. To prevent overfitting, each model's feature selection was subjected to 5-fold and 10-fold cross-validation, conducted 50 times. The optimal number of iterations and features for the final genetic algorithm was determined by calculating the $RMSE$ and R^2 of the external samples. For the model F_{ml} , the optimal number of iterations was 46, with the optimal features being Hvar, H98, and HCD (Figure 4a). For the model F_{ms} , the optimal number of iterations was 49, with the optimal features being PHV, DVire1, EVI, and EVire1 (Figure 4b). For the model Q_{lg} , the optimal number of iterations was 39, with the optimal features being rh90, rh10, rh50, and cover (Figure 4c). For the model G_{ls} , the optimal number of iterations was 11, with the optimal features being KNDVire1, PHV, EVI, EVire1, and Vol (Figure 4d). For the model Q_{gs} , the optimal number of iterations was 20, with the optimal features being PHV, Hlx, EVire1, Vol, EVire3, KNDVire1, and PHH (Figure 4e).

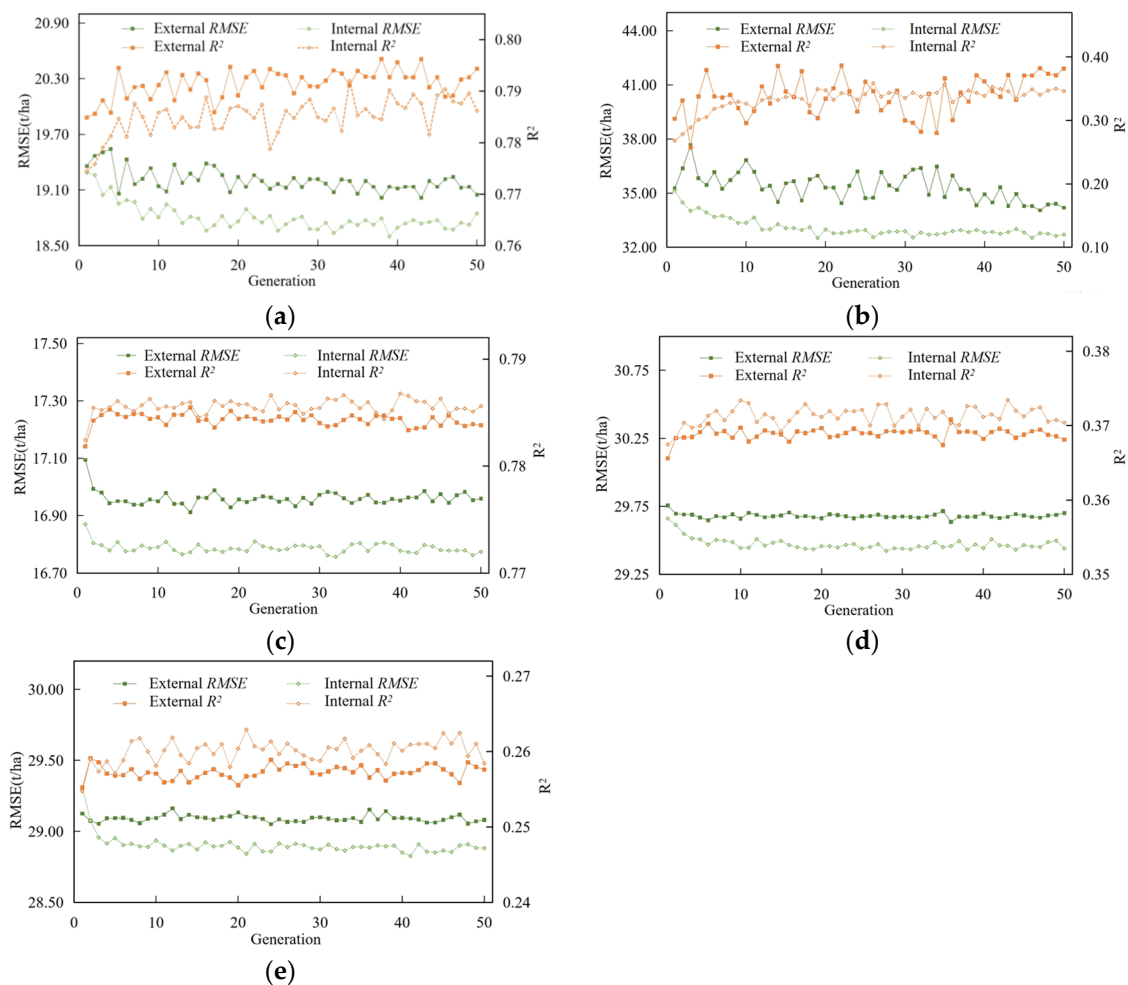


Figure 4. Variation in the feature selection iteration curves using genetic algorithms according to F_{ml} (a), F_{ms} (b), Q_{lg} (c), G_{ls} (d), and Q_{gs} (e).

3.2. Model Fitting Results and Accuracy Evaluation

Independent samples were employed to validate each model. The results indicate that models F_{ml} and Q_{lg} exhibit higher accuracy; both R_{adj}^2 values exceed 0.80, with estimation accuracies close to 80%. This suggests that using ALS-based AGB and GEDI-based AGB for multi-stage methods is feasible (Figure 5a,b). For the single-stage model using the CMB method, the model R_{adj}^2 is 0.37, with an RMSE of 33.95 t/ha and an accuracy of 63.28% (Figure 5c). In the case of multi-stage models, as the sample size increases, the precision of the GHMB model is improved, with an R_{adj}^2 reaching 0.38, an RMSE of 33.72 t/ha, and an accuracy reaching 63.53% (Figure 5d). However, the model accuracy of the 3sGHMB method decreases, with an R_{adj}^2 dropping to 0.27, an RMSE of 36.58 t/ha, and an accuracy of 60.43% (Figure 5f). The ratio of the nugget to the sill ($C_0/(C_0 + C_1)$), which represents the spatial dependence structure, is estimated for RK-GHMB and RK-3sGHMB. In general, a ratio below 25% indicates a strong spatial dependence structure. A ratio between 25% and 75% suggests a moderate spatial dependence structure, while a ratio exceeding 75% implies a weak spatial dependence structure [32]. In this study, all the ratios are less than 75% (Table 4, Figure 6). The global Moran's I of the two models is 0.20 and 0.40, with both p -values being 2.2×10^{-16} , also indicating the presence of spatial autocorrelation effects. The optimal residual variogram model forms for both models are exponential, with a range of 2625.28 m and 16,069.84 m, respectively, which means that the RK model could have a limited effect beyond these distances (Table 4). After considering spatial correlation, the model accuracy improved. The R_{adj}^2 of the RK-GHMB method increased to 0.60, with an RMSE reduction to 27.07 t/ha and an improvement in accuracy to 70.72% (Figure 5e). The R_{adj}^2 of the RK-3sGHMB method increased to 0.55, with an RMSE reduction to 28.55 t/ha and an accuracy improvement to 69.13% (Figure 5g). The fitting formulas for each model are presented in Table 5.

