

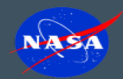
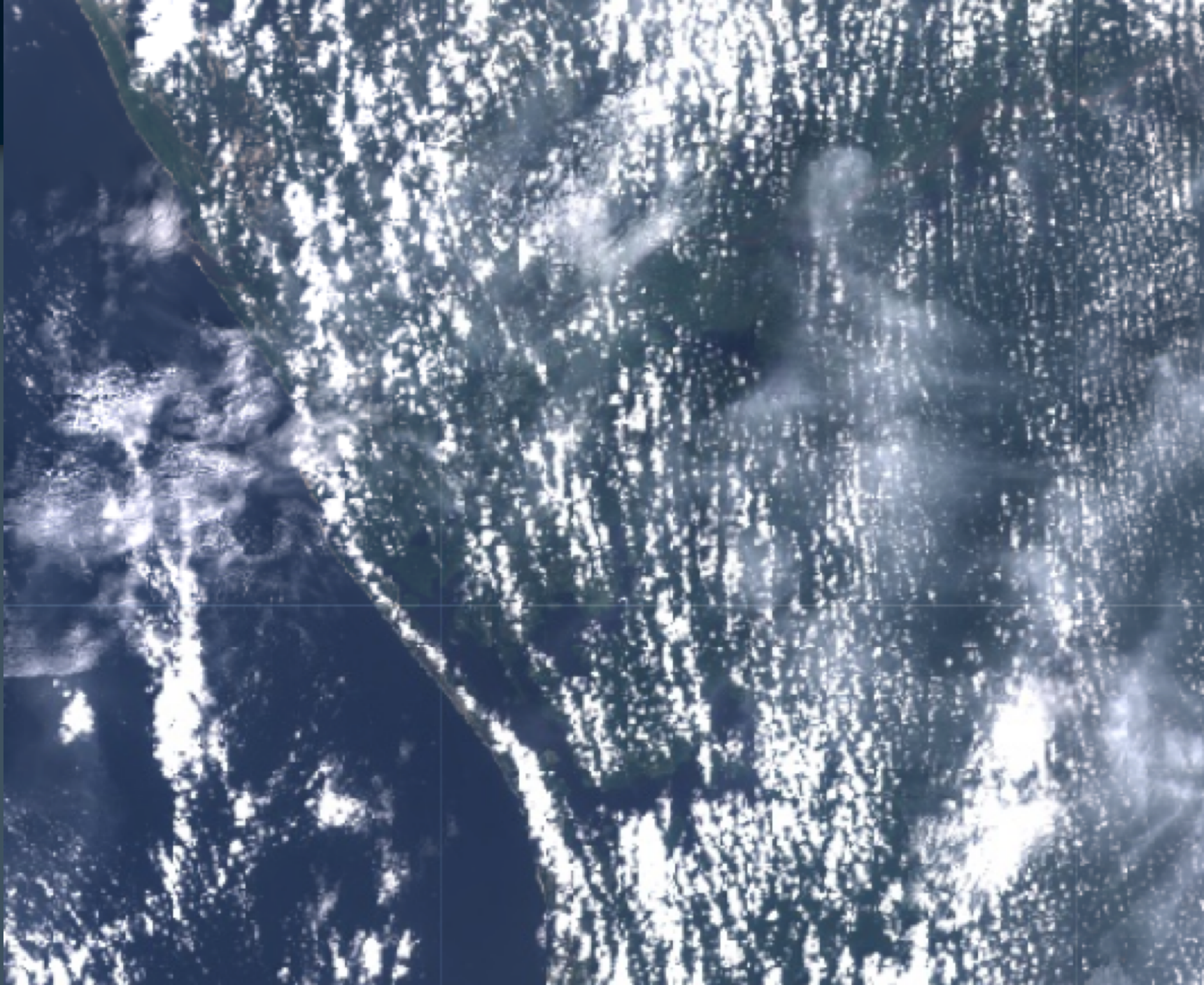
Sampling Designs for SAR-Assisted Forest Biomass

Content created by Hans-Erik Anderson
USDA Forest Service, Pacific Northwest Research Station





- Need: Timely, accurate information on forest carbon stocks and changes to reduce carbon emissions from the forest sector
- Problem: Traditional forest inventory sampling designs are difficult/impossible to implement
 - Large numbers of well distributed field plots = \$\$\$
- Solution! : New sampling designs using RS data!





- Need: Timely, accurate information on forest carbon stocks and changes to reduce carbon emissions from the forest sector
- Problem: Traditional forest inventory sampling designs are difficult/impossible to implement
 - Large numbers of well distributed field plots = \$\$\$
- Solution! : New sampling designs using RS data!
- Updated solution! : New sampling designs using SAR!



- Sampling designs and statistical modelling/estimation frameworks are increasingly sought with the following properties:
 1. Understand, quantify, and communicate uncertainty
 2. Be efficient!
 1. The less field plots the better!
 3. Provide flexibility to accommodate a variety of field plot configurations and remote sensing data acquisition strategies/resolutions