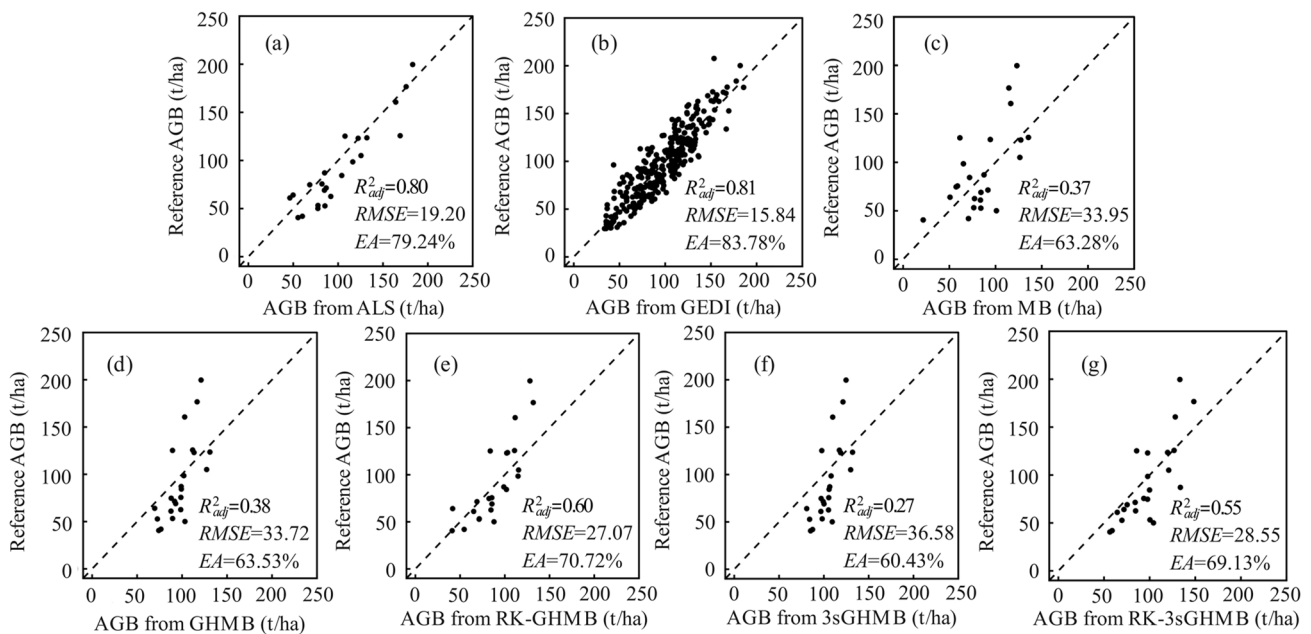


Figure 5. Model accuracy evaluation of F_{ml} (a), Q_{lg} (b), MB (c), GHMB (d), RK-GMB (e), 3sGHMB (f), and RK-3sGHMB (g) using independent samples.

Table 4. Fitting accuracy of the semivariogram for RK-3sGHMB and RK-GHMB.

Residual Source	Fitting Method	Nugget (C_0/C_0^*)	Partial Sill (C_1/C_1^*)	Range (m)	Ratio (%)	R^2_{adj}	RMSE (t/ha)	EA (%)
RK-3sGHMB	Exponential	328.35	525.55	16,069.84	38.45	0.55	28.55	69.13
	Spherical	414.00	420.82	7109.70	49.59	0.49	30.67	66.83
	Gaussian	516.99	302.11	10,615.19	63.12	0.48	30.70	66.79
	Linear	747.71	747.71	/	/	0.38	38.22	58.7
RK-GHMB	Exponential	120.02	539.06	2625.28	18.21	0.60	27.07	70.72
	Spherical	338.20	533.07	2019.45	38.82	0.53	29.21	68.41
	Gaussian	380.09	533.08	2288.70	41.62	0.41	32.84	64.48
	Linear	610.96	610.96	/	/	0.41	32.85	64.47

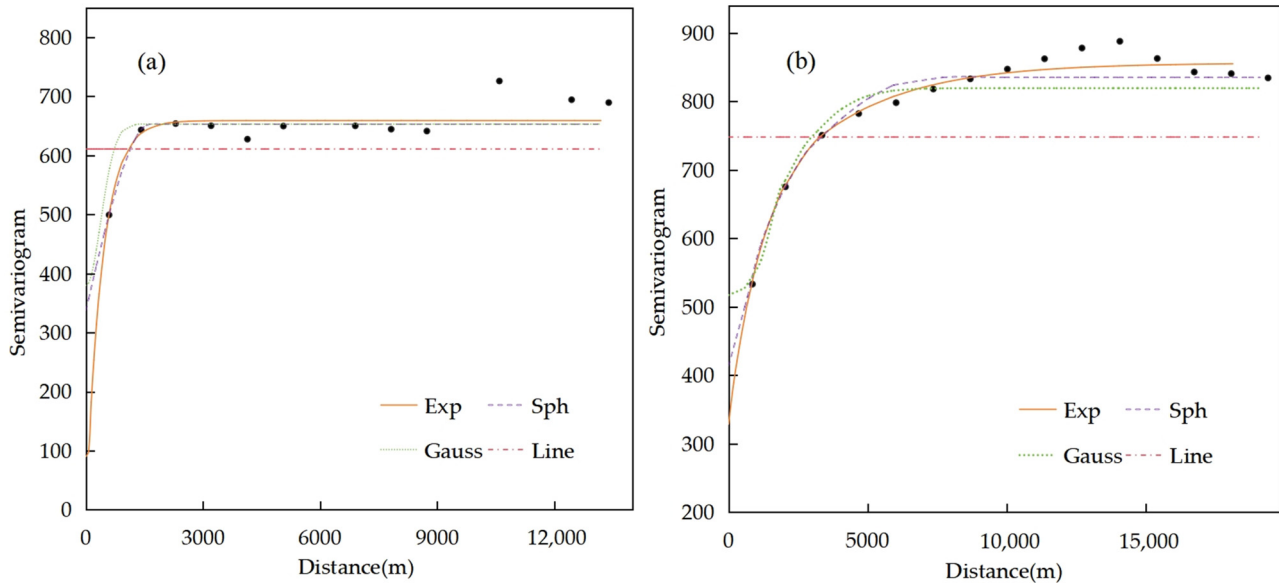


Figure 6. Residual variogram and fitted model with exponential, gaussian, spherical, and linear fitting using RK-GHMB (a) and RK-3sGHMB (b).

Table 5. Fitted forms for models F_{ms} , F_{ml} , G_{ls} , Q_{lg} , and Q_{gs} .

Model	Model Formula	R^2_{adj}	RMSE (t/ha)	EA(%)
F_{ms}	$201.09464 / (1 + \exp(-0.70266 \times EVI - 0.09913 \times PHV + 21.18109 \times DVire1 - 0.14215 \times EVIre1 - 9.69895))$	0.37	33.95	63.28
F_{ml}	$299.30254 / (1 + \exp(-0.074878 \times Hvar + 0.026074 \times H98 - 0.006368 \times CD + 1.947972))$	0.80	19.20	79.24
G_{ls}	$189.4068 / (1 + \exp(-29.7195 \times EVI - 0.1698 \times PHV + 30.0420 \times KNDVire1 + 41.6752 \times Vol + 18.6999 \times EVIre1 + 17.9654))$	0.38	33.72	63.53
Q_{lg}	$288.317726 / (1 + \exp(-0.094833 \times rh90 - 0.046170 \times rh50 + 0.059449 \times rh10 + 0.003274 \times Cover + 3.446146))$	0.81	15.84	83.78
Q_{gs}	$199.54404 / (1 + \exp(20.57206 \times Hlx - 9.80876 \times EVIre3 + 7.25995 \times KNDVire1 + 4.68176 \times EVIre1 - 0.18229 \times PHV - 0.08263 \times PHH + 0.02713 \times Vol + 5.37781))$	0.27	36.58	60.43

3.3. Model Estimation and Uncertainty Analysis

We estimated the uncertainty at the pixel level using the five inferences; the spatial distribution is shown in Figure 7a–e. It can be seen that the uncertainty range of the CMB is 2–48.3 t/ha. The uncertainty of the forest AGB low-value area from CMB is lower than that of other models, while the uncertainty in high-value areas is larger. The uncertainty range of the GHMB is 10–32.22 t/ha, and that of the 3sGHMB model is 26.09–41.40 t/ha. The uncertainty of the RK-GHMB and RK-3sGHM, which, when considering spatial correlation, ranges from 0.61 to 32.69 t/ha and from 18.06 to 38.63 t/ha, respectively. Their uncertainty is significantly lower than that of GHMB and 3sGHMB, especially in areas near the samples;

the uncertainty of RK-3sGHMB is the smallest, but the uncertainty in the low-value area is larger (Figure 7f).

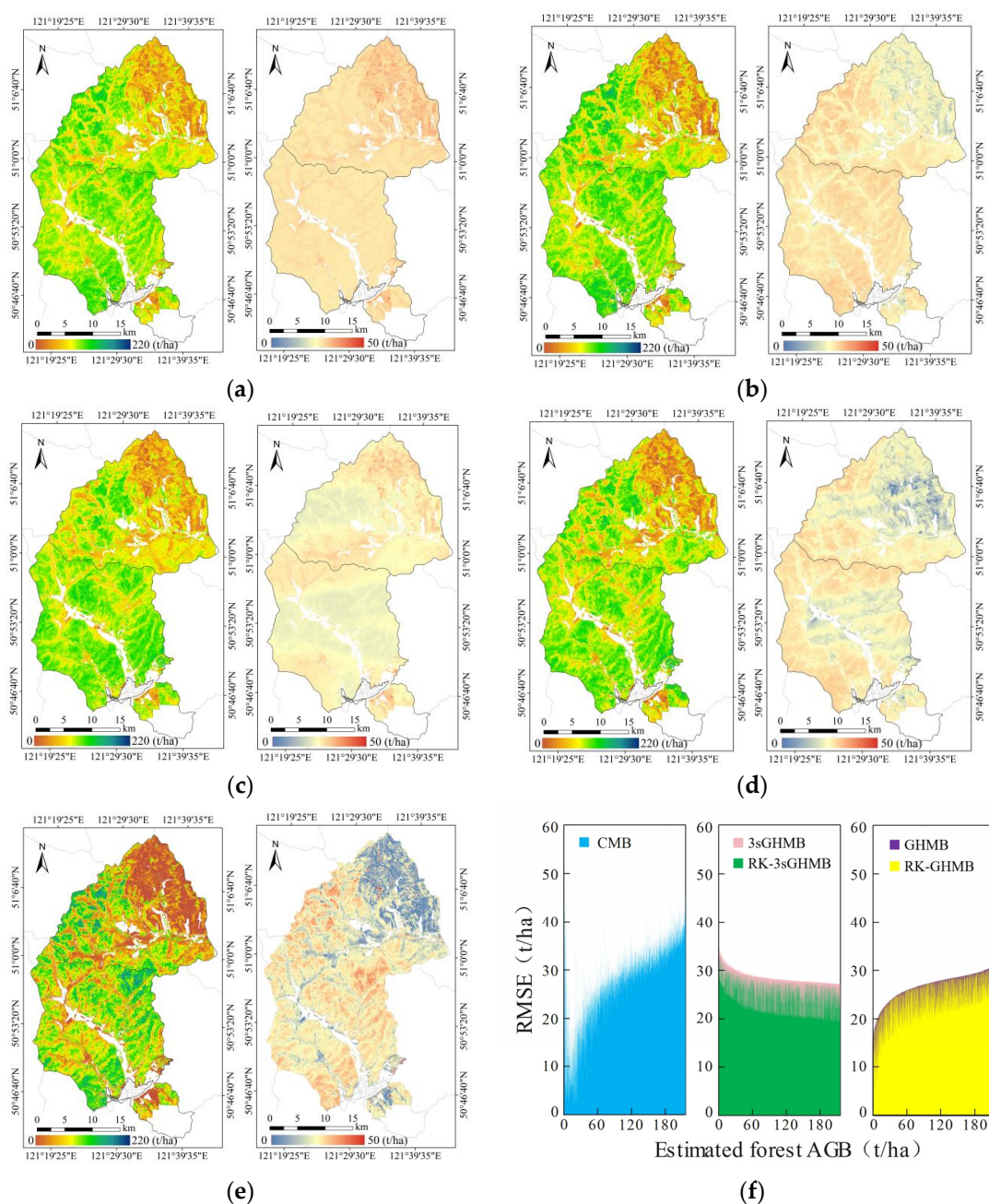


Figure 7. The results for forest AGB and uncertainty: the spatial distribution of forest AGB and uncertainty estimated by 3sGHMB (a), GHMB (b), RK-3sGHMB (c), RK-GHMB (d), and CMB (e); (f) pixel-level uncertainty ($RMSE$) with respect to predicted AGB.