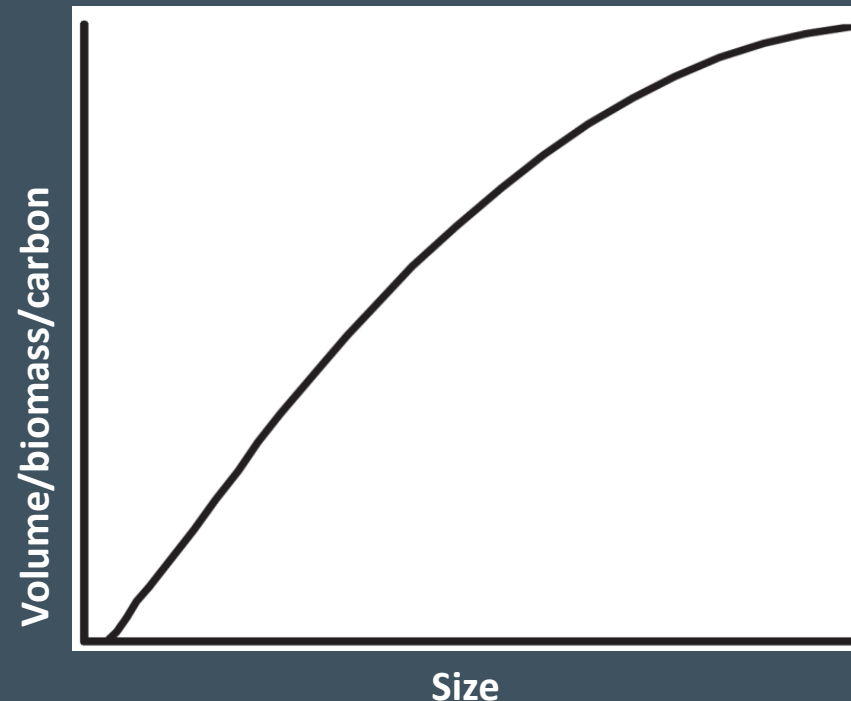


- Three primary sources of variability in the context of a forest carbon inventory and monitoring system:
 - 1) Measurement error
 - 2) Modelling error
 - Allometry
 - Relationships with auxiliary information
 - 3) Sampling error



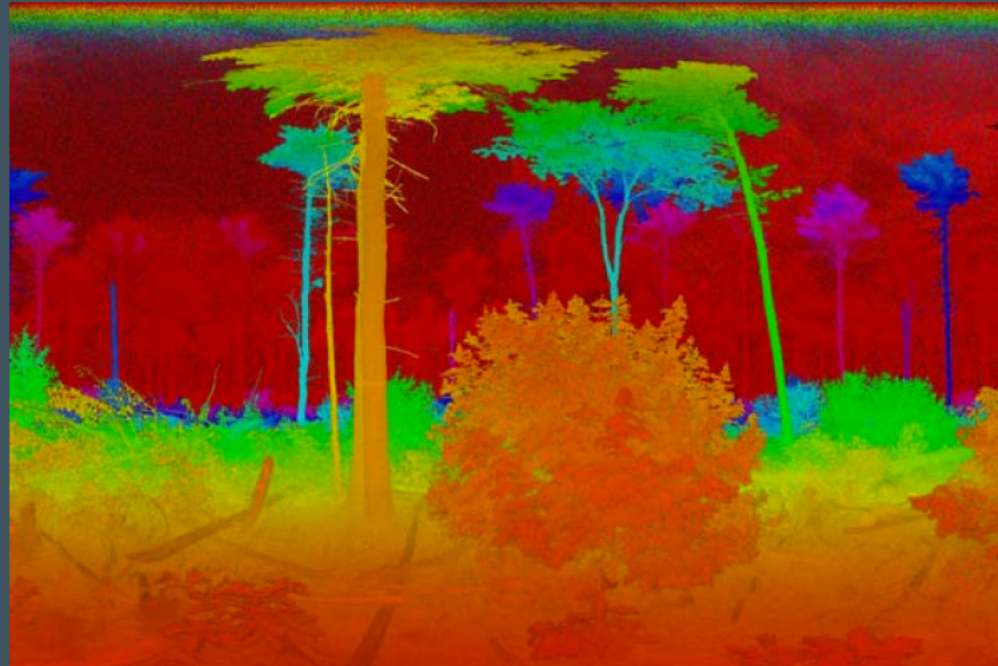
- Discrepancy between a recorded field measurement and the expected value of the measurement as defined by documented protocol
 - Introduced through inadequate training or lack of adherence to protocol.
 - In practice, measurement error is usually assessed and mitigated through quality assurance/quality control (QA/QC) procedures
 - In many cases, assumed to be minimal in comparison to the measurement itself (Gregoire and Valentine, 2008, p. 32).

- In the context of forest carbon monitoring using remote sensing, modelling error is introduced in two ways:
 1. The use of allometric models to estimate tree-level biomass/carbon using various tree measurements (diameter at breast height, height, etc.)
 2. The use of models relating the remotely-sensed measurement (SAR backscatter, air photo-derived canopy height and cover, etc.) to the plot-level biomass/carbon



Modeling error

- Uncertainty due to allometric modelling remains the most difficult source of error to account for in large scale carbon monitoring programs (Duncanson et al., 2017).
 - Emergence of new technologies, such as terrestrial laser scanning (Calders et al., 2015) hold promise for improving the efficiency of field measurements & assessment of uncertainty
- National forest inventory programs often do not explicitly account for this error in official reports.





- Due to its sensitivity to forest biomass, global coverage, and capability to penetrate cloud cover, L-band satellite radar has been used extensively as an auxiliary source of data to support forest monitoring programs across a range of biomes
- L-band dual-polarization (HH, HV) backscatter well-correlated with forest biomass up to approximately 150 Mg/ha, particularly useful for assessment of low-biomass forests (high-latitude boreal biome, tropical semi-arid, savanna forests)
 - **Generalizing relationships b/n radar backscatter and biomass across forest types is inadvisable – radar backscatter from a forest scene is a function of numerous forest structural characteristics (stem density, height, stem diameter) as well as other scene properties (soil moisture, slope, etc.) with varying correlation to tree biomass.**
- Additional forest structure information, perhaps obtained from lidar or repeat-pass interferometry can help to decouple the complex relationships between backscatter and forest structural attributes that can obscure the biomass-backscatter signal at higher biomass levels.



- Typically the only direct measurements are tree parameters (DBH, height, etc)
 - Estimate biomass/carbon via allometry - very limited portion of the landscape
 - Remote sensing provides a more comprehensive picture of forest conditions across a region.
- In this workshop, we explore sampling approaches that utilize a combination of field data and auxiliary information – including wall-to-wall satellite SAR imagery and sampled high-resolution (e.g. lidar) in multi-level inventory designs – to estimate support forest monitoring programs.
 - Obtaining the required precision for carbon estimation within the limitations of the resources



1. Develop a simulated artificial population to facilitate evaluation/comparison of various sampling designs & estimators
2. Calculate elementary statistical estimators (simple random sampling & post-stratification)
3. Calculate model-assisted and model-based estimators with one – two source(s) of auxiliary data (sampled and wall-to-wall)

- The multi-level estimators for forest biomass presented in this workshop provide a range of options for design of carbon monitoring programs, including:
 - ***Model-assisted*** approaches requiring probability samples for all levels of the design that provide design-unbiased estimators
 - ***Model-based*** approaches that may be less expensive to implement due to the lack of requirement for a probability sample, but at the cost of a possibly biased estimator if the model is incorrectly specified.