From the population mean estimation results and uncertainty, it can be seen that the precision of the multi-stage method is superior to that of single-stage inference, considering that the spatial correlation of the multi-stage method shows better accuracy than that of the original models. Among the models, the precision of the population mean by RK-3sGHMB is the highest, with the precision of the Upper Yangge and Tidal Slag Forest Farms reaching 93.48% and 95.39%, respectively, with the combined population accuracy of the two farms being 94.44% (Table 6).

Table 6. The population results and uncertainties for regions.

Region	Methods	Population Mean (t/ha)	Variance in the Population Mean (t/ha)	SE (t/ha)	<i>p</i> (%)
Upper Yangge Qi Forest Farm	CMB	80.43	30.27	5.50	86.60
	GHMB	83.63	19.43	4.41	89.67
	RK-GHMB	81.50	9.99	3.16	92.40
	3sGHMB	84.20	29.33	5.42	87.39
	RK-3sGHMB	79.16	6.94	2.63	93.48
Tidal Slag Forest Farm	CMB	98.66	28.98	5.38	89.31
	GHMB	102.75	17.20	4.15	92.09
	RK-GHMB	100.79	8.77	2.96	94.24
	3sGHMB	101.07	24.35	4.93	90.43
	RK-3sGHMB	103.72	5.96	2.44	95.39
Total	CMB	89.55	29.63	5.44	87.96
	GHMB	93.19	18.32	4.28	90.88
	RK-GHMB	91.15	9.38	3.06	93.32
	3sGHMB	92.64	26.84	5.18	88.91
	RK-3sGHMB	91.44	6.45	2.54	94.44

4. Discussion

The results of this study indicate that although multi-stage inferences are trained with a larger number of samples in the second or third modeling stage, only the accuracy of GHMB is improved, while the accuracy of 3sGHMB is lower than that of CMB (Figure 5). The primary reason for this may be related to the propagation of errors. Although two-stage inference also experiences error propagation, the features of ALS are strongly correlated with plot AGB. By utilizing ALS features to build a model to plot AGB, to predict the corresponding AGB at other locations in the sampled ALS strips, we thereby expanded the sample size for the second stage of modeling, while keeping the expanded samples more uniformly distributed in space and more representative of the population. The errors introduced by error propagation can be compensated for by the increased number of units in the expanded samples; hence, the model accuracy is still improved, which is consistent with the conclusions of existing research [14]. In the case of three-stage inference that does not consider spatial relationships, the final model is affected by the error propagation from the first two models. Whether this is the cause of the reduced model accuracy requires further investigation.

For population mean estimation, as the AGB sample size used in modeling increases, the representativeness of the sample improves, leading to higher precision and more stable estimation results for the population (Figure 8), which is consistent with existing research [14,17,64,65]. It is important to highlight that all the samples are obtained through random sampling, as shown in Figure 8. For two-stage and three-stage inferences, the change in sample quantity pertains exclusively to the final stage [33]. Concurrently, when estimating models in the other stages, all available samples are utilized to ensure optimal model accuracy [14]. Models that account for spatial correlation exhibit higher accuracy (Table 6, Figure 8), which is closely related to the sample size employed in modeling. Specifically, when the sample size is less than approximately 489, methods that consider spatial relationships demonstrate a comparable accuracy in population mean estimation to those that do not account for spatial correlation. In contrast, when the sample size exceeds approximately 489, the accuracy of models incorporating spatial relationships gradually improves with the increasing sample size. This trend is attributed to the inherent characteristics of geostatistical methods [66]. Under identical sample size conditions, the CMB has a lower population estimation uncertainty than other inferences; this is due to the error propagation in multi-stage models [17]. In this study, the sample size used for the CMB model was only 103, and the curve for CMB model variation and sample size only reached 103. However, Saarela et al. [67] also showed that under the same sample size

conditions, the uncertainty of the CMB is less than that of multi-stage inferences. When the sample size is less than approximately 2998, the precision of population estimation results of two-stage inference (GHMB and RK-GHMB) is always better than that of three-stage inferences (3sGHMB and RK-3sGHMB). When the sample size exceeds approximately 2998, as the sample size increases, the estimation results for RK-3sGHMB begin to surpass those of GHMB but remain lower than those of RK-GHMB. When the sample size exceeds a certain value, the RK-3sGHMB model estimation results are optimal among all methods, which is also the reason why the precision of RK-3sGHMB for population estimation is the highest in this study, indicating that geostatistical methods can effectively improve model estimation accuracy. Although the accuracy of the 3sGHMB model is lower than that of CMB, when the sample size exceeds approximately 1000, the population estimation precision is better than that of the CMB model with only 103 samples, but it can never compare with GHMB, RK-GHMB, and RK-3sGHMB.

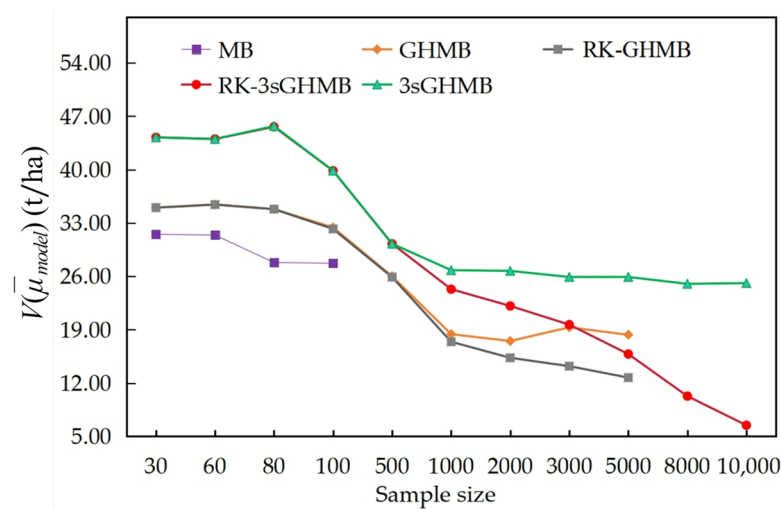


Figure 8. The influence of sample size used for $V(\bar{\mu}_{model})$ estimation using MB, GHMB, RH-GHMBa, 3sGHMB, and RK-3sGHMB.

Based on the aforementioned analysis, we can categorize the application scenarios for inferences in this study as follows. (1) The integration of limited sample plots into comprehensive wall-to-wall data, including, but not limited to, ALS data, optical imagery, and SAR data, as exemplified in Case A, represents a widely adopted framework for AGB estimation. This framework is particularly well suited to CMB. However, the costs associated with field plot surveys are substantial, which can represent a significant barrier to implementation. In an ideal scenario without budgetary constraints, wall-to-wall ALS data would be the optimal choice due to their high accuracy and detailed information. Nonetheless, obtaining full-coverage ALS data over extensive areas is often impractical due to logistical and financial limitations. (2) For smaller study areas (counties, townships, forest farms, etc.), the strategic sampling of ALS data (collected by flying over some strips or blocks in the study area) can reduce the sample size of field plot data required. Utilizing sampled ALS data to estimate target variables can expand the sample size, potentially reducing the cost associated with field plot acquisition and enhancing the representativeness of the samples, particularly for regions that are inaccessible to humans. This framework is amenable to both GHMB and RK-GHMB. It is imperative to note that the RK-GHMB method can only manifest its advantages when there is an adequate number of ALS data to expand the modeling samples in the second stage. Furthermore, if actual AGB data at the satellite waveform LiDAR footprint locations can be procured, then the ALS data used for sample modeling could be substituted with more cost-effective satellite LiDAR data,

such as GEDI, ICESat-2, TECIS, and GF-7. However, obtaining actual values for footprint locations in practice remains challenging. (3) In the context of larger study areas (such as at the national, provincial, or municipal level) and with budgetary constraints on ALS data acquisition, employing a limited number of field plots in conjunction with partial sampling ALS data linking the field plots to satellite footprints to augment sample size can improve the uniformity and representativeness of sample distribution. In such scenarios, the RK-3sGHMB exhibits certain advantages.

The results of this study demonstrate that RK-3sGHMB can improve the estimation accuracy of model fitting and the precision of population prediction and can enhance the rationality of the uncertainty estimation results; however, geostatistical methods were only applied in the final stage. Given that the second-stage model also exhibits spatial correlation, future research will explore a three-stage inference estimation approach that accounts for spatial correlation in both the second and third stages. However, the derivation of the uncertainty estimation formula for the second stage using the RK model is too complex and challenging to implement. The spatial filtering estimation method for feature vectors could serve as a primary direction for future research [68,69]. This method has the potential to reduce the model formula complexity and improve the operational efficiency of the model.

5. Conclusions

In the present study, we successfully implemented a three-stage forest AGB estimation method, RK-3sGHMB, by integrating RK into the third stage of 3sGHMB, and conducted a comparative analysis with the CMB, GHMB, RK-GHMB, and 3sGHMB methods, leading to the following conclusions. (1) When the sample data show no spatial correlation or the model estimation does not take spatial correlation into account, multi-stage inference does not demonstrate a distinct advantage, yielding a lower estimation accuracy under equivalent sample size conditions compared with single-stage inference. (2) The integration of geostatistical methods can significantly enhance the model accuracy, which is correlated with the number of samples. As the sample size in the last stage increases, the precision of the estimation results consistently improves. With an adequate number of GEDI samples, the RK-3sGHMB model proposed in this study exhibits superior population estimation precision compared to other models, offering certain advantages in the estimation of forest resource parameters across broader areas.

Author Contributions: Conceptualization, E.C. and X.D.; methodology, X.D. and E.C.; validation, X.D.; formal analysis, X.D.; investigation, X.D., E.C., L.Z. and Y.M.; data curation, X.D., Y.F. and J.W.; writing—original draft preparation, X.D.; writing—review and editing, E.C., L.Z. and X.D.; visualization, X.D.; funding acquisition, E.C. All authors contributed to interpreting the results and to the improvement of the article. All authors have read and agreed to the published version of the manuscript.

Funding: This study was funded by the National Key R&D Program of China (grant number: 2023YFF1303900).

Data Availability Statement: The field data, ALS data, and P-SAR data are not publicly available to protect the privacy of private landowners.