- Traditional forest inventories: Design Based Inference
 - field plots were randomly distributed
 - each unit in the population of interest has a positive probability of being selected in a sample (“inclusion probability”)
 - the population is considered fixed and all uncertainty in the estimation of a population parameter (total biomass, volume, etc.) is due to variability between randomly-drawn samples from the population.
- In reality - Probabilities of selection can vary across the population (cost reduction or increase the statistical precision of the estimates).
- **Model-assisted** inference is a means of using models to improve precision of estimates within the design-based inferential paradigm.



- Each value from an element in the population is considered a realization of a random variable with a specific probability distribution.
- All population-level values (e.g. total or mean biomass, etc.) are also considered random variables.
- Uncertainty in the estimation of a population parameter is due to randomness in the values observed for each population element.
- The validity of inferences in the model-based paradigm are not dependent upon a random (probability) sample
- MB approach can be applied in situations where collecting a sufficiently large probability sample of field plots is either too expensive or logistically difficult
 - Estimation within small areas
 - Remote regions lacking transportation infrastructure



- Both:
 - offer alternative mechanisms for inference from samples to populations
 - are based on a notion of statistical model
- The **key difference** is whether a *statistical model* considered an **unknown** construct or a **known** one.
 - **Assumptions** about an underlying *statistical model* is what essentially differentiates model-based approach and design-based approach
- In particular, in a model-based approach can be defined as an approach, where statistical model is unknown (hence, the presence of the word "model", as it's the *focus* of discovery). Correspondingly, a design-based approach can be defines as one, where statistical model is known and the *focus* is on a study/experiment design.



- Results from model-based and model-assisted approaches are difficult to compare directly
- Design-based/model-assisted inference:
 - Advantage – Design-based estimators can be considered unbiased (for large sample sizes) regardless of the model that is used
 - In a regional or national forest inventory context, unbiasedness is critical and the design-based approach may be more appropriate
 - Disadvantage – Requirement of adequate probability sample can be very difficult or costly in large remote regions, or small areas
- Model-based inference:
 - Advantage – there is no requirement that the field plots be a probability sample
 - Disadvantage – inferences in the model-based context are conditional on the model and may produce severely-biased estimators in cases where the model is developed using an unrepresentative sample.

Exercise 1: Simulating an Artificial Population



- Simulation can be a useful approach to gain insight into the statistical properties of various survey estimators, especially in the case of somewhat complex, multi-level sampling designs
- In statistics, simulation is used to assess the performance of a method, typically when there is a lack of theoretical background. With simulations, the statistician knows and controls the truth.
- Advantages:
 - Provides the empirical estimation of sampling distributions
 - Studies the misspecification of assumptions in statistical procedures
 - Determines the power in hypothesis tests, etc.
- When generating a simulated population, it is desirable to include realistic correlations between the response variable (e.g. biomass) and the predictor variables (remote sensing metrics)
 - For example, while a multivariate normal distribution can be used to model correlation between several variables, we may prefer that these variables have more realistic marginal distributions (gamma, exponential, etc.).



- A copula function is a useful mathematical tool to simulate a population with specified multivariate correlation structure and marginal distributions
- While an in-depth discussion of copula models is outside the scope of this workshop, they essentially allow for expressing multivariate distributions in terms of their corresponding univariate marginal distributions and a copula function.

Basics of a copula function



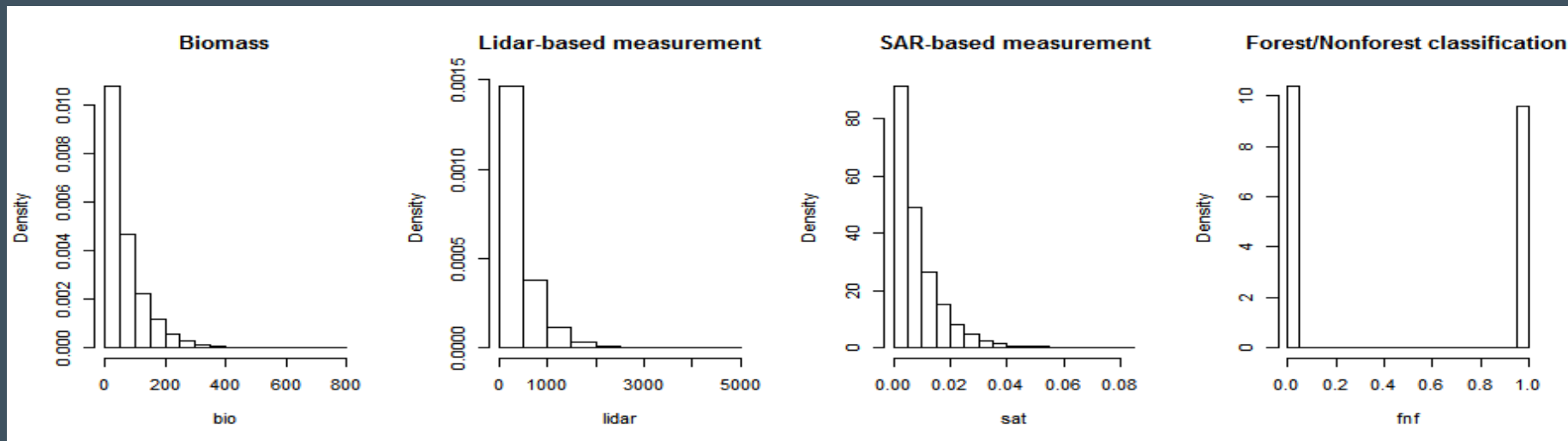
- We want to generate two correlated random variables, x and y
- An easy way would be to sample from a multivariate Gaussian (normal) distribution with the specified correlation, but suppose we want to specify non-Gaussian marginal distributions, $f(x)$ and $g(y)$, for x and y
- Here's how to use a Gaussian copula to achieve this:
 - Generate a, b from a multivariate Gaussian distribution with the desired correlation
 - Transform them into uniformly distributed (but still correlated) variables using the cumulative Gaussian distribution function $u = \Phi(a), v = \Phi(b)$.
 - Transform again to $x = F^{-1}(u), y = G^{-1}(v)$
 - We now have variables with the correct correlation and marginal distributions!