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. Dixon, R.K.; Brown, S.; Houghton, R.A.; Solomon, A.M.; Trexler, M.C.; Wisniewski, J. Carbon pools and flux of global forest Ecosystems. *Science* **1994**, *263*, 185–190. [[CrossRef](#)] [[PubMed](#)]
2. Hese, S.; Lucht, W.; Schmullius, C.; Barnsley, M.; Dubayah, R.; Knorr, D.; Neumann, K.; Riedel, T.; Schröter, K. Global biomass mapping for an improved understanding of the CO₂ balance—the earth observation mission carbon-3D. *Remote Sens. Environ.* **2005**, *94*, 94–104. [[CrossRef](#)]
3. Chen, Q. Modeling aboveground tree woody biomass using national-scale allometric methods and airborne lidar. *ISPRS J. Photogramm. Remote Sens.* **2015**, *106*, 95–106. [[CrossRef](#)]
4. Askne, J.I.H.; Soja, M.J.; Ulander, L.M.H. Biomass estimation in a boreal forest from TanDEM-X data, lidar DTM, and the interferometric water cloud model. *Remote Sens. Environ.* **2017**, *196*, 265–278. [[CrossRef](#)]
5. Ståhl, G.; Saarela, S.; Schnell, S.; Holm, S.; Breidenbach, J.; Healey, S.P.; Patterson, P.L.; Magnussen, S.; Næsset, E.; McRoberts, R.E.; et al. Use of models in large-area forest surveys: Comparing model-assisted, model-based and hybrid estimation. *For. Ecosyst.* **2016**, *3*, 5. [[CrossRef](#)]
6. McRoberts, R.E.; Nasset, E.; Gobakken, T.; Chirici, G.; Condés, S.; Hou, Z.; Saarela, S.; Chen, Q.; Ståhl, G.; Walters, B.F. Assessing components of the model-based mean square error estimator for remote sensing assisted forest applications. *Can. J. For. Res.* **2018**, *48*, 642–649. [[CrossRef](#)]
7. Marshak, C.; Simard, M.; Duncanson, L.; Silva, C.A.; Denbina, M.; Liao, T.-H.; Fatoyinbo, L.; Moussavou, G.; Armston, J. Regional tropical aboveground biomass mapping with L-Band Repeat-Pass Interferometric Radar, Sparse lidar, and multiscale superpixels. *Remote Sens.* **2020**, *12*, 2048. [[CrossRef](#)]
8. Saarela, S.; Varvia, P.; Korhonen, L.; Yang, Z.; Patterson, P.L.; Gobakken, T.; Næsset, E.; Healey, S.P.; Ståhl, G. Tree-phase hierarchical model-based and hybrid inference. *MethodsX* **2023**, *11*, 102321. [[CrossRef](#)]
9. Saarela, S.; Wästlund, A.; Holmström, E.; Mensah, A.A.; Holm, S.; Nilsson, M.; Fridman, J.; Ståhl, G. Mapping aboveground biomass and its prediction uncertainty using LiDAR and field data, accounting for tree-level allometric and LiDAR model errors. *For. Ecosyst.* **2020**, *7*, 43. [[CrossRef](#)]
10. Magnussen, S.; Nord-Larsen, T.; Nielsen, T.R. Lidar supported estimators of wood volume and aboveground biomass from the Danish national forest inventory (2012–2016). *Remote Sens. Environ.* **2018**, *211*, 146–153. [[CrossRef](#)]
11. Silva, C.A.; Hudak, A.T.; Vierling, L.A.; Loudermilk, E.L.; O'Brien, J.J.; Hiers, J.K.; Jack, S.B.; Gonzalez-Benecke, C.; Lee, H.; Falkowski, M.J.; et al. Imputation of individual longleaf pine (*pinus palustris* mill.) tree attributes from field and LiDAR data. *Can. J. Remote Sens* **2016**, *42*, 554–573. [[CrossRef](#)]
12. Xu, C.; Manley, B.; Morgenroth, J. Evaluation of modelling approaches in predicting forest volume and stand age for small-scale plantation forests in New Zealand with RapidEye and LiDAR. *Int. J. Appl. Earth Obs. Geoinf.* **2018**, *73*, 386–396. [[CrossRef](#)]
13. Lim, K.; Treitz, P.; Wulder, M.; St-Onge, B.; Flood, M. LiDAR remote sensing of forest structure. *Prog. Phys. Geogr.* **2003**, *27*, 88–106. [[CrossRef](#)]
14. Xie, B.; Cao, C.; Xu, M.; Bashir, B.; Singh, R.P.; Huang, Z.; Lin, X. Regional forest volume estimation by expanding LiDAR samples using multi-sensor satellite data. *Remote Sens.* **2020**, *12*, 360. [[CrossRef](#)]
15. McRoberts, R.E.; Næsset, E.; Gobakken, T. Inference for lidar-assisted estimation of forest growing stock volume. *Remote Sens. Environ.* **2013**, *128*, 268–275. [[CrossRef](#)]
16. McRoberts, R.E.; Næsset, E.; Saatchi, S.; Quegan, S. Statistically rigorous, model-based inferences from map. *Remote Sens. Environ.* **2022**, *279*, 113028. [[CrossRef](#)]
17. Saarela, S.; Holm, S.; Grafström, A.; Schnell, S.; Næsset, E.; Gregoire, T.G.; Nelson, R.F.; Ståhl, G. Hierarchical model-based inference for forest inventory utilizing three sources of information. *Ann. For. Sci.* **2016**, *73*, 895–910. [[CrossRef](#)]
18. Saarela, S.; Holm, S.; Healey, S.P.; Patterson, P.L.; Yang, Z.; Andersen, H.E.; Dubayah, R.O.; Qi, W.; Duncanson, L.I.; Armston, J.D.; et al. Comparing frameworks for biomass prediction for the Global Ecosystem Dynamics Investigation. *Remote Sens. Environ.* **2022**, *278*, 113074. [[CrossRef](#)]
19. Holm, S.; Nelson, R.; Ståhl, G. Hybrid three-phase estimators for large-area forest inventory using ground plots, airborne lidar, and space lidar. *Remote Sens. Environ.* **2017**, *197*, 85–97. [[CrossRef](#)]
20. Pang, Y.; Yu, T.; Jia, W.; Liang, X.; Li, Z.; Fu, A.; Wu, F.; Liu, X.; Zhang, X.; Huang, J.; et al. TECIS: The first mission towards forest carbon mapping by combination of lidar and multi-angle optical observations. In *IGARSS 2023-2023 IEEE International Geoscience and Remote Sensing Symposium*; IEEE: Pasadena, CA, USA, 2023; pp. 1890–1893. [[CrossRef](#)]
21. Varvia, P.; Saarela, S.; Maltamo, M.; Packalen, P.; Gobakken, T.; Næsset, E.; Ståhl, G.; Korhonen, L. Estimation of boreal forest biomass from ICESat-2 data using hierarchical hybrid inference. *Remote Sens. Environ.* **2024**, *311*, 114249. [[CrossRef](#)]
22. Wulder, M.A.; White, J.C.; Bater, C.W.; Coops, N.C.; Hopkinson, C.; Chen, G. Lidar plots—a new large-area data collection option: Context, concepts, and case study. *Can. J. Remote Sens.* **2012**, *38*, 600–618. [[CrossRef](#)]

23. Margolis, H.A.; Nelson, R.F.; Montesano, P.M.; Beaudoin, A.; Sun, G.; Andersen, H.E.; Wulder, M.A. Combining satellite lidar, airborne lidar, and ground plots to estimate the amount and distribution of aboveground biomass in the boreal forest of North America. *Can. J. For. Res.* **2015**, *45*, 838–855. [[CrossRef](#)]
24. Narine, L.L.; Popescu, S.C.; Malambo, L. Using ICESat-2 to estimate and map forest aboveground biomass: A first example. *Remote Sens.* **2020**, *12*, 1824. [[CrossRef](#)]
25. Varvia, P.; Korhonen, L.; Bruguère, A.; Toivonen, J.; Packalen, P.; Maltamo, M.; Saarela, S.; Popescu, S.C. How to consider the effects of time of day, beam strength, and snow cover in ICESat-2 based estimation of boreal forest biomass? *Remote Sens. Environ.* **2022**, *280*, 113174. [[CrossRef](#)]
26. Guerra-Hernández, J.; Narine, L.L.; Pascual, A.; Gonzalez-Ferreiro, E.; Botequim, B.; Malambo, L.; Neuenschwander, A.; Popescu, S.C.; Godinho, S. Aboveground biomass mapping by integrating ICESat-2, SENTINEL-1, SENTINEL-2, ALOS2/PALSAR2, and topographic information in mediterranean forests. *GISci. Remote Sens.* **2022**, *59*, 1509–1533. [[CrossRef](#)]
27. Zakeri, F.; Mariethoz, G. A review of geostatistical simulation models applied to satellite remote sensing: Methods and applications. *Remote Sens. Environ.* **2021**, *259*, 112381. [[CrossRef](#)]
28. Fayad, I.; Baghdadi, N.; Bailly, J.S.; Barbier, N.; Gond, V.; Herault, B.; El Hajj, M.; Fabre, F.; Perrin, J. Regional scale rain-forest height mapping using regression-kriging of spaceborne and airborne LiDAR data: Application on French Guiana. *Remote Sens.* **2016**, *8*, 240. [[CrossRef](#)]
29. Pouladi, N.; Møller, A.B.; Tabatabai, S.; Greve, M.H. Mapping soil organic matter contents at field level with cubist, random forest and kriging. *Geoderma* **2019**, *342*, 85–92. [[CrossRef](#)]
30. Silveira, E.M.; Santo, F.D.E.; Wulder, M.A.; Júnior, F.W.A.; Carvalho, M.C.; Mello, C.R.; Mello, J.M.; Shimabukuro, Y.E.; Terra, M.C.N.S.; Carvalho, L.M.T.; et al. Pre-stratified modelling plus residuals kriging reduces the uncertainty of aboveground biomass estimation and spatial distribution in heterogeneous savannas and forest environments. *For. Ecol. Manag.* **2019**, *445*, 96–109. [[CrossRef](#)]
31. Chen, L.; Ren, C.; Zhang, B.; Wang, Z. Multi-sensor prediction of stand volume by a hybrid model of support vector machine for regression kriging. *Forests* **2020**, *11*, 296. [[CrossRef](#)]
32. Zhao, J.; Zhao, L.; Chen, E.; Li, Z.; Xu, K.; Ding, X. An improved generalized hierarchical estimation framework with geostatistics for mapping forest parameters and its uncertainty: A case study of forest canopy height. *Remote Sens.* **2022**, *14*, 568. [[CrossRef](#)]
33. Chen, F.; Hou, Z.; Saarela, S.; McRoberts, R.E.; Ståhl, G.; Kangas, A.; Packalen, P.; Li, B.; Xu, Q. Leveraging remotely sensed non-wall-to-wall data for wall-to-wall upscaling in forest inventory. *Int J Appl Earth Obs Geoinf.* **2023**, *119*, 103314. [[CrossRef](#)]
34. Li, C.; Zhang, W.; Li, Z.; Chen, E.; Tian, X. Retrieval of forest above-ground biomass using multi-source data in Genhe, Inner Mongolia. *J. Beijing For. Univ. (Chin. Ed.)* **2016**, *38*, 64–72. [[CrossRef](#)]
35. Zhou, G.; Yin, G.; Tang, X. *Biomass Equation for Forest Ecosystems in China-Carbon Storage*; China Science Publishing: Beijing, China, 2018; pp. 41–80.
36. Pang, Y.; Li, Z.; Ju, H.; Lu, H.; Jia, W.; Si, L.; Guo, Y.; Liu, Q.; Li, S.; Liu, L.; et al. LiCHy: The CAF's LiDAR, CCD and hyperspectral integrated airborne observation system. *Remote Sens.* **2016**, *8*, 398. [[CrossRef](#)]
37. Isenburg, M. LAStools—Efficient Tools for LiDAR Processing. Version 2.0.1. 2022. Available online: <https://rapidlasso.de/lastools-220107> (accessed on 25 March 2022).
38. McGaughey, R.J. *FUSION/LDV: Software for LiDAR Data Analysis and Visualization*; January 2021-FUSION Version 4.20; United States Department of Agriculture: Washington, DC, USA, 2021. Available online: http://forsys.cfr.washington.edu/software/fusion/FUSION_manual.pdf (accessed on 16 June 2021).
39. Cao, L.; Xu, T.; Shen, X.; She, G. Mapping biomass by integrating Landsat OLI and airborne LiDAR transect data in subtropical forests. *J. Remote Sens.* **2016**, *20*, 665–678. [[CrossRef](#)]
40. Donoghue, D.N.M.; Watt, P.J.; Cox, N.J.; Wilson, J. Remote sensing of species mixtures in conifer plantations using LiDAR height and intensity data. *Remote Sens. Environ.* **2007**, *110*, 509–522. [[CrossRef](#)]
41. Fan, Y.; Zhao, L.; Chen, E.; Xu, K.; Zhang, W.; Ma, Y. Evaluation of forest stock estimation ability of high resolution airborne multi-band PolSAR in cold temperate coniferous forests. *Natl. Remote Sens. Bull.* **2024**, *28*, 2525–2539. [[CrossRef](#)]
42. Zhao, L.; Chen, E.; Li, Z.; Zhang, W.; Gu, X. Three-step semi-empirical radiometric terrain correction approach for PolSAR data applied to forested areas. *Remote Sens.* **2017**, *9*, 269. [[CrossRef](#)]
43. Yamaguchi, Y.; Moriyama, T.; Ishido, M.; Yamada, H. Four-component scattering model for polarimetric SAR image decomposition. *IEEE T. Geosci. Remote.* **2005**, *43*, 1699–1706. [[CrossRef](#)]
44. Spracklen, B.; Spracklen, D.V. Detrending of structural characteristics of old-growth forest in Ukmine using spaceborne LiDAR. *Remote Sens.* **2021**, *13*, 1233. [[CrossRef](#)]
45. Lin, X.; Xu, M.; Cao, C.; Dang, Y.; Bashir, B.; Xie, B.; Huang, Z. Estimates of forest canopy height using a combination of ICESat-2/ATLAS data and stereo-photogrammetry. *Remote Sens.* **2020**, *12*, 3649. [[CrossRef](#)]
46. Liu, L.; Wang, C.; Nie, S.; Zhu, X.; Xi, X.; Wang, J. Analysis of the influence of different algorithms of GEDI L2A on the accuracy of ground elevation and forest canopy height. *J. Univ. Chin. Acad. Sci.* **2022**, *39*, 502–511. [[CrossRef](#)]