- In this example, we use a copula function to simulate an artificial population where each element has:
 - Forest/nonforest classification
 - Biomass (Mg/ha)
 - Lidar-based measurement (function of lidar-derived height and cover)
 - SAR-based measurement (function of HH and HV backscatter)
- Realistic marginal distributions and correlation structure between remote sensing measurements and field-based biomass were developed based on an analysis of airborne lidar, SAR, and field biomass data from a site in interior Alaska (Andersen et al. 2013).
- To introduce realistic spatial heterogeneity across the simulated area, a binary random field (200×200 grid cells) was used to generate an image with a simulated spatial distribution of “forest” and “nonforest” areas.
- The grid cells within the simulated forest/nonforest image were then populated with elements from the simulated population generated using the copula function.
- In this way, each element in the image had a value for forest/nonforest, biomass, lidar, and SAR, and the simulated population had realistic marginal distributions, correlation structure and spatial pattern of forest cover.

Complete exercise 1

Simulated marginal distributions of biomass, lidar-based measurements, SAR-based measurements, and forest/nonforest classification. Exponential distributions used to model biomass, lidar, and SAR variables; Bernoulli distribution used to model forest/nonforest class.



Design-based estimation: Simple random sampling

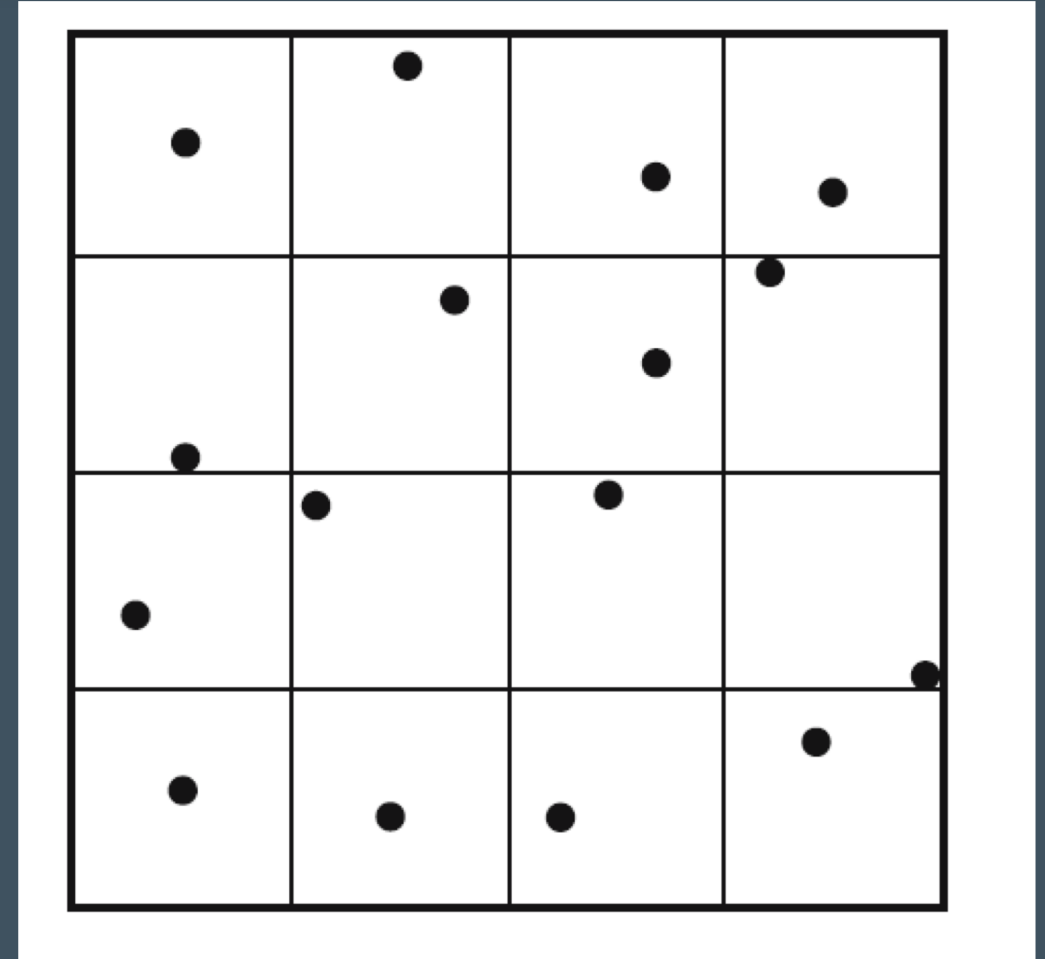


- Simple random sampling (SRS) represents the most fundamental type of design-based sampling and is often used as the basis of comparison for more complex sampling designs
- Given a probability sample of elements of size n from a population of size N , where a forest attribute of interest (Y_i) is obtained for each element i , the SRS estimator of the population mean is given by the sample mean:

$$\hat{\mu}_{SRS} = \bar{Y} = \frac{1}{n} \sum_{j=1}^n Y_i$$

the variance estimator is given by:

$$\hat{V}(\hat{\mu}_{SRS}) = \frac{1}{n(n-1)} \sum_{i=1}^n (Y_i - \bar{Y})^2$$





- The statistical properties of the various estimators can be assessed using the simulated population
- At each iteration:
 1. A random sample of elements is drawn from the population
 2. The point estimator *mean* and the variance estimation are calculated.
 3. Since we know the actual population mean, we can also calculate:
 - The mean percent bias of the point estimator
 - The relative standard error of the point estimator
 - Coverage probability of the 95% confidence interval for the point estimator