47. Dorado-Ronda, I.; Pascual, A.; Godinho, S.; Silva, C.A.; Botequim, B.; Ro-dríguez-González, P.; González-Ferreiro, E.; Guerra-Hernández, J. Assessing the accuracy of GEDI data for canopy height and aboveground biomass estimates in Mediterranean forests. *Remote Sens.* **2021**, *13*, 2279. [[CrossRef](#)]
48. Qi, W.; Saarela, S.; Armston, J.; Ståhl, G.; Dubayah, R. Forest biomass estimation over three distinct forest types using TanDEM-X InSAR data and simulated GEDI lidar data. *Remote Sens. Environ.* **2019**, *232*, 111283. [[CrossRef](#)]
49. Dubayah, R.; Blair, J.B.; Goetz, S.; Fatoyinbo, L.; Hansen, M.; Healey, S.; Hofton, M.; Hurtt, G.; Kellner, J.; Luthcke, S.; et al. The Global Ecosystem Dynamics Investigation: High-resolution laser ranging of the Earth's forests and topography. *Sci. Remote Sens.* **2020**, *1*, 100002. [[CrossRef](#)]
50. Silva, C.A.; Hamamura, C.; Valbuena, R.; Hancock, S.; Cardil, A.; Broadbent, E.N.; Almeida, D.R.A.; Silva, J.; Klauber, C. rGEDI: NASA's Global Ecosystem Dynamics Investigation (GEDI) Data Visualization and Processing. Version 0.1.9. Available online: <https://CRAN.R-project.org/package=rGEDI> (accessed on 22 October 2020).
51. Camps-Valls, G.; Campos-Taberner, M.; Moreno-Martínez, Á.; Walther, S.; Duveiller, G.; Cescatti, A.; Mahecha, M.D.; Muñoz-Marí, J.; García-Haro, F.J.; Guanter, L. A unified vegetation index for quantifying the terrestrial biosphere. *Sci. Adv.* **2021**, *7*, eabc7447. [[CrossRef](#)]
52. Rouse, J.; Haas, R.H.; Deering, D.; Schell, J.; Harlan, J. *Monitoring the Vernal Advancement and Retrogradation (Green Wave Effect) of Natural Vegetation*; Technical Report; NASA/GSFC Type III Final Report; NASA/GSFC: Greenbelt, MD, USA, 1974; Corpus ID: 129198382.
53. Shoko, C.; Mutanga, O. Examining the strength of the newly-launched Sentinel 2 MSI sensor in detecting and discriminating subtle differences between C3 and C4 grass species. *ISPRS J. Photogramm. Remote Sens.* **2017**, *129*, 32–40. [[CrossRef](#)]
54. Richardson, A.J.; Wiegand, C.L. Distinguishing vegetation from soil background information. *Photogramm. Eng. Remote Sens.* **1977**, *43*, 1541–1552, Corpus ID: 126604551.
55. Zhang, X.; Shen, H.; Huang, T.; Wu, Y.; Guo, B.; Liu, Z.; Luo, H.; Tang, J.; Zhou, H.; Wang, L.; et al. Improved random forest algorithms for increasing the accuracy of forest aboveground biomass estimation using Sentinel-2 imagery. *Ecol. Indic.* **2024**, *159*, 111752. [[CrossRef](#)]
56. Liu, H.Q.; Huete, A. A feedback based modification of the NDVI to minimize canopy background and atmospheric noise. *IEEE T. Geosci. Remote.* **1995**, *33*, 457–465. [[CrossRef](#)]
57. Jordan, C.F. Derivation of leaf-area index from quality of light on the forest floor. *Ecology* **1969**, *5*, 663–666. [[CrossRef](#)]
58. Huete, A.R. A soil-adjusted vegetation index (SAVI). *Remote Sens. Environ.* **1988**, *25*, 295–309. [[CrossRef](#)]
59. Gitelson, A.A.; Kaufman, Y.J.; Merzlyak, M.N. Use of a green channel in remote sensing of global vegetation from EOS-MODIS. *Remote Sens. Environ.* **1996**, *58*, 289–298. [[CrossRef](#)]
60. McRoberts, R.E.; Domke, G.M.; Chen, Q.; Næsset, E.; Gobakken, T. Using genetic algorithms to optimize k-Nearest Neighbors configurations for use with airborne laser scanning data. *Remote Sens. Environ.* **2016**, *184*, 387–395. [[CrossRef](#)]
61. Kuhn, M. Building predictive models in R using the caret package. *J. Stat. Softw.* **2008**, *28*, 1–26. [[CrossRef](#)]
62. McRoberts, R.E.; Næsset, E.; Hou, Z.; Ståhl, G.; Saarela, S.; Esteban, J.; Travaglini, D.; Mohammadi, J.; Chirici, G. How many bootstrap replications are necessary for estimating remote sensing-assisted, model-based standard errors? *Remote Sens. Environ.* **2023**, *288*, 113455. [[CrossRef](#)]
63. Tao, S.; Labrière, N.; Calders, K.; Fischer, F.J.; Rau, E.; Plaisance, L.; Chave, J. Mapping tropical forest trees across large areas with lightweight cost-effective terrestrial laser scanning. *Ann. For. Sci.* **2021**, *78*, 103. [[CrossRef](#)]
64. Saarela, S.; Schnell, S.; Grafström, A.; Tuominen, S.; Nordkvist, K.; Hyypä, J.; Kangas, A.; Ståhl, G. Effects of sample size and model form on the accuracy of model-based estimators of growing stock volume. *Can. J. For. Res.* **2015**, *45*, 1524–1534. [[CrossRef](#)]
65. Li, C.; Yu, Z.; Dai, H.; Zhou, X.; Zhou, M. Effect of sample size on the estimation of forest inventory attributes using airborne LiDAR data in large-scale subtropical areas. *Ann. For. Sci.* **2023**, *80*, 40. [[CrossRef](#)]
66. Voss, S.; Zimmermann, B.; Zimmermann, A. Detecting spatial structures in throughfall data: The effect of extent, sample size, sampling design, and variogram estimation method. *J. Hydrol.* **2016**, *540*, 527–537. [[CrossRef](#)]
67. Saarela, S.; Holm, S.; Healey, S.P.; Andersen, H.E.; Petersson, H.; Prentius, W.; Patterson, P.L.; Næsset, E.; Gregoire, T.G.; Ståhl, G. Generalized hierarchical model-based estimation for aboveground biomass assessment using GEDI and landsat data. *Remote Sens.* **2018**, *10*, 1832. [[CrossRef](#)]
68. Griffith, D.A. Spatial autocorrelation and eigenfunctions of the geographic weights matrix accompanying geo-referenced data. *Can. Geogr.* **1996**, *40*, 351–367. [[CrossRef](#)]
69. Getis, A.; Griffith, D.A. Comparative spatial filtering in regression analysis. *Geogr. Anal.* **2002**, *34*, 130–140. [[CrossRef](#)]

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.