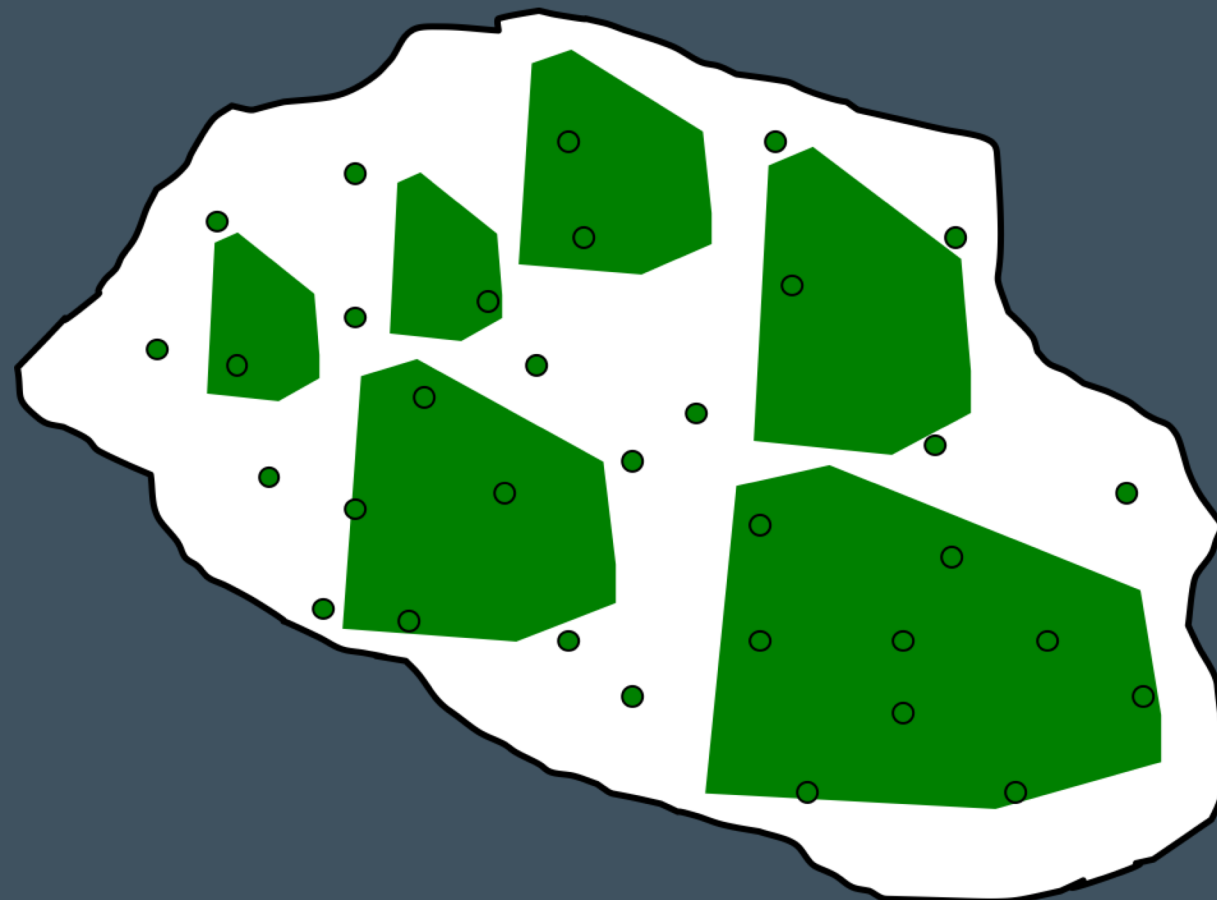
- The coverage probability provides an indication of how reliable (i.e. unbiased) the variance estimator is for a parameter.
- A coverage probability (95% CP) near 95% is an indicator that the 95% confidence intervals (CIs) calculated using this variance estimator are reliable
- Coverage probabilities less than 95% indicate that the calculated 95% CIs are giving a falsely precise estimate of uncertainty, while coverage probabilities greater than 95% indicate that the 95% CIs obtained from this estimator are overly conservative
- The coverage probability of the 95% confidence interval for the point estimator:

$$Prob \left(\hat{\mu} - z_{0.025} \sqrt{\hat{V}(\hat{\mu})} < \mu < \hat{\mu} + z_{(0.975)} \sqrt{\hat{V}(\hat{\mu})} \right) \times 100\%$$

Complete Exercise 2

Exercise 3: Post-stratification

1. Stratify the population into (hopefully) homogenous groups and
2. Estimate the inventory parameter as a weighted average of the strata means

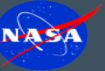


Complete Exercise 3

Regression estimators - Design Based / Model Assisted



- Model-assisted estimators essentially provide a means to use models based on auxiliary data (e.g. remote sensing) to improve inferences within the design-based inferential framework
 - In other words, random (probability) sampling at all levels in the design is the basis for all inference.
- Model-assisted regression estimators are based on a model of the relationship between the forest attribute of interest (e.g. biomass/carbon), Y , and a vector X , of auxiliary variables, formulated as:
$$Y_i = f(X_i; \beta) + \varepsilon_i$$
 - Where $f(X_i; \beta)$ expresses the mean of Y given observation of X , the β s are parameters to be estimated, and ε_i is a random residual term.
- In practice, we don't observe the entire population but estimate the parameters of the regression relationship $\hat{\beta}$ based on a sample of the population.
- We then utilize this regression model and observed vector of auxiliary variables to predict the inventory attribute for a particular unit of the population: $\hat{Y}_i = f(X_i; \hat{\beta})$.





- A regression estimator for the population mean, $\hat{\mu}_{ma,1}$, when a single source of wall-to-wall auxiliary information (e.g. satellite SAR) is given by the following expression:

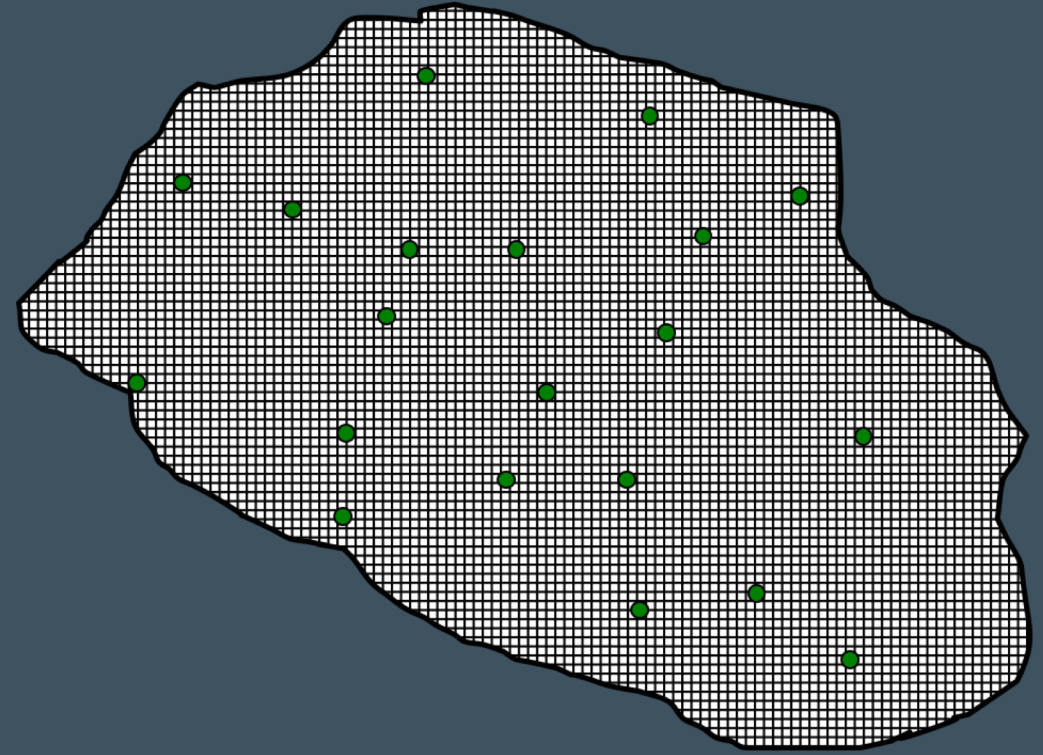
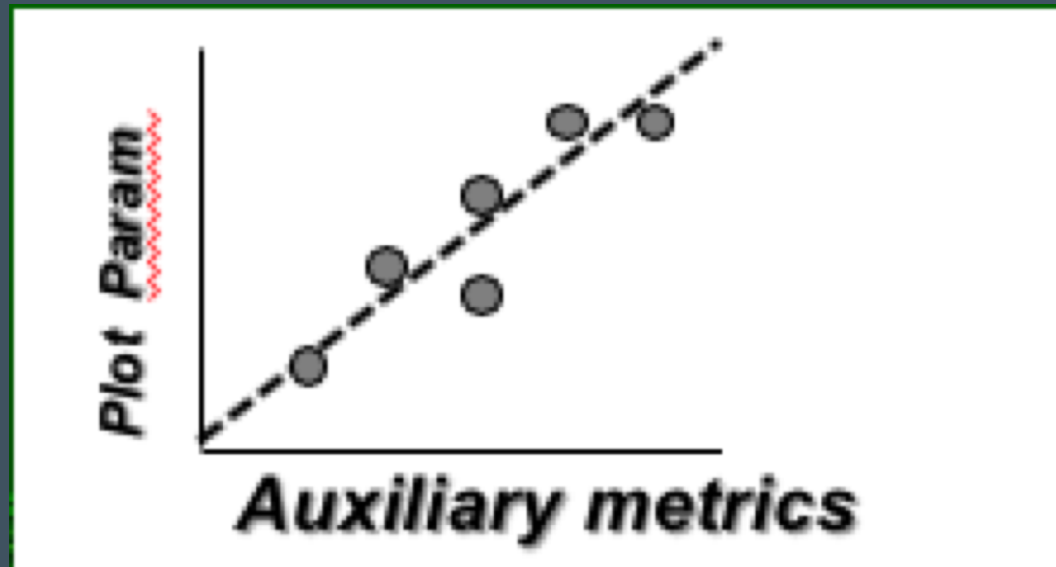
$$\hat{\mu}_{ma,1} = \frac{1}{N} \sum_{j=1}^N \hat{y}_j + \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

- Where N is the population size, n is the sample size, and \hat{Y}_i are predictions from regression model.
 - The first right-hand term in this equation is the sum of the model predictions for the entire population
 - The second right-hand term is a correction term which, when added to the first term, compensates for model bias.
- The regression estimator can be expressed in different forms, but the above formulation is the easiest form to interpret in our context, since the model predictions \hat{Y}_i are based on remotely-sensed imagery or measurements and the second term is the mean of the residuals observed at the field plots.
- We can therefore see that the degree to which the relationship with X explains variability in Y will determine the gain in precision from using the regression estimator as opposed to the SRS estimator.
 - Post-stratification – where population-level strata proportions are used to improve precision of an estimate in the estimation (rather than the design) stage – is a special case of regression estimation where the predictors are categorical variables (for example, satellite image-based land cover classes).

Level 1: Inexpensive measurements using auxiliary data

- Remote sensing plots
- Mapped information

Level 2: Subsample of field plots



Complete Exercise 4



- A few subtle points:
 - In the model-based literature you often find the model-based predictor expressed as the sum of the sampled units plus the estimated sum of the non-sampled units. Since the difference is very small for a large population, in practice we can ignore
 - Likewise, since the difference between the sum of the pixel-level predictions and the sum of the pixel-level predicted means is very small, in practice we can ignore the residual terms in calculation of variance
- It should be noted that when using internal models developed from a SRS sample at all levels of the sampling design, the correction term in the model-assisted estimator goes to zero and the model-based estimator will yield virtually the same point estimate and variance estimator as the model-assisted estimator.
 - However, the assumptions behind these estimators differ and provide more flexibility in the application of the model-based estimator (e.g. application to non-probability samples)
 - Care must be taken to ensure that models are based on a representative (if not random) sample to reduce bias in the point and variance estimators.



- Following McRoberts et al. (2010) and Saarela et al. (2016), if Y is the random variable (AGB) with a mean μ and standard deviation σ the observed AGB value at the i^{th} pixel (y_i) can be represented as:

$$y_i = \mu_i + \epsilon_i$$

- Where $\epsilon_i \sim N(0, \sigma^2)$, the mean AGB at the i^{th} pixel is given by $\mu_i = f(X_i; \beta)$ which is estimated by $\hat{\mu}_i = f(X_i; \hat{\beta})$, and X_i is the lidar-based predictor variable at the i^{th} pixel and $\hat{\beta}$ is the vector of p predicted regression coefficients.

$$\hat{\mu}_{MB} = \frac{1}{N} \sum_{i=1}^N \hat{\mu}_i$$

$$\hat{V}(\hat{\mu}_{MB}) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \widehat{Cov}(\hat{\mu}_i, \hat{\mu}_j)$$



- The model-based estimate of mean AGB over the entire area is (in matrix notation):

$$\widehat{\mu}_U = \iota'_U X_U \widehat{\beta}$$

- Where ι'_U is a N-length column vector where every element equals 1/N, X_U is a $N \times (p+1)$ matrix of satellite auxiliary variables available for each element in the population U.
- The variance of the model-based mean AGB estimate is given by (in matrix notation):

$$V(\widehat{\mu}_U) = \iota'_U X_U V_{\widehat{\beta}} X'_U \iota_U$$

- Where $V_{\widehat{\beta}}$ is the variance-covariance matrix for the regression model parameter estimates $\widehat{\beta}$. For example, in the case of $p=2$, $V_{\widehat{\beta}}$ is given by:

$$\begin{bmatrix} \widehat{V}(\widehat{\beta}_0) & \widehat{Cov}(\widehat{\beta}_0, \widehat{\beta}_1) \\ \widehat{Cov}(\widehat{\beta}_1, \widehat{\beta}_0) & \widehat{V}(\widehat{\beta}_1) \end{bmatrix}$$

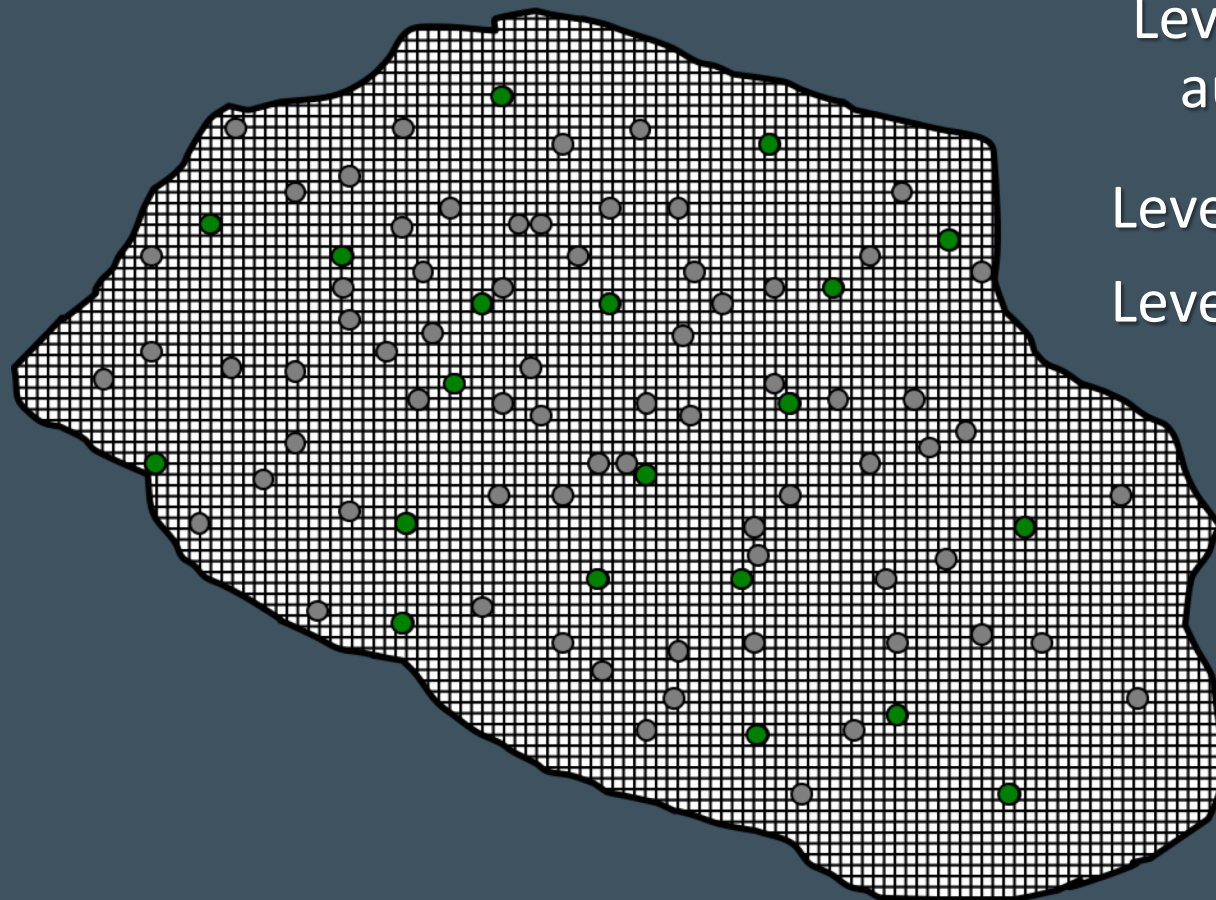
Complete exercise 5

Sampling designs with multiple sources of auxiliary data



- In some cases two types of auxiliary information are available, where one (e.g. satellite SAR imagery) is collected wall-to-wall and another type of (more expensive and higher resolution) remotely-sensed data is collected in a sampling mode.
- For example, multi-level sampling design may consist of:
 1. A large sample of relatively inexpensive photo-interpreted plots distributed over an area of interest, with
 2. Detailed, relatively expensive, field measurements of the attribute of interest (e.g. tree biomass/carbon) collected on a subsample of these photo plots, and
 3. Free, or very inexpensive, satellite image data (SAR) available over the entire area.
- Depending on application and how the data were collected this type of multi-level sampling design can be approached from a *model-based or model-assisted* inferential standpoint.

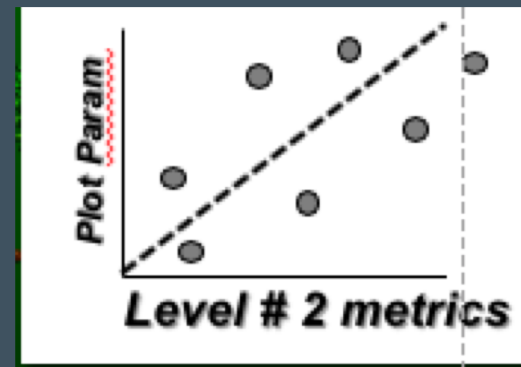
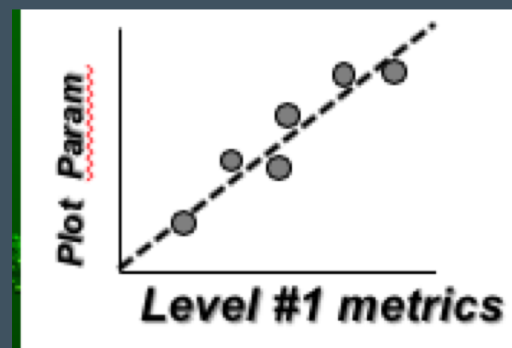
Estimators with two sources of auxiliary information



Level 1: Inexpensive measurements using auxiliary data (X_U)

Level 2: Subsample of inexpensive RS plots (X_1)

Level 3: Subsample of expensive field plots (Y)



Sampling designs with multiple sources of auxiliary data: Model-assisted



- A model-assisted estimator of mean aboveground tree biomass can be developed using field plot data and two sources of auxiliary data in the following manner:
 1. A vector X_U of remote sensing-derived variables that are known for all N elements in the population (U) and
 2. A vector X_1 of remote sensing-derived variables that are known only for the elements in the first phase sample of n_1 units, and the inventory attribute of interest, Y_2 , is only measured on a relatively small second-phase subsample of n_2 photo plots.
- As a specific example, the X_{Ui} variables may represent satellite image data (e.g. SAR HV/HH backscatter) that is available wall-to-wall over the entire study area, and the X_{1i} variables represent photo plot measurements (average tree height, cover, forest type) that are only available at a sample of locations distributed over the area of interest.



- Regression analysis is used to develop a linear model for predicting biomass from photo-based measurements:

$$Y_{1i} = f(X_{1i}; \beta_1) + \varepsilon_{1i}$$

- Satellite-derived predictor variables are used to predict biomass using satellite-based measurements:

$$Y_{Ui} = f(X_{Ui}; \beta_U) + \varepsilon_{Ui}$$

Complete exercise 6

Sampling designs with multiple sources of auxiliary data: Model-based



- A model-based approach to utilizing auxiliary data collected at multiple levels was developed by Saarela et al. (2016)
- As in the previous example of model-based estimator, the relationship between the inventory attribute, Y , which is a random variable (AGB) with a mean μ and standard deviation σ , the observed mean AGB value at the i^{th} pixel (Y_i) can be represented as: $Y_i = f(X_i; \beta) + \epsilon_i$ where $\epsilon_i \sim N(0, \sigma^2)$.
- The mean AGB at the i^{th} pixel is given by $\hat{\mu}_i = f(X_i; \beta)$ which is estimated by $\hat{\mu}_i = f(X_i; \hat{\beta})$, where X_i is the set of lidar-based predictor variables available for the second phase sample n_2 of the population and $\hat{\beta}$ is the vector of p predicted regression coefficients.
- This linear model is used to estimate the mean AGB at every pixel in the first phase sample n_1 : $\hat{\mu}_{1i} = f(X_{1i}; \hat{\beta})$
- In this hierarchical modelling framework, a second model is developed relating the satellite-based predictor variables Z_U available over the entire population to the $\hat{\mu}_1$ predictions available within the first phase sample:

$$\hat{\mu}_{1i} = f(Z_{1i}; \alpha_1) + \omega_{1i}$$

where $\omega_{1i} \sim N(0, \sigma^2)$ and α_1 is the vector of model coefficients linking lidar-estimated AGB values and the satellite predictor variables estimated by $\hat{\alpha}_1$.

Sampling designs with multiple sources of auxiliary data: Model-based



- The model-based estimate of mean AGB over the entire area is

$$\widehat{\mu}_U = \iota'_U \mathbf{Z}_U \widehat{\alpha}_1$$

- The variance of the model-based mean AGB estimate is given by:

$$\widehat{V}(\widehat{\mu}_U) = \iota'_U \mathbf{Z}_U \mathbf{V}_{\widehat{\alpha}_1} \mathbf{Z}'_U \iota_U$$

- Where $\mathbf{V}_{\widehat{\alpha}_1}$ is the variance-covariance matrix for the regression model parameter estimates by $\widehat{\alpha}_1$ given by:

$$\mathbf{V}_{\widehat{\alpha}} = \frac{\widehat{\omega}'_1 \widehat{\omega}_1}{M-q-1} (\mathbf{Z}'_1 \mathbf{Z}_1)^{-1} + (\mathbf{Z}'_1 \mathbf{Z}_1)^{-1} \mathbf{Z}'_1 [\mathbf{X}_1 \mathbf{V}_{\widehat{\beta}} \mathbf{X}'_1] \mathbf{Z}_1 (\mathbf{Z}'_1 \mathbf{Z}_1)^{-1}$$

- Where $\mathbf{V}_{\widehat{\beta}}$ is the variance-covariance matrix for the regression model parameter estimates $\widehat{\beta}$ and $\widehat{\omega}_1 = \mathbf{X}_1 \widehat{\beta} - \mathbf{Z}_1 \widehat{\alpha}_1$ is a n_1 length vector of model residuals.

Do exercise 7

Toward an optimal sample for model-based estimation: Lidar-informed stratified sampling for carbon monitoring



- From a model-based standpoint, it can be desirable to collect a sample that is well-distributed across the range of predictor variables
 - Especially important in fitting nonlinear models
- One approach is to stratify the population using lidar structural classes (e.g. height, cover), then distribute field plots as a stratified random sample
 - Ensures an adequate sample is collected across range of structures
 - Introduces additional complexity into estimators since inclusion probabilities will vary



- The IPCC has specified good practice as it pertains to the concept of REDD+ forest monitoring as inventory design that “neither over- nor under-estimates so far as can be judged, and in which uncertainties are reduced as far as is practicable.” (GFOI, 2016).
- This guidance essentially promotes the implementation of monitoring programs that maximize precision of estimates, while minimizing bias, within the constraints of available resources